



PROBING THE TOP-YUKAWA COUPLING BY
SEARCHING FOR ASSOCIATED HIGGS BOSON
PRODUCTION WITH A SINGLE TOP QUARK AT
THE CMS EXPERIMENT

Zur Erlangung des akademischen Grades eines
DOKTORS DER NATURWISSENSCHAFTEN
von der Fakultät für Physik des
Karlsruher Instituts für Technologie (KIT)

genehmigte

DISSERTATION

von

Dipl.-Phys. Simon Fink
aus Ludwigsburg

Mündliche Prüfung: 08. Juli 2016

Referent: Prof. Dr. Th. Müller
Institut für Experimentelle Kernphysik

Korreferent: Prof. Dr. G. Quast
Institut für Experimentelle Kernphysik

Introduction

It is the curiosity of man and the urge to know how things work and what they are made of that drives humanity and science in particular. The discovery of the Higgs boson in 2012 [1, 2] is the latest milestone in elementary particle physics, the study of the smallest building blocks of nature, making the underlying theory framework, the standard model of particle physics, self-consistent. Simultaneously, phenomenons like the existence of dark matter or the matter-antimatter asymmetry in the universe make it evident, that physics beyond the standard model must exist. Unlike the situation before the Higgs boson discovery, where unphysical predictions of the theory made it clear that either the Higgs boson would be discovered or another mechanism would have to emerge, such a landmark, and thus a guaranteed discovery, is missing in today's particle physics. If nature is especially cruel, this could mean that new physics phenomena would not appear at human reachable energies. In this situation the main responsibility lies with the experimental physicists to lead the way by searching for deviations from the prediction. An obvious first step would include measuring the youngest member and linchpin of the standard model, the Higgs boson, in excruciating precision, as possible deviations from the prediction could hint towards what kind of new physics is realized in our world and where to look for it. An intriguing possibility, where such deviations could surface, is the coupling of the Higgs boson to top quarks. This coupling of the two heaviest elementary particles can be directly accessed by studying the production of a single top quark which is produced in association with a Higgs boson. Exactly this process is at the core of this thesis, as a search for it is conducted at the LHC at $\sqrt{s} = 8 \text{ TeV}$ and 13 TeV .

This first chapter of this thesis gives a short introduction into the standard model of particle physics and the interactions and particles described by it. The special focus lies on the Higgs boson and its underlying mechanism. Additionally, a theoretical motivation into the associated production of Higgs bosons with single top quarks is provided.

In the second chapter an overview of the experimental setup is provided, namely the LHC accelerator complex and the CMS detector.

The third chapter covers the generation of Monte Carlo simulation samples, necessary tools to interpret the measured data. This chapter also describes the journey from measured electrical signal in the detector to a fully reconstructed physical object.

In the fourth chapter an overview of the methodology used during the course of the analysis is given. Multivariate analysis tools are employed to separate signal events from background events. Additionally, the principles used in the statistical inference of the analysis are described.

In the fifth and sixth chapter the searches for the associated production of single top quarks and Higgs bosons are conducted at $\sqrt{s} = 8 \text{ TeV}$ and 13 TeV , respectively.

Finally, in the last chapter a conclusion, as well as an outlook for the associated single top and Higgs boson production for the coming years is provided.

Contents

1	Theory Essentials and Current Status	1
1.1	The Standard Model of Particle Physics	1
1.1.1	Gauge Bosons	2
1.1.2	Fermions	3
1.2	The Higgs Boson	5
1.2.1	The Higgs Mechanism	5
1.2.2	Experimental Qualities of the Higgs Boson	7
1.2.3	Production Mechanisms	7
1.2.4	Decay Channels	8
1.3	Associated Production of Higgs Boson and Single Top Quark	8
1.3.1	CP -mixing in tH	11
1.3.2	Experimental Status in tH Production	11
2	The CMS Experiment at the LHC	15
2.1	The Large Hadron Collider	15
2.2	The Compact Muon Solenoid	19
2.2.1	Tracking System	19
2.2.2	Electromagnetic Calorimeter	22
2.2.3	Hadron Calorimeter	22
2.2.4	Superconducting Magnet	24
2.2.5	Muon System	24
2.2.6	Trigger System	25
2.2.7	Computing Model	26
3	Event Generation, Simulation & Reconstruction	29
3.1	Event Generation	29
3.1.1	Hard Scattering	30
3.1.2	Parton Shower	32
3.1.3	Hadronization and Decay of Unstable Particles	33
3.1.4	Underlying Event and Pileup	33
3.1.5	Monte Carlo Generators	33
3.2	Detector Simulation	35
3.3	Event Reconstruction	35
3.3.1	Particle-Flow Algorithm	35

3.3.2	Vertex and Track Reconstruction	36
3.3.3	Muon Reconstruction	37
3.3.4	Electron Reconstruction	37
3.3.5	Photon and Hadron Reconstruction	37
3.3.6	Jet Reconstruction	38
3.3.7	Missing Transverse Energy	42
4	Multivariate & Statistical Methods	45
4.1	Multivariate Methodology	45
4.1.1	Artificial Neural Networks	45
4.1.2	Boosted Decision Trees	46
4.1.3	Overtraining	48
4.1.4	Variable Ranking	48
4.2	Statistical Inference	49
4.2.1	Maximum Likelihood Parameter Estimation	49
4.2.2	Nuisance Parameter Treatment	50
4.2.3	Exclusion Limits	51
4.2.4	Asymptotic Limits	52
5	Search for Associated Production of Single Top Quark and Higgs Boson at $\sqrt{s} = 8$ TeV	53
5.1	Analysis Strategy	53
5.2	The tHq Process	55
5.3	Background Processes	55
5.4	Datasets	59
5.4.1	Heavy Flavor Splitting	60
5.5	Physics Objects	60
5.5.1	Primary Vertices	60
5.5.2	Muons	60
5.5.3	Electrons	61
5.5.4	Jets	61
5.5.5	Missing Transverse Energy	62
5.5.6	W Boson Reconstruction	62
5.6	Monte Carlo Corrections	63
5.6.1	Pileup Reweighting	63
5.6.2	Lepton Efficiency Scale Factors	63
5.6.3	Top Quark p_T Reweighting	64
5.6.4	b-tagging Efficiency Correction	64
5.6.5	Jet Pseudorapidity in the Forward Region	65
5.7	Event Selection	66
5.8	Event Reconstruction	69
5.8.1	tHq Reconstruction	69
5.8.2	$t\bar{t}$ Reconstruction	74
5.8.3	Evaluation of Reconstruction Methods	79

5.9	Event Classification	81
5.10	Systematic Uncertainties	90
5.11	Results	93
5.11.1	Fit of Final Discriminator	93
5.11.2	Analysis of Nuisance Parameters	96
5.11.3	CL _S Limit	97
5.12	Combination with Other Decay Channels	97
5.13	Summary	101
6	Search for Associated Production of Single Top Quarks and Higgs Bosons at $\sqrt{s} = 13$ TeV	105
6.1	Analysis Developments	106
6.2	Analysis Strategy	107
6.3	Signal Processes	108
6.3.1	t -channel	108
6.3.2	tW -channel	108
6.4	Background Processes	109
6.5	Datasets	109
6.6	Object Definitions	111
6.6.1	Primary Vertices	112
6.6.2	Muons	112
6.6.3	Electrons	112
6.6.4	Jets	112
6.6.5	Missing Transverse Energy	112
6.6.6	W Boson Reconstruction	113
6.7	Monte Carlo Corrections	113
6.7.1	Pileup Reweighting	113
6.7.2	Lepton Efficiency Scale Factors	113
6.7.3	CSV Shape Reweighting	114
6.7.4	Jet Pseudorapidity in the Forward Region	114
6.8	Event Selection	115
6.9	Event Reconstruction	116
6.9.1	tHq Reconstruction	116
6.9.2	$t\bar{t}$ Reconstruction	120
6.9.3	Evaluation of Reconstruction Methods	121
6.10	Event Classification	127
6.11	Systematic Uncertainties	136
6.12	Results	139
6.12.1	Fit of Final Discriminator	139
6.12.2	Analysis of Nuisance Parameters	139
6.12.3	CL _S Limits	142
6.13	Search for $C\mathcal{P}$ -mixing in tHq	146
6.14	Summary	150

7 Conclusion and Outlook	153
A Appendix - Search for tHq Production at $\sqrt{s} = 8$ TeV	157
B Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV	167
List of Figures	185
List of Tables	189
Bibliography	191

1. Theory Essentials and Current Status

Particle physics has been an evolving field for decades with experimental and theoretical developments in constant interplay made possible by a common framework. The standard model of particle physics, the model describing almost all interactions of known elementary particles, reached its culmination with the discovery of the Higgs boson in 2012 [1, 2]. The new-found Higgs boson is the confirmation for the existence of a Higgs field, a field postulated 50 years earlier coincidentally by François Englert and Robert Brout [3]; Peter Higgs [4]; and Thomas Kibble, Gerald Guralnik and Carl Hagen [5, 6]. Although the standard model is an extremely successful effective model, it is clear that it can not be the terminus of particle physics, due to its inability to describe observed phenomena, such as the existence of dark matter, or due to conceptual problems, like the inclusion of the gravitational force into the model.

Knowing that physics beyond the standard model (BSM) must exist, it is mandatory to study the properties of the Higgs boson in excruciating detail. One of the presumable points where deviations from the expectation could turn up is the coupling of the Higgs boson to the heaviest elementary particle, the top quark. Exactly this coupling parameter is the subject of study in this thesis and an introduction to the theoretical model describing this parameter is given in this chapter.

The first part of this chapter gives a short overview of the standard model of particle physics and introduces its particle content. In the second part the Higgs boson and its theoretical basis are described. Additionally, the current status of the Higgs boson analyses currently performed at the LHC are presented. The last chapter gives a theoretical motivation of the search for the associated production of single top quarks with Higgs bosons, the topic of this thesis. The information provided in this chapter is obtained from various text books [7, 8], if not otherwise stated.

1.1. The Standard Model of Particle Physics

The standard model is able to precisely predict fields and interactions with the introduction of only a few elementary particles. The theoretical model underlying the standard model is the quantum field theory (QFT), a theory so complex that a thorough description would be by far out of scope of this thesis, but which is subject of many different textbooks. In this chapter a more phenomenological approach is chosen, whereas essential details are provided. For the sake of simplicity natural units $\hbar = c = 1$ are used throughout this thesis.

The standard model is the consolidation of electroweak theory [9–11], which is described by a $SU(2)_L \otimes U(1)_Y$ symmetry group, and quantum chromodynamics (QCD), described by a $SU(3)_C$ symmetry group. The QFT describes nature with the use of particle and force fields. The quantum excitations of these fields can be interpreted as particles and are henceforth

addressed as such.

According to Noether's theorem [12] each symmetry of the underlying physics results in a conserved charge. The particles incorporated in the standard model each contain a set of these conserved quantum numbers, such as the color charge, the weak isospin, the electric charge, and the spin. The most general ordering scheme is based on the spin: particles with a half-integer spin are known as fermions, the excitations of the fermionic fields, and particles with an integer spin, known as bosons. The SM bosons themselves consist of the gauge bosons, the excitations of the fields of the fundamental interactions, and the scalar Higgs boson, the boson associated with the Higgs field. Fermions and gauge bosons are explained in the following subsections, whereas the Higgs boson will be explained separately in more detail as it plays a central role in this thesis.

1.1.1. Gauge Bosons

The mediator particles of the forces incorporated in the standard model are the gauge bosons. With a spin of 1 the gauge bosons obey the Bose-Einstein statistics, allowing particles to occupy the same state with same quantum numbers. An overview of the gauge bosons in the SM can be found in Table 1.1.

The photons are the mediators of the electromagnetic (EM) force, which is described by the theory of quantum electrodynamics (QED). Photons couple to the electromagnetic charge of particles, but carry no such charges themselves, and due to their masslessness, the range of the electromagnetic interaction is infinite.

The gluons are the mediator particles of the strong force and are also massless, just like photons. The range of the strong force, however, is very much limited in contrary to the EM force, as gluons carry one unit of color and anticolor charge each, what causes a self-coupling of gluons. This leads to an effect known as *confinement*, where the field lines of two divergent quarks form a narrow tube, as the gluons are attracted to each other, resulting in a linearly increasing potential for particles flying apart.

The gauge bosons of the weak force are the W^\pm and Z bosons. While W bosons are electrically charged and couple to the third component of the weak isospin I_3 of particles, Z bosons are electrically neutral and couple to the weak hypercharge $Y = 2q - I_3$, a linear combination of electric charge and weak isospin. The heavy masses of $m_Z = 91.2 \text{ GeV}$ and $m_W = 80.4 \text{ GeV}$ for Z bosons and W bosons, respectively, limit the range of the weak interaction severely to scales below 10^{-16} m .

A huge success was the unification of electromagnetic and weak interaction by the Glashow-Weinberg-Salam theory of electroweak interactions. Electromagnetic and weak interactions are regarded as two aspects of one unified interaction and W bosons, Z boson and photons are actually mixed states of the bosons of the electroweak interaction. However, the mediators of the electroweak force are predicted to be massless, which is in clear contrast to the experimentally determined high masses of W and Z boson. This conundrum is solved via the electroweak symmetry breaking (EWSB) introduced by the Higgs mechanism, which is explained in more detail in Section 1.2.1.

Table 1.1.: The gauge bosons of the standard model with their corresponding interactions and their masses are listed. Mass values are taken from Reference [13]. Photon and gluon are predicted to be massless and no evidence for a non-zero mass has been found so far.

Gauge Boson	Interaction	Mass (GeV)
γ (photon)	Electromagnetic	–
g (gluon)	Strong	–
W^\pm	Weak	80.385 ± 0.015
Z		91.188 ± 0.002

1.1.2. Fermions

Fermions described in the standard model can be divided into two classes, the leptons and quarks, with three generations of particles each. An overview of the fermionic particle content of the standard model can be found in Table 1.2.

Each generation consists of two closely related particles, for the fermions one electromagnetically charged particle and its corresponding neutral neutrino: the electron e and the electron neutrino ν_e in the first, the muon μ and the muon neutrino ν_μ in the second, and the tau τ and the tau neutrino ν_τ in the third generation. For each particle there exists also an antiparticle with inverted quantum numbers.

The generations are ordered by the masses of the charged leptons, with the electron being the lightest and the τ being the heaviest lepton. Neutrinos were long thought to be massless, but the observation of neutrino flavor oscillations [14–16] proved irrefutably that at least two neutrinos have a non-zero mass, whereas the actual values remain to be measured and an extension to the SM has to be found which is able to generate these masses.

The electron is the only stable charged lepton, muons and taus decay after their short lifetimes of $\tau_\mu = 2.2 \cdot 10^{-6}$ s and $\tau_\tau = 2.9 \cdot 10^{-13}$ s [13], respectively. Charged leptons are affected by the electromagnetic and the weak force, as well as by gravity, which is many orders of magnitude weaker than the other forces and is hence neglected. As neutrinos solely interact via the weak force they are very difficult to measure.

The second class are the quarks, particles that carry color charge and thus interact, additionally to the electromagnetic and weak interactions, also via the strong force. The three generation of quarks are composed of the up (u) and down quark (d) in the first, the charm (c) and strange quark (s) in the second, and the top (t) and bottom quark (b) in the third generation. As for the leptons, there is an antiquark for each quark with opposite quantum numbers. Quarks of the up-type have a positive charge of $+\frac{2}{3}e$ and down-type quarks a negative charge of $-\frac{1}{3}e$. Contrary to the gluon, quarks only possess exactly one unit of color charge, or anti-color for the case of antiquarks.

The confinement already mentioned above prohibits the existence of particles with a net color

1. Theory Essentials and Current Status

Table 1.2.: Overview of the fermionic particle content of the standard model with their respective quantum numbers.

Fermion	Generation			Electric Charge (e)	Color	Weak Isopin
	I	II	III			
Leptons	ν_e	ν_μ	ν_τ	-	-	$+1/2$
	e	μ	τ	-1		$-1/2$
Quarks	u	c	t	$+2/3$	r g b	$+1/2$
	d	s	b	$-1/3$		$-1/2$

charge, hence quarks only appear in hadrons, bound states of quarks, predominantly consisting of either three quarks (baryons), or $q\bar{q}$ (mesons). If the distance between two quarks and thus the energy becomes too large, a new pair of quark and antiquark is produced from the vacuum. However, the *asymptotic freedom* of the strong interaction causes quarks to act quasi-freely at short distances. Although initially prohibited, interactions via the weak force show a non-zero transition amplitude between quarks of different generations. The weak force couples not to the mass eigenstates of the quarks but to a set of flavor eigenstates which is obtained by rotating the mass eigenstates with the Cabibbo-Kobayashi-Maskawa (CKM) matrix [17, 18]. The diagonal elements of the CKM matrix are the dominant elements keeping the admixture of different flavors small.

The Top Quark

The top quark, discovered in 1995 [19, 20] at the Tevatron in Chicago, is - with a mass of $m_t = 173.21$ GeV [13] - the heaviest elementary particle in the standard model. Due to its large mass and short lifetime top quarks do not take part in the hadronization and do not form bound states but decay beforehand, thereby conserving the quantum number information of the initial quark.

The mixture in the quark sector is minimal for the third generation of top and bottom quark, leading to an almost exclusive top quark decay into a W boson and a bottom quark.

Top quarks can either be produced in pairs via the strong interaction, in pp colliders predominantly via gluon-fusion process, or singly via the weak interaction. The top quark pair production constitutes the dominant background for the analyses of this thesis. The single top quark production can be categorized into three production channels: The t -channel, the s -channel and the tW -channel. Single top quark production in the t - [21] and in the tW -channel [22] have been successfully observed at the LHC, whereas the s -channel is only measured at the LHC with a signal significance of 2.5 standard deviations [23].

1.2. The Higgs Boson

The only scalar boson in the standard model is the Higgs boson. In the following a short description of the Brout-Engler-Higgs mechanism is provided, as well as an overview of the experimental properties of the Higgs boson necessary for this thesis. For convenience the Brout-Engler-Higgs mechanism will only be addressed as Higgs mechanism in the remainder of this thesis.

1.2.1. The Higgs Mechanism

The introduction of a Higgs field and its excitations, the Higgs bosons, allows the spontaneous breaking of the electroweak symmetry, thereby lifting the restriction to accredit the gauge bosons of the weak interaction with non-zero masses, which would otherwise lead to violation of gauge invariance. This chapter gives a brief overview of the mechanism, a far more detailed description can be found in the literature [8].

In the mechanism a complex scalar field

$$\phi = \frac{1}{\sqrt{2}} \begin{pmatrix} \phi^+ \\ \phi^0 \end{pmatrix}$$

is introduced and its potential is chosen to have the form

$$V(\phi^\dagger \phi) = -\mu^2 \phi^\dagger \phi + \frac{\lambda}{2} (\phi^\dagger \phi)^2$$

with the parameters μ and λ defining the form of the potential. The ground state of the field is dependent on the shape of the potential: Whereas the vacuum expectation value would be $\langle 0 | \phi | 0 \rangle = 0$ for a positive μ^2 , a local maximum is found at $|\phi| = 0$ for a negative μ^2 . A simplified visualization of such a potential with only two degrees of freedom can be seen in Figure 1.1. The potential itself is symmetric under $\mathcal{SU}(2)_L \otimes \mathcal{U}(1)_Y$ transformations and has an infinite number of minima located at

$$\phi^\dagger \phi = \frac{-\mu^2}{2\lambda} = \frac{v^2}{2}.$$

The distinct choice of one ground state, which can be chosen without loss of generality as

$$\phi_0 = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ v \end{pmatrix},$$

leads to the breaking of this symmetry and a breaking of all of the $\mathcal{SU}(2)_L \otimes \mathcal{U}(1)_Y$ generators, with an exception of the generator of the $\mathcal{U}(1)_Q$ symmetry. With the utilization of the unitary gauge above (see Reference [8]), small perturbations about this minimum can be parametrized with a hermitian field $H(x)$ as

$$\phi(x) = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ v + H(x) \end{pmatrix}. \quad (1.1)$$

Inserting this into a standard Lagrangian density, the underlying equation describing the different field interactions as defined in Reference [8], yields a mass term for the Higgs boson of

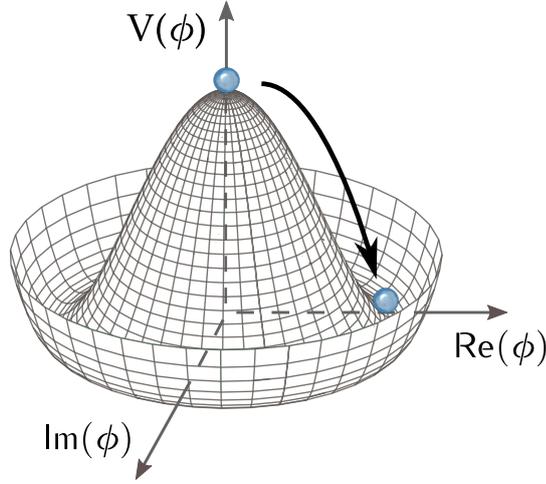


Figure 1.1.: Simplified visualization of the potential of a complex scalar field ϕ with two degrees of freedom and parameters as described in Section 1.2.1. The potential has a local maximum for $\phi = 0$, but an infinite number of minima located at $|\phi| = v$. By choosing one distinct minimum, e. g. $\text{Re}(\phi) = v$, $\text{Im}(\phi) = 0$, the symmetry of the initial potential is broken.

$m_H = \sqrt{\lambda/2}v^2$ and mass terms for W^+ , W^- and Z fields. This is in accordance with the Goldstone theorem [24], although the theorem predicts one massless Goldstone boson per spontaneously broken generator. In the breaking of this gauge symmetry the Goldstone bosons are *eaten* by the gauge bosons, attributing them with masses and a longitudinal polarization. The photon remains massless due to the unbroken generator of the $\mathcal{U}(1)_Q$ symmetry.

The Higgs mechanism successfully lends mass to the gauge bosons of the weak force, but so far the inclusion of fermion mass terms in the form of $m_f\bar{\psi}\psi$ into the SM Lagrangian is still permitted as they are not invariant under $\mathcal{SU}(2)_L \otimes \mathcal{U}(1)_Y$ transformations. However, the inclusion of the Higgs field as written in Equation 1.1 introduces terms of the form $(vy_f/\sqrt{2})\bar{\psi}\psi$ in the Lagrangian that are independent of x and where y_f is the Yukawa-coupling of a fermion f to the Higgs field. Given these parameters, a mass term for the fermion can be identified as

$$m_f = \frac{y_f v}{\sqrt{2}}.$$

The application of the principle of Yukawa interactions to the quark sector of the SM is aggravated by the mismatch of flavor eigenstates and mass eigenstates of quarks. This leads to the theoretical deduction of the CKM matrix already explained earlier in Section 1.1.2.

The incorporation of neutrino masses into the SM is still a hot topic among physicists, as models able to generate these mass terms exist, like the seesaw mechanism or the Majorana nature of neutrinos, but lack experimental proof.

With a measured top quark mass of $m_t = 173.2 \text{ GeV}$ the top-Yukawa coupling is strikingly found to be $y_t \sim 1$, whereas other Yukawa couplings are of the order of 10^{-2} or less. The top-Yukawa coupling is of special interest when searching for BSM physics, as it could hint to the scale of new physics [25] and is even a crucial parameter in the stability of the vacuum [26].

1.2.2. Experimental Qualities of the Higgs Boson

The focus of Higgs boson analyses since its discovery shifted towards the measurement of the properties of the long elusive boson. As a free parameter in the standard model the Higgs boson mass is crucial for subsequent predictions about the boson itself. The most precise result is obtained in a combination of the $H \rightarrow ZZ$ and $H \rightarrow \gamma\gamma$ measurements by the CMS and ATLAS collaborations [27], and the mass is estimated to be

$$m_H = 125.09 \pm 0.21(\text{stat.}) \pm 0.11(\text{syst.}) \text{ GeV.}$$

With this knowledge the corresponding width of the Higgs boson is predicted to be $\Gamma_H \sim 4 \text{ MeV}$, whereas the most accurate measurement of width is able to constrain it to $\Gamma_H < 13 \text{ MeV}$ [28]. The spin of the Higgs boson is also subject of many analyses within CMS and ATLAS and the $J^{CP} = 1^\pm$ hypotheses, which would be forbidden for a standard model Higgs boson by the Landau-Yang theorem [29, 30], and several of the considered $J^{CP} = 2$ hypotheses are already excluded at 95% C.L. [31]. The determination of the CP properties of the Higgs boson proves much more difficult, as the observed particle could in principle consist of any mixture of CP -even and CP -odd components.

1.2.3. Production Mechanisms

The Higgs boson can be produced in several different ways which are explained in the following.

Gluon-Gluon Fusion The gluon-gluon fusion is the main production channel at the LHC due to the high momentum fraction of the proton carried by gluons in the pp collider. The massless gluons can not couple directly to the Higgs boson, hence this production is mediated via a loop of virtual fermions. The magnitude of the contribution scales with m_f^2 , causing the top quark to constitute the main contribution.

In gluon-gluon fusion the Higgs boson is produced at leading order without any secondary particles, which makes this a well suited channel for searches for processes with distinct final states, such as $H \rightarrow \gamma\gamma$ or $H \rightarrow ZZ \rightarrow 4\ell$.

Vector Boson Fusion The production mechanism with the second largest production cross section is the vector boson fusion (VBF), where the Higgs boson is produced in association with two light quarks, which are predominantly found in the forward region of the detector.

Higgsstrahlung In the Higgsstrahlung process the Higgs boson is emitted by one of the weak gauge bosons. The decay products of this vector boson can serve as a trigger object to suppress multijet background processes for Higgs boson decay channels like $H \rightarrow b\bar{b}$. The possible presence of a top quark loop in the production in next-to-leading order also leads to a small sensitivity to the relative sign of the top-Yukawa coupling.

Associated Top Quark Pair Production The production of a Higgs boson in association with a top quark pair leads to a variety of final states, depending on the Higgs boson decays and the individual top quark decays. The $t\bar{t}H$ production is an important background process to the analyses of this thesis, and more information can be found in Chapter 5.3.

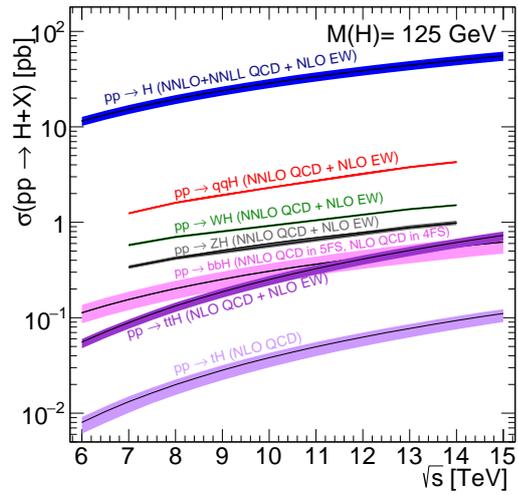


Figure 1.2.: The dependence of the Higgs boson production cross sections on the center-of-mass energy is shown. In blue the total Higgs boson production cross section is depicted. It is also visible that $t\bar{t}H$ and tH production are the processes with the highest cross section increase among all production modes. The diagram is taken from Reference [32].

1.2.4. Decay Channels

After a precise measurement of the Higgs boson mass the decay branching ratios (\mathcal{BR}) can be predicted accurately. The largest branching ratio of 58.4%¹ is predicted for the decay into a pair of bottom quarks, as the heaviest particles whose on-shell production is kinematically allowed in the decay. While allowing for the investigation of the most produced Higgs bosons, the strong interactions of the bottom quarks lead to a tough to analyze multijet environment, which can only be overcome with a sensible choice of the Higgs boson production mechanism. This decay channel is also the subject of study for the analyses of this thesis.

Other important decay channels are the decay into a pair of photons, $H \rightarrow \gamma\gamma$ (0.2% \mathcal{BR}), and into a pair of Z (26.4% \mathcal{BR}) or W bosons (21.5% \mathcal{BR}). With their distinct signatures of two photons or possibly four leptons the $H \rightarrow \gamma\gamma$ and $H \rightarrow ZZ$ decay channels were the two driving forces behind the Higgs boson discovery.

Another very important decay channel is the decay into two tau leptons (6.3% \mathcal{BR}), as it was the first evidence of a coupling of the Higgs boson to fermions.

Other decay channels either have a very small cross section (e.g. $H \rightarrow \mu\mu$) or have decay products that are almost impossible to detect (e.g. $H \rightarrow gg$).

1.3. Associated Production of Higgs Boson and Single Top Quark

The Higgs boson production mechanism studied in this thesis is the associated production with a single top quark. Analogous to the single top production three production channels can be

¹Branching ratios are taken from Reference [33] for a given Higgs boson mass of $m_H = 125.09$ GeV.

introduced: the t -channel process tHq, the tW-channel process tHW, and the s -channel process tHb. The dominant Feynman diagrams for the tHq and tHW production channels can be found in Figure 1.3. Just as in single top quark production at the LHC, the t -channel production has the highest cross section [34, 35] of all three channels, with

$$\sigma(pp \rightarrow \text{tHq})_{\text{SM}}^{8\text{TeV}} = 18.28^{+0.42}_{-0.38} \text{ fb} \quad \text{and} \quad \sigma(pp \rightarrow \text{tHq})_{\text{SM}}^{13\text{TeV}} = 70.96^{+3.00}_{-4.81} \text{ fb},$$

taken from References [34, 36], and is one of the main interests in the analyses of this thesis. The tHW process, which is only considered in the analysis at $\sqrt{s} = 13 \text{ TeV}$, is attributed with the second highest cross section of

$$\sigma(pp \rightarrow \text{tHW})_{\text{SM}}^{13\text{TeV}} = 15.61^{+0.73}_{-1.06} \text{ fb}.$$

The cross section of the s -channel process tHb with $\sigma(pp \rightarrow \text{tHb})_{\text{SM}}^{13\text{TeV}} = 2.8 \text{ fb}$ [37] is negligibly small. A theoretical introduction into the peculiarities of the single top and Higgs boson production, which is also subject of many theoretical studies [34, 35, 38–44], is given in this chapter.

The standard model predicts a large destructive interference between each of the two diagrams of Figure 1.3(a) and 1.3(b), where the Higgs boson is either emitted from the top quark or the W boson, respectively. Possible deviations of the predicted couplings of the Higgs boson to the top quark or the W boson can potentially lead to a significant increase in the production cross section. In order to quantify the deviation from the prediction two real dimensionless scaling factors are introduced:

$$C_f = y_f / y_f^{\text{SM}} \quad \text{and} \quad C_V = g^{\text{HVV}} / g_{\text{HVV}}^{\text{SM}},$$

where y_f is the Yukawa coupling to fermion f , g_{HVV} the coupling of the Higgs boson to a boson V and SM denotes the predicted value by the standard model. The generalized Yukawa coupling to fermions C_f is by far dominated by the top-Yukawa coupling C_t , hence in the following only C_t is used. The Mandelstam variables of the $\text{Wb} \rightarrow \text{tH}$ hard scattering are

$$s = (p_W + p_b)^2, \quad t = (p_W - p_H)^2 \quad \text{and} \quad u = (p_W - p_t)^2.$$

In the high-energy, hard-scattering regime, satisfying $s, -t, -u \gg m_t^2, m_W^2, m_H^2$, the scattering amplitude is given by

$$\mathcal{A} = \frac{g}{\sqrt{2}} \left[(C_t - C_V) \frac{m_t \sqrt{s}}{m_W v} A \left(\frac{t}{s}, \varphi; \xi_t, \xi_b \right) + \left(C_V \frac{2m_W s}{v} \frac{1}{t} + (2C_t - C_V) \frac{m_t^2}{m_W v} \right) B \left(\frac{t}{s}, \varphi; \xi_t, \xi_b \right) \right],$$

where φ is the azimuthal angle of the Higgs boson around the axis, defined as parallel to the direction of the incoming W boson. For simplicity the Higgs boson mass and the bottom quark mass are neglected and terms are left out that vanish in the high-energy limit. The explicit expressions for A and B with their spinors ξ_t and ξ_b can be found in Reference [38].

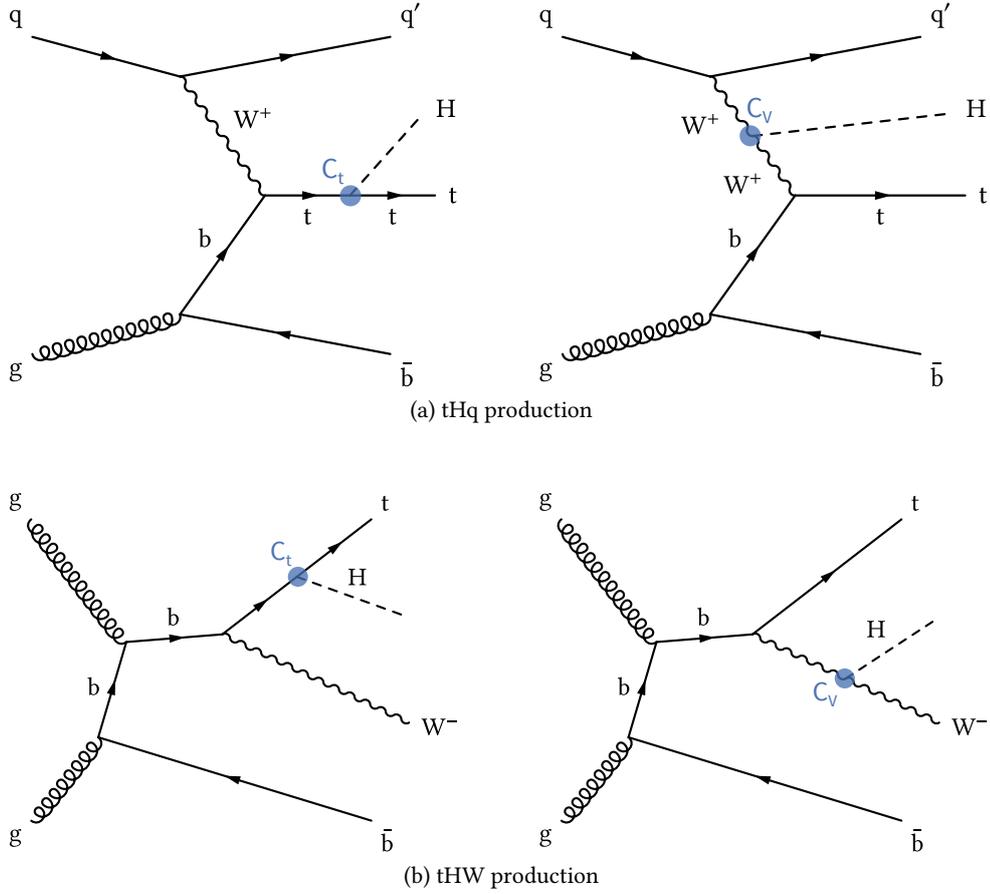


Figure 1.3.: Feynman graphs of the tHq (a) and tHW (b) production mechanism. The graphs are shown in the so-called four-flavor scheme which prohibits the presence of a b quark in the proton and hence requires an initial gluon splitting into $b\bar{b}$ (see Section 3.1.2). The interference between the two diagrams of each channel ultimately leads to the sensitivity of tH production to the sign of C_t .

The standard model predicts $C_t = C_V = 1$ causing the cancellation of the first term and the amplitude to be constant for large s . However, any deviation of C_t and C_V from the prediction comes with an increase of the amplitude, which then increases with \sqrt{s} . The cross section increase of the tHq and tHW processes when moving to other parameters in the C_V - C_t parameter plane can be calculated by

$$\sigma(pp \rightarrow tHq) \sim 3.4 \cdot C_t^2 + 3.56 \cdot C_V^2 - 5.96 \cdot C_t C_V$$

$$\sigma(pp \rightarrow tHW) \sim 1.84 \cdot C_t^2 + 1.57 \cdot C_V^2 - 2.41 \cdot C_t C_V,$$

as stated in Reference [45]. A proof that unitarity is not violated for $C_t = -C_V$ up to a cutoff scale $\Lambda \sim 9.3$ TeV can be found in Reference [38].

The scaling C_V can be fixed to positive values without loss of generality, since $\sigma_{(C_V=1, C_t=-1)} = \sigma_{(C_V=-1, C_t=1)}$.

1.3.1. $C\mathcal{P}$ -mixing in $t\bar{H}$

Recent measurements of the $C\mathcal{P}$ -properties of the Higgs boson disfavor the case of a pseudoscalar state but are not able to exclude it. Interesting for the analyses of this thesis is the assumption of a $C\mathcal{P}$ -violating $Ht\bar{t}$ coupling, on which no stringent constraints can be put yet. Several theoretical studies concerning the $C\mathcal{P}$ -properties of the Higgs boson [40, 41, 43] are available. Assuming a generic spin-0, $C\mathcal{P}$ -symmetry violating particle X_0 with SM-like coupling to the W boson, the effective Lagrangian below the electroweak symmetry breaking scale can be written as

$$\mathcal{L}_0^t = -\bar{\psi}_t (\cos \alpha \kappa_{Ht\bar{t}} g_{Ht\bar{t}} + i \sin \alpha \kappa_{At\bar{t}} g_{At\bar{t}} \gamma_5) \psi_t X_0. \quad (1.2)$$

Here α is the $C\mathcal{P}$ -mixing phase and $g_{Ht\bar{t}} = g_{At\bar{t}} = m_t/v = y_t/\sqrt{2}$, with the Higgs vacuum expectation value $v \sim 246$ GeV. Analogous to C_t and C_V , the dimensionless rescaling parameters $\kappa_{Ht\bar{t}}$ and $\kappa_{At\bar{t}}$ are introduced. This parametrization allows for an easy interpolation between the $C\mathcal{P}$ -even, recovered for $\alpha = 0^\circ$, and the $C\mathcal{P}$ -odd couplings, recovered for $\alpha = 90^\circ$. By setting $\kappa_{Ht\bar{t}}$ to one and α to 0° the SM is recovered.

As elaborated in Reference [43] the choice of

$$\kappa_{Ht\bar{t}} = 1, \quad \kappa_{At\bar{t}} = 2/3$$

leads to a gluon-fusion cross section independent of the $C\mathcal{P}$ -mixing phase α , hence leaving it basically unconstrained by the current analyses. A visualization of the cross section for the t -channel production $t\bar{X}_0q$ and the X_0 production together with a top quark pair $t\bar{t}X_0$ can be found in Figure 1.4. It is apparent that the production in association with a single top quark is sensitive to the exact value of α , whereas the $t\bar{t}X_0$ production cross section is degenerate under $\alpha \rightarrow \pi - \alpha$ and is exceeded by the $t\bar{X}_0q$ production cross section from $\alpha \sim 60^\circ$ on.

A search for the associated single top production with such a spin-0 particle with a $C\mathcal{P}$ -violating coupling to the top quark is performed in Section 6.13 of this thesis. As a variation of α also comes with change of kinematics of all involved objects (see Figure 1.5), a complete analysis is performed exploiting the shapes of kinematic distributions to reach the highest possible sensitivity at 21 different α values.

1.3.2. Experimental Status in $t\bar{H}$ Production

Measurements from CMS and ATLAS have been able to put weak constraints on the scaling factors C_t and C_V via indirect measurements. Especially analyses employing the gluon-fusion Higgs boson production mechanism and searching for the $H \rightarrow \gamma\gamma$ decay mode are sensitive to C_t due to interference of either top quarks or W bosons in the loop needed for coupling the Higgs boson to massless particles. Under the strong assumption that no BSM particles take part in the loops of this diagram the coupling parameter of the Higgs boson to photons can be expressed as stated in Reference [45] as

$$C_\gamma^2 \approx 0.07 \cdot C_t^2 + 1.61 \cdot C_V^2 - 0.68 \cdot C_t C_V.$$

A flipped sign of C_t would hence lead to an increase of the $H \rightarrow \gamma\gamma$ branching fraction of ~ 2.4 , which is not supported by current measurements.

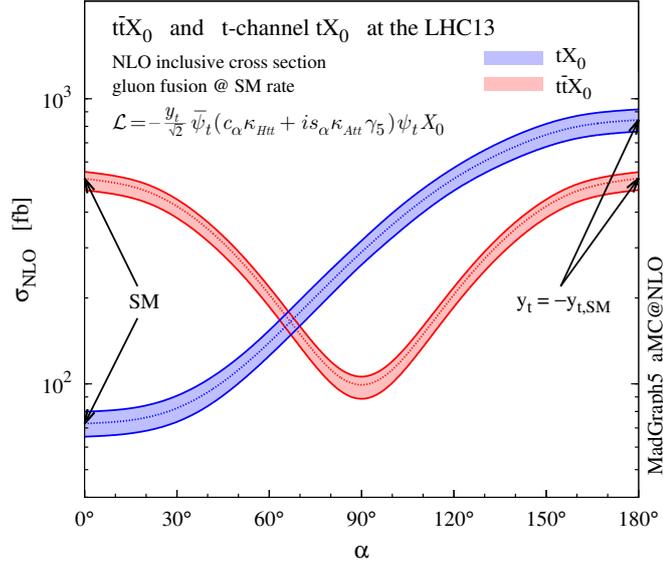


Figure 1.4.: The production cross section for a \mathcal{CP} violating X_0 particle produced in association with a top quark pair (red) and a single top quark (blue) as a function of \mathcal{CP} -mixing angle α . The diagram is taken from Reference [43].

CMS and ATLAS perform global fits, where the results of many different analyses are combined and the different coupling parameters of the Higgs boson are constrained. The most recent measurement is given by a combination of CMS and ATLAS data [45] and the constraints on the C_V - C_t plane can be seen in Figure 1.6². The colored planes show the allowed 68% C.L. regions for different Higgs boson decay modes and it is apparent that, although four out of five considered decay modes have their respective best-fit value for a negative κ_f , the combination favors a point close to the SM prediction. This is mainly caused by the before mentioned incompatibility of the measured $\mathcal{BR}(H \rightarrow \gamma\gamma)$ and the predicted enhancement for $C_t = -1$. However, if the inclusion of non-standard model particles in loop diagrams is allowed, the dependence of C_γ on C_t and C_V becomes largely unknown, making any statement about the sign of C_t much more imprecise. The tH production is almost exclusively able to directly probe the sign of the top-Yukawa coupling.

ATLAS and CMS have followed different paths in the search for tH at $\sqrt{s} = 8$ TeV. The search for $tH \rightarrow \gamma\gamma$ in the ATLAS collaboration [46] considers the tHq process as a signal contribution, thereby achieving sensitivity on the sign of C_t . CMS performs several dedicated tH analyses directly searching for different Higgs boson decay channels. Out of these the analysis searching for $H \rightarrow b\bar{b}$ is presented at $\sqrt{s} = 8$ TeV in Section 5 and at $\sqrt{s} = 13$ TeV in Section 6.

²The combination uses another common, but different nomenclature for the dimensionless scaling parameters with $\kappa_f = C_f$.

1.3. Associated Production of Higgs Boson and Single Top Quark

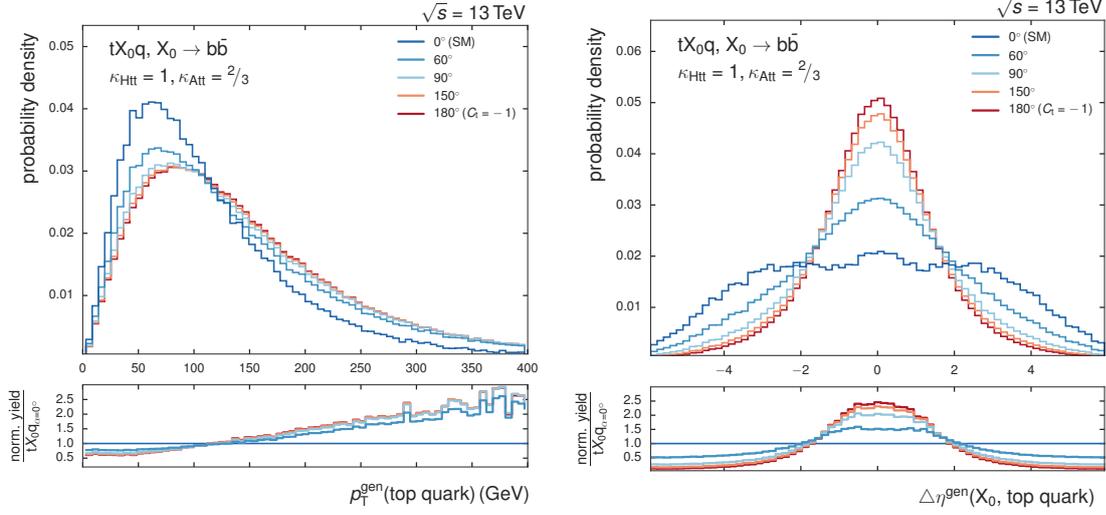


Figure 1.5.: Distributions of the transverse momentum of the top quark and the difference of the pseudo-rapidity of top quark and Higgs boson in the tHq process for different CP -mixing angles.

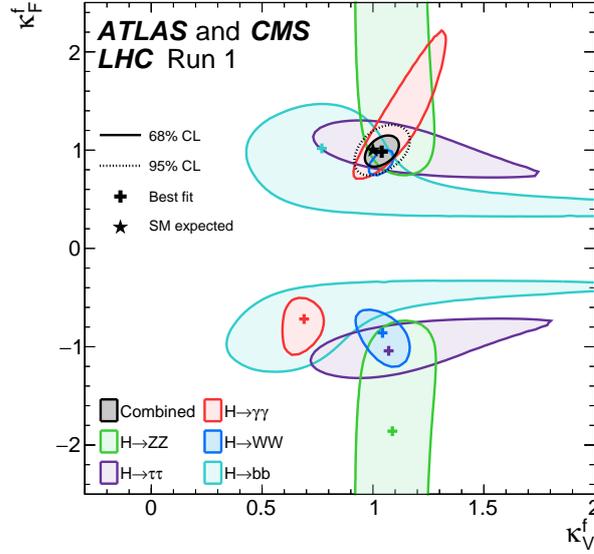


Figure 1.6.: Combined coupling fits of CMS and ATLAS, under the assumption that no BSM particles are resolved in the virtual loops and all coupling parameters for fermions are equal. Whereas four out of five Higgs decay mechanisms favor a negative coupling factor of the Higgs boson to fermions, the measurements are more incompatible for negative values of κ_f , leading to a best-fit value very close to the SM prediction. This incompatibility stems mainly from the $H \rightarrow \gamma\gamma$ decay which is indirectly sensitive to the relative sign of κ_f and κ_V . The diagram is taken from Reference [45].

2. The CMS Experiment at the LHC

In the last century particle physicists have pushed the energy frontier to ever higher regions to discover heavier and heavier particles. Whereas Rutherford studied the properties of the proton in a small laboratory, today's hunt for new particles requires huge research facilities. Energies close to levels reached shortly after the Big Bang are needed to artificially produce the heaviest elementary particles of the standard model.

The Large Hadron Collider (LHC) located at the Conseil Européen pour la Recherche Nucléaire (CERN) in Geneva, Switzerland, provides energies large enough to produce Higgs bosons, top quarks and possibly even heavier particles.

Alongside the acceleration ring four large detectors are located. The two multi-purpose detectors ATLAS (**A** **T**oroidal **L**HC **A**pparatu**S**) [47] and CMS (**C**ompact **M**uon **S**olenoid) [48] are equipped with technology to cover a broad range of different physics interactions. The LHCb detector [49] is specifically designed to measure rare B meson decays and increase the precision in the field of flavor physics. The ALICE (**A** **L**arge **I**on **C**ollider **E**xperiment) [50] detector is especially well suited to record heavy ion collisions.

The analyzed data of the Chapters 5 and 6 has been recorded with the CMS detector. Both, the accelerating structure and the detector, will be explained in detail in the following chapter.

2.1. The Large Hadron Collider

Located roughly 100 m underground the Franco-Swiss border near Geneva lies the Large Hadron Collider. With a circumference of 26.7 km it is for the time being the largest ring accelerator in the world [51]. The LHC itself is the main accelerator of the acceleration complex at CERN, which can be seen in its entirety in Figure 2.1. After 14 years of meticulous planning and building the LHC saw its first collisions on 10th of September 2008. Due to technical complications, the start of continuous data taking was deferred to 2010 and the consensus was reached that the first proton collisions would take place at a lower center-of-mass energy of 7 TeV, half of the original design at 14 TeV. In the first set of years of continuous data taking, also called LHC Run I, the center-of-mass energy was constant at 7 TeV for 2010 and 2011 and was increased to 8 TeV for 2012. The time between consecutive colliding proton bunches, the so-called bunch spacing, was fixed to 50 ns. The time during the first Long Shutdown (LS1) in 2013 and 2014 was utilized to upgrade certain parts of accelerator and detectors and carry out incurred repairs. The LHC Run II started in 2015 with an increased energy per beam of 6.5 TeV, resulting in a center-of-mass energy of 13 TeV, and a bunch spacing of 25 ns.

In its main run mode, protons travel through different pre-accelerators with increasing energies until they reach the LHC. The single proton bunches are split into two contra-rotating

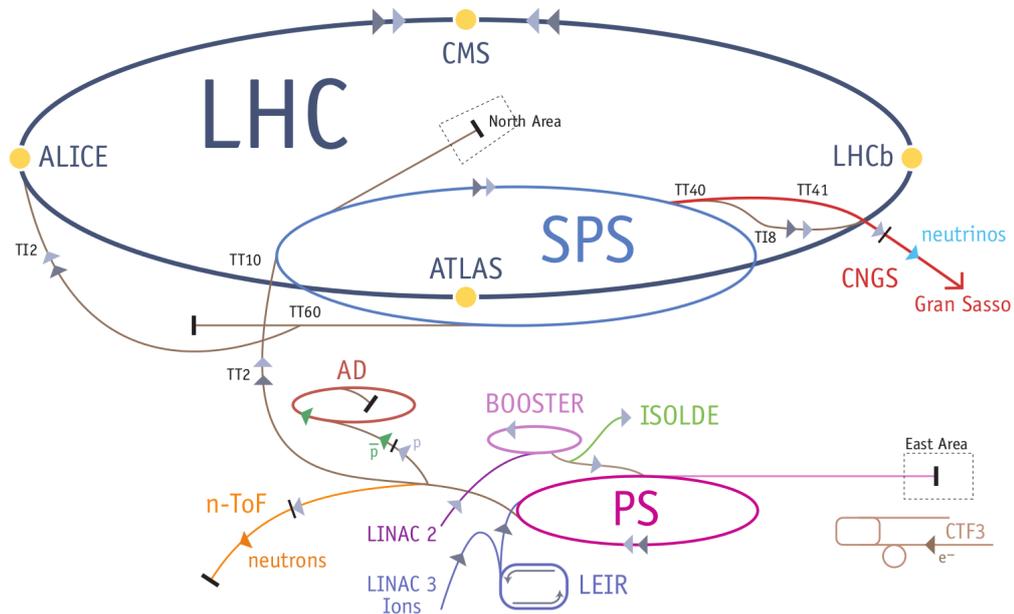


Figure 2.1.: Overview of the LHC accelerator complex, taken from [52]. For the standard proton-proton run mode, protons are accelerated in the LINAC2 accelerator, followed by the PSB (BOOSTER), PS and SPS before they are injected into the LHC. They are brought to collision at the four main experiment sites. When in heavy ion run mode, lead ions of vaporized lead are accelerated in the LINAC3 and the Low Energy Ion Ring (LEIR) before entering the same route as the protons.

beams and their energy is increased further until the planned energy is reached and subsequently the beams are brought to collision.

In its alternative run mode lead ions pass the accelerator chain and heavy ion collisions or mixed proton-lead collisions can be studied.

The acceleration chain starts with a simple hydrogen gas bottle. In a duoplasmatron the hydrogen atoms are stripped of their electrons and the resulting protons are focused and accelerated to an energy of 750 keV in a radio frequency quadrupole (RFQ). The linear accelerator LINAC2 is the next step in the chain where the proton energy reaches 50 MeV. In the Proton Synchrotron Booster (PSB), the Proton Synchrotron (PS) and the Super Proton Synchrotron (SPS) the beam energy is further increased to 1.4 GeV, 26 GeV and 450 GeV, respectively. From the SPS the proton bunches are injected into the LHC in two separate beam pipes allowing for contra-rotating beams. The beams are accelerated to their final energy and are brought to collision at four predetermined crossing points. At these collision points the detectors ALICE, ATLAS, CMS and LHCb are located and observe the products of the collisions.

Eight radio frequency cavities per beam are installed and ensure the longitudinal acceleration of the beam. Inside the LHC 1232 dipole magnets create the magnetic field needed to keep the beam circulating. A cross section of such a dipole magnet can be seen in Figure 2.2. The

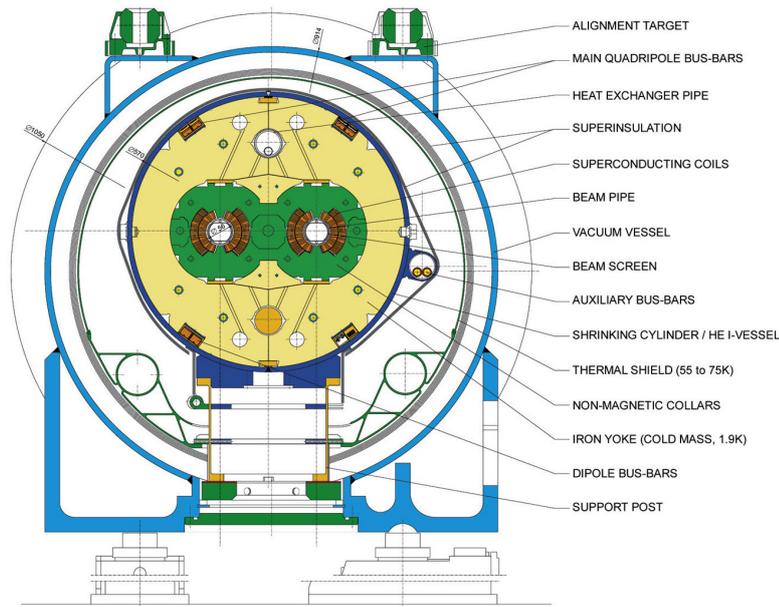


Figure 2.2.: The cross section of one of the 1232 dipole magnets of the LHC. In the center the two separate beam pipes are visible needed for the contra-rotation of the two positively charged proton beams. Taken from [53].

dipole magnets lead to a defocussing of the beam which is countered by 392 quadrupole magnets.

At the design energy of 7 TeV per beam a field strength of 8.33 T would be necessary. The protons are accelerated in an oscillating electro-magnetic field forcing the once continuous proton beam into a bunch structure.

A cooling system operating with super-fluid helium is crucial to keep the magnets at these enormous field strengths in their superconducting states. The respective numbers are taken from Reference [51]. As particle collisions are non-deterministic processes a large number of collisions must be recorded and analyzed to be able to make a statement about the properties of certain particles. As a measure for the collision rate the instantaneous luminosity is used. It can be calculated as

$$L = fn \frac{N_a N_b}{4\pi\sigma_x\sigma_y},$$

with the revolution frequency f , the number of bunches per beam n , the number of particles per bunch N_a and N_b and the Gaussian transverse sizes of the proton bunches σ_x and σ_y . With the cross section σ for a given process the number of events per second can be calculated as

$$\frac{dN}{dt} = L\sigma$$

The integrated luminosity $L_{\text{int}} = \int L dt$ is a measure for the total collected data. In Figure 2.3 the annually by CMS recorded integrated luminosities for Run I can be seen.

The Long Shutdown 1 after Run I offered time for consolidation of different elements of the

2. The CMS Experiment at the LHC

acceleration complex. Work on the accelerators included a consolidation of 10,000 splices, the connections between the superconducting magnets, at the LHC. Additionally, replacements for worn parts, as well as additional shielding in strategic locations, have been installed [54, 55]. At the start of Run II in 2015 the CMS detector faced some early problems with different detector components hindering the early data taking period. Combined with additional problems on the accelerator side lead to lower recorded luminosities than what was expected for 2015. The recorded and delivered luminosities can be found in Figure 2.3. Downtimes and other complications at the detector side lead to a lower recorded luminosity than what was delivered by the accelerator. The datasets which are analyzed in this thesis correspond to integrated luminosities of 20.2 fb^{-1} at $\sqrt{s} = 8 \text{ TeV}$ and 2.3 fb^{-1} at $\sqrt{s} = 13 \text{ TeV}$, respectively.

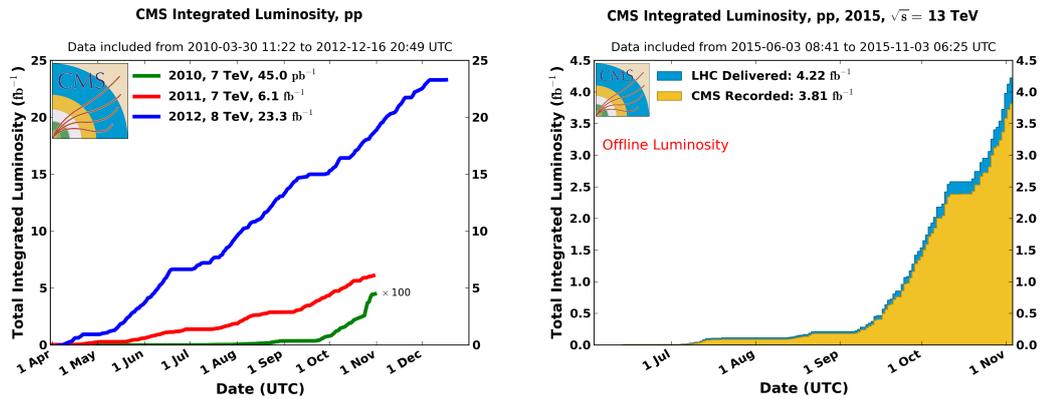


Figure 2.3.: An overview of the collected data in Run I and Run II, respectively. On the left side the total integrated luminosity for the complete Run I can be found. It can be seen by the enhanced green line that in 2010 only 45 pb^{-1} at $\sqrt{s} = 7 \text{ TeV}$ were recorded, as it was made sure that every detail in the LHC was working properly. In 2011 the integrated luminosity grew constantly until 6.4 fb^{-1} were gathered. After the upgrade to a slightly higher center of mass energy of $\sqrt{s} = 8 \text{ TeV}$ CMS collected its major part of data in 2012.

On the right side the integrated luminosity collected in 2015 during Run II is shown. Although a total of 3.81 fb^{-1} have been recorded by the CMS detector, only 2.3 fb^{-1} are utilizable in this thesis due to a problem with the solenoid at that time. These diagrams are taken from [56].

2.2. The Compact Muon Solenoid

The Compact Muon Solenoid detector is a multi-purpose detector located at Point 5 on the LHC accelerator ring. The technical details of the detector are described in great detail in Reference [57]. The information for this section is taken from there, if not otherwise stated. The detector consists of different layers of subdetectors. In the following sections the different components are described from inner subdetectors outwards.

Coordinate System

A right-handed coordinate system is used in CMS with its origin at the interaction point. The x -axis points towards the center of the LHC, the y -axis points upwards. Inevitably, the z -axis points westward, pointing from the CMS detector towards the Jura mountains, tangential to the beam axis. As interactions are invariant under rotations around the beam pipe a polar coordinate system is commonly used. The azimuthal angle ϕ measures the angle from the x -axis in the x - y -plane. The polar angle θ is measured from the positive z -direction.

In particle physics it is prevalent to use the rapidity instead of the polar angle as the rapidity is invariant under Lorentz transformations. It can be calculated as

$$y = \frac{1}{2} \ln \left(\frac{E + p_z}{E - p_z} \right)$$

with the energy E and the momentum in direction of the z -axis p_z . Although correct, handling with energies and momenta makes the rapidity a cumbersome unit. With only minor differences to the rapidity and a simple dependence on the polar angle the usage of the pseudorapidity is much more prevalent. It is defined as

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

The interval for η ranges from 0 for a particle flying vertical to the beam pipe to ∞ for a particle flying alongside the beam pipe. Another commonly used measure is the angular distance ΔR between two objects in the detector. It is defined as

$$\Delta R = \sqrt{(\Delta\eta)^2 + (\Delta\phi)^2}.$$

2.2.1. Tracking System

The tracking system is the innermost part of the detector with the purpose of measuring the tracks of all electromagnetically charged particles produced in a collision. It encloses the beam pipe and the collision point with a length of 5.8 m and a diameter of 2.5 m. A precise measurement of the track curvature due to the magnetic field of the solenoid allows for an estimation of the transverse momentum of the particle and its charge. The tracks of particles from collisions of other particles in the same bunch or collisions from previous bunches will still be visible in the detector, thus making it necessary to measure tracks close to the interaction point as precisely as possible. An accurate determination of the track origin allows for the reconstruction

of their respective vertex.

The detector consists of silicon semiconductors. Charged particles traversing the silicon produce electron-hole pairs. These charges are then collected at the ends of the semiconductor. These electric signals are amplified and read out counting as a hit in the specific detector segment. The tracking system consists of two subsystems, the silicon pixel and the silicon strip detector.

Silicon Pixel Tracker

Three 53 cm long concentric layers of pixel modules and an endcap disk on each side make up the pixel tracker. The pixel tracker is composed of 1440 modules with a total of 66 million silicon pixels. Each pixel is 285 μm thick and spans $100 \times 150 \mu\text{m}$ in $r - \phi$ and z , respectively. The pixel tracker achieves a resolution of $10 \times 20 \mu\text{m}$, making it well suited for the reconstruction of vertices. The active material covers the pseudorapidity range of $|\eta| < 2.5$.

In the upcoming Long Shutdown 2 the pixel tracker will be replaced with an upgraded system comprising four layer of pixel modules [58]. As there was not sufficient space for the fourth layer, the pixel detector was removed during the LS1 and the beam pipe, which is closely encompassed by the pixel detector was replaced by a narrower one. The extraction of the pixel detector made a recovery of previously broken electronic channels possible.

Silicon Strip Tracker

The strip detector environs the pixel detector and consists of 15,148 modules resulting in a total of 9.3 million silicon strips, which cover a total area of about 198 m^2 . The strip detector contains four modules itself: the Tracker Inner Barrel (TIB), Tracker Outer Barrel (TOB), Tracker Inner Disks (TID) and Tracker End-Caps (TEC). Each strip in the four layers of the TIB is 320 μm thick and has a pitch between 80 and 120 μm . Strips with a thickness of 500 μm and a strip pitch from 120 μm to 180 μm are encompassed in the TOB. A single charged particle traversing the strip detector is registered up to ten times, each with a single point resolution of 30 μm to 50 μm .

The TEC consists of nine layers of silicon strips with a mean pitch between 96 μm and 143 μm . The 3 disks of the TID are located in the gap between TIB and TEC. The strips of the TID as well as the three innermost layers of the TEC have a thickness of 320 μm , whereas the rest of the TEC has a thickness of 500 μm .

An accurate description of the layout can be found in Figure 2.4.

In both barrel and endcap stereo modules are encompassed, which allow the 3-D reconstruction of hit positions. These modules actually consist of two back-to-back modules, where one module is rotated through a stereo angle. This way additional information about the missing coordinate can be obtained. These stereo modules attain a resolution in the otherwise inaccessible third dimension ranging from 230 μm to 530 μm .

In order to reduce the effects of radiation damage and therefore increase the longevity, the tracker system is operated in a colder environment at -15°C instead of $+4^\circ \text{C}$ in Run II of the LHC. In order to prohibit condensation on the cooling circuits and detectors, controlling the humidity inside the detector by blowing in dry gas is crucial in Run II [59].

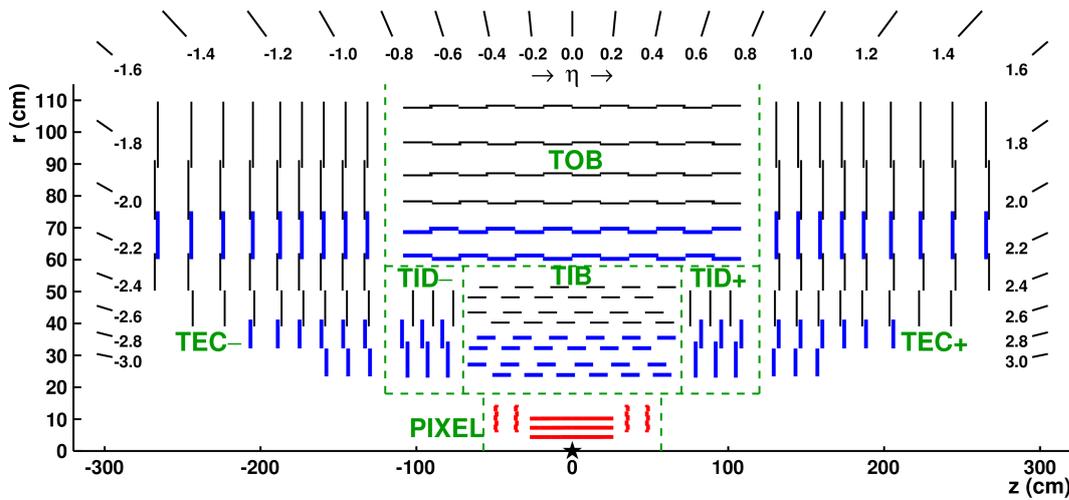


Figure 2.4.: Overview of the CMS tracking system taken from [60]. The pixel tracker (PIXEL) is shown in red surrounding the interaction point (star). The Tracker Inner Barrel (TIB) and Disks (TID) occupy the space $r < 55$ cm and $|z| < 118$ cm. The barrel consists of four layers of silicon strip modules in the TIB and the three layers in each disk. These modules achieve a position resolution in $r-\phi$ of 13 - 38 μm . The Tracker Outer Barrel (TOB) covers $r > 55$ cm and $|z| < 118$ cm and incorporates six modules. Its position measurement resolution is approximately 18 - 47 μm . The Tracker Endcaps (TEC) cover the region $124 < |z| < 282$ cm. Each of the endcaps consists of nine disks with seven concentric rings of silicon strip modules per disk. A similar resolution as in the TOB is reached. The blue lines show the location of the stereo modules.

2.2.2. Electromagnetic Calorimeter

The electromagnetic calorimeter (ECAL) is located outside of the tracking system. Its purpose is to stop electromagnetically interacting particles and measure the energy that is deposited in the ECAL by the particle while doing so. Particles which enter the electromagnetic calorimeter produce an electromagnetic shower due to bremsstrahlung and pair production. The scintillator emits light proportionally to the deposited energy which is subsequently transported through the crystal and read out by photodiodes. The electric signal is amplified and allows for an accurate deduction of the native particle energy. A detailed depiction of the ECAL can be found in Figure 2.5.

The ECAL is a hermetic and homogeneous lead tungstate (PbWO_4) crystal calorimeter. Advantageously, lead tungstate has a short radiation length $X_0 = 8.8$ mm and a small Molière radius $R_M = 22$ mm. The ECAL contains 75,848 crystals and is arranged in a central barrel section (EB) which is closed off by two endcaps (EE). The coverage of the ECAL ranges to a pseudorapidity up to $|\eta| = 3.0$, with the EB covering the central part $|\eta| \leq 1.479$ and the EE the forward part $1.653 \leq |\eta| \leq 3.0$. The crystals are positioned slightly off their direct line to the interaction point to avoid particle tracks hiding in the crystal gaps.

The EB is made up of 36 supermodules containing 1700 crystals each, where as the endcaps consist of two structures, shortly named for their shape as “dees”, with 3662 crystals each.

The EB crystals are 230 mm long, corresponding to approximately 26 radiation lengths, and a transverse front face size of 2.2×2.2 cm². With a length of 220 mm the EE crystals are slightly shorter but have a larger front face size of 2.86×2.86 cm².

Also incorporated into the ECAL is the preshower detector (ES) which improves the separation of actual photons and non-prompt photons of a π^0 decay. It consists of silicon strip sensors and lead absorber parts and is located in front of the endcaps.

During the LS1 a problem with the heating system in the ES was discovered which led to a refurbishment of the ES in general [61]. Additionally, the EE also underwent minor repairs.

2.2.3. Hadron Calorimeter

The hadron calorimeter (HCAL) is a sampling calorimeter with alternating brass absorber and plastic scintillator layers. Brass was chosen as main absorber material as its high density leads to a short interaction length of $\lambda_I = 16.42$ cm and because of its non-magnetivity. The hadron calorimeter works similarly to the ECAL, measuring the particle energies by absorbing them in the material. With its much higher interaction length the HCAL is better suited for measuring the energy of hadrons which do not get absorbed by the ECAL. They produce hadron showers due to inelastic scattering with the detector material. The fraction of energy that is deposited in the scintillator is used to estimate the total energy of the particle. The scintillation light is transported with wavelength shifting optical fibers and subsequently read out by hybrid photodiodes. The sampling structure of the HCAL causes a worse energy resolution than that of the homogeneous ECAL.

The HCAL consists of four different subdetectors: The Hadron Barrel (HB), Hadron Outer (HO),

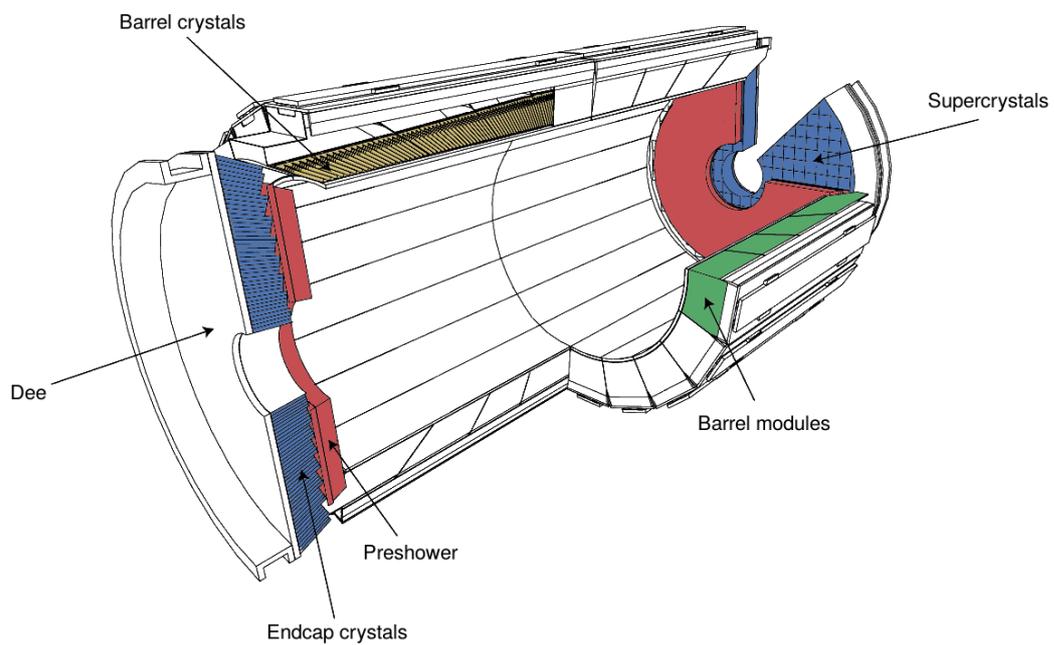


Figure 2.5.: A three-dimensional overview of the CMS ECAL. The central barrel section contains 36 supermodules (green) with each supermodule containing 1700 crystals each (yellow). The barrel endcap region contains four dees with 3662 crystals each (blue). The preshower detectors (red) are located on each side in front of the dees. The original diagram is taken from [48].

Hadron Endcap (HE) and the Hadron Forward (HF) subdetector.

The HB and HF are located inside of the superconducting solenoid. The HB covers the pseudo-rapidity range $0 < |\eta| < 1.4$, the HE covers $1.3 < |\eta| < 3.0$. In order to catch energy passing the HB and the magnet coil the scintillators of the HO are located outside the solenoid cryostat. The HF uses quartz fiber and steel as scintillator and absorption material, respectively. It covers the forward $|\eta|$ range between 3.0 and 5.2, which is particularly important for the analysis of this thesis.

During the LS1 several of the photodetectors of the HCAL were replaced including all of the photodetectors in the HO.

2.2.4. Superconducting Magnet

The tracker system of CMS exploits the curvature of the tracks of charged particles to determine their momentum. The bending of trajectories of charged particles is achieved by a high magnetic field. In the CMS detector the magnetic field is generated by its superconducting solenoid which is built around the calorimetry system. The magnet itself has a diameter of 6.3 m, a length of 12.5 m and can store up to 2.6 GJ at a design field strength of 4 T. The solenoid was operated at a field strength of 3.8 T during data-taking periods relevant for this thesis.

The solenoid is kept at a temperature of 4.7 K in its superconducting state by a 220 t cold mass. The magnetic flux is returned through a 10,000 t iron yoke with an outer diameter of 14 m. Together with the embedded muon system the iron yoke embodies the outermost layer of the CMS detector. The iron also lends structural stability to the complete detector.

2.2.5. Muon System

The eponymous muon system is embedded in the return yoke in the outermost layer of CMS. Together it takes up more than 70% of the total CMS volume.

Muons are the only charged particles traversing all previous detector layers due to their low ionization energy loss and thus provide a unique signature in the detector.

Based on the required large surface coverage and the different radiation exposure, three different gaseous detector types are employed in the muon system. Drift tube chambers (DT), which are generally better suited for low occupancy environments, are located in the central region with $|\eta| < 1.2$, as the muon rate as well as the neutron induced background is low and the residual magnetic field is small. The 250 drift chambers are filled with an Ar/CO₂ gas mixture and are arranged in four layers.

With high muon rate, high neutron induced backgrounds and high residual magnetic field, the endcaps pose a vastly different environment. Cathode strip chambers (CSC) are deployed in this region up to $|\eta| < 2.4$. A mixture of Ar/CO₂/CF₄ is used in the 468 CSCs which make up the muon endcaps.

Both detector types are assisted by three disks of resistive plate chambers (RPC), which are operated in avalanche mode. The RPCs combine a good time resolution of only one nanosecond with a good spatial resolution. A schematic overview of the muon system can be found in Figure 2.7.

During the LS1 the muon system also underwent maintenance where problematic electronics

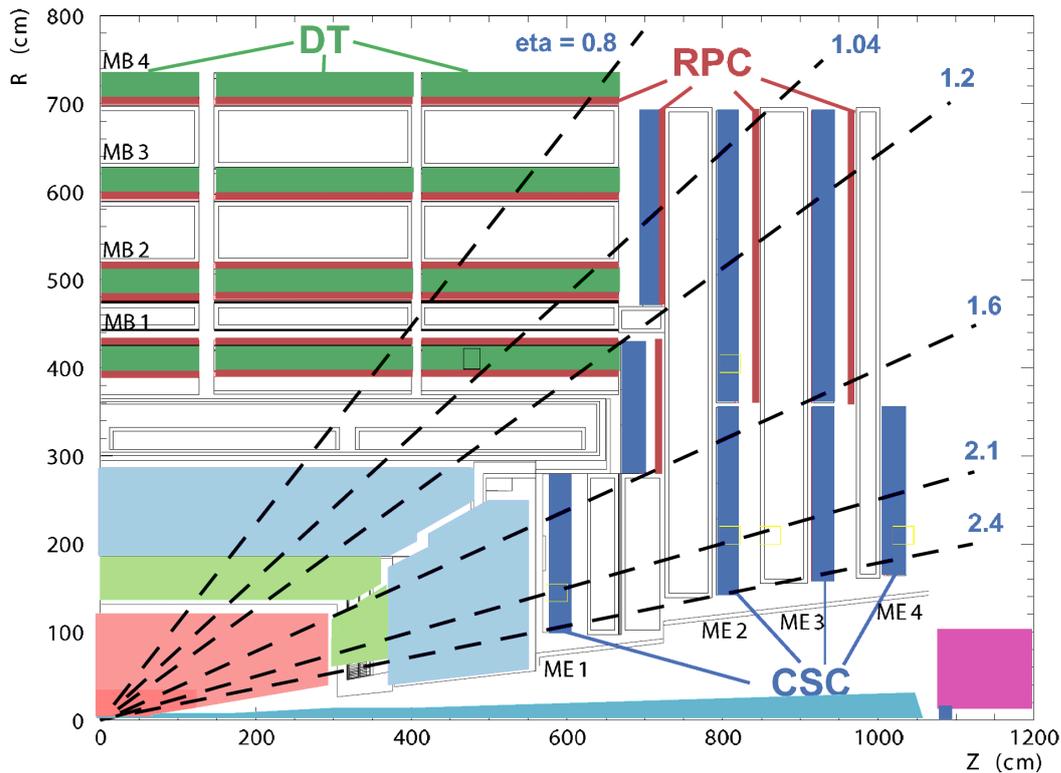


Figure 2.6.: An overview of the CMS muon system. The diagram is adapted from [57]. The three different detector types employed in the muon system are shown. The drift tube chambers (DT, green) cover the central region, whereas the cathode strip chambers (CSC, blue) cover the forward region. Both detector types are assisted by resistive plate chambers (RPC, red).

or detector parts have been repaired. Furthermore, a fourth disk of RPCs was installed as it was originally designed in Reference [57] and an additional set of 72 CSCs was installed in the muon system.

2.2.6. Trigger System

Storing every collision is impossible at rates of 20 MHz (i.e. 50 ns bunch crossing interval in Run I) or 40 MHz (i.e. 25 ns bunch crossing interval in Run II) for the currently available hardware. Combined with a high rate of physics processes that are not of interest for the current research a drastic rate reduction is necessary.

The trigger system is designed to reach a rate reduction factor of 10^6 and thus bringing the data rate down to a manageable level. For an event to be recorded for offline analysis it has to pass the two main levels of the trigger system, the Level-1 trigger (L1) and High Level Trigger (HLT).

The Level-1 trigger is able to achieve a reduction of the event rate to less than 100 kHz. It is comprised of programmable hardware like FPGAs and ASICs. Information from the hit patterns

of the muon chambers and the energy deposits in the calorimeter form the basis for decisions of the Level-1 trigger, whether an event is discarded or passed on to the HLT.

The Level-1 trigger is partly housed directly in the detector itself and partly in the CMS underground control room.

The HLT consists of multiple software filters executed on event data in a processor farm. The first step of the HLT also uses information from the muon detectors and the calorimeters, later steps take advantage of the first provisionally reconstructed objects. Again, a set of selection criteria is applied deciding whether an event is stored for further analysis or rejected. The processable event rate is increased by only reconstructing the minimum amount of detector information needed for the trigger selection [62].

2.2.7. Computing Model

Although the trigger system reduces the event rate drastically the storage and distribution of the selected data still provides an immense challenge. The Worldwide LHC Computing Grid (WLCG) was constructed to cope with the sheer size of the data as well as the high rate at which the data needs to be transferred. An overview of the grid structure of the WLCG can be found in Figure 2.7. The center of the WLCG consists of the Tier-0 site at CERN where the raw recorded data is stored directly. After a first processing the data is then stored redundantly at Tier-1 centers around the world. The Tier-1 sites also offer custodial storage for simulated data samples.

The smaller Tier-2 sites offer additional storage space and the Tier-3 centers provide computing and storage services for end users for their analyses. A well performing computing infrastructure is crucial for analyses, such as the one described in this thesis. The existence of powerful resources allows analysts to optimize their analyses, with almost no compromises due to computing expenses.

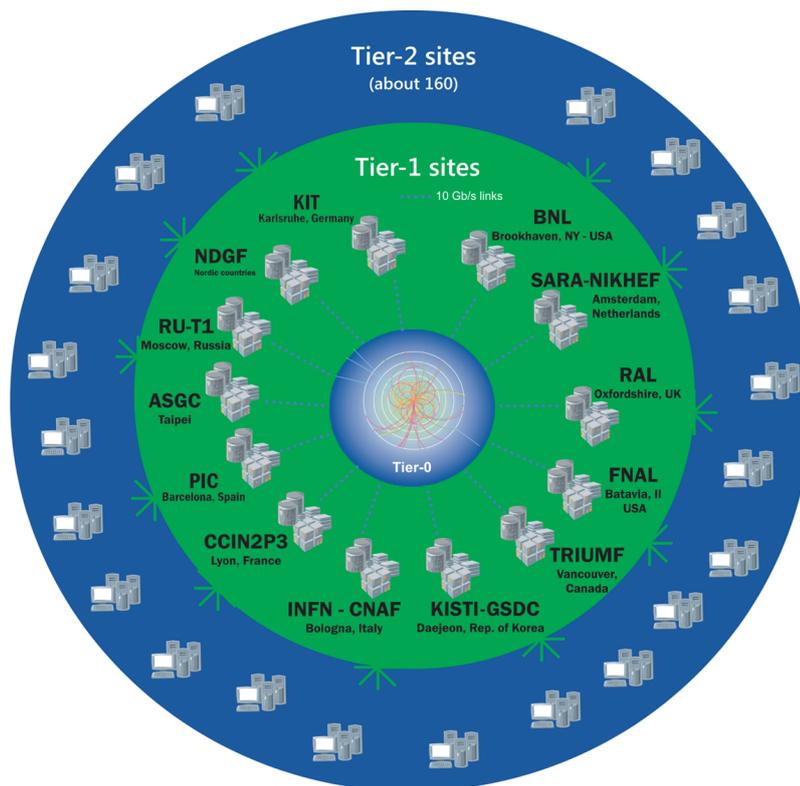


Figure 2.7.: Overview of the Worldwide LHC Computing Grid (WLCG). GridKa, the Tier-1 center at KIT, Karlsruhe, is the only Tier-1 site in Germany. The picture is adapted from [63].

3. Event Generation, Simulation & Reconstruction

Although the standard model (see Section 1.1) is presently the most accurate description of particle physics, the quantum nature of these interactions makes it impossible to precisely predict single processes in a proton-proton collision. However, the application of the rules of the standard model allows for calculation of probabilities, such as decay or interaction probabilities. This renders the study of single events almost pointless, but allows for the prediction of distributions of different observables by studying samples with many events.

In order to find deviations from the predicted, millions of events are simulated with so-called Monte Carlo generators applying an approach based on random sampling. These simulated events are subsequently subjected to an accurate virtual modeling of the CMS detector.

The details of event generation and a description of all involved event generation packages are the focus of the first half of this chapter.

The second part covers the reconstruction of physical objects based on recorded electronic signals by the detector. These signals must be read out, combined and validated in order to map the measured signals to actual particles which appeared in the detector. The techniques used to reconstruct physical objects are explained in the latter part of this chapter.

3.1. Event Generation

The quantum mechanical nature of high energy particle physics does not allow for accurate predictions of single collisions. Yet by employing Monte Carlo methods, a class of computational algorithms based on random sampling, a high number of simulated events mimic the behavior of a high number of actually recorded events. Whereas single simulated events might never happen exactly that way in the detector, the complete stochastic sample of events is able to reproduce distributions of physical observables.

The generation of events can be treated as several steps independent of each other, a concept known as factorization. Firstly, the actual process of the hard scattering is simulated, restricted to the processes which are the subject of investigation. The final parton content of the hard interaction is then subjected to the parton shower which adds radiations from accelerated particles carrying color charges to the event. Due to confinement, particles with a net color charge are not physical and are recombined in the hadronization step. The decay of unstable particles is also simulated in the same stage. An overview of all processes involved in the event generation can be seen in Figure 3.1.

Each of these parts is explained in detail in the following section. The information from this

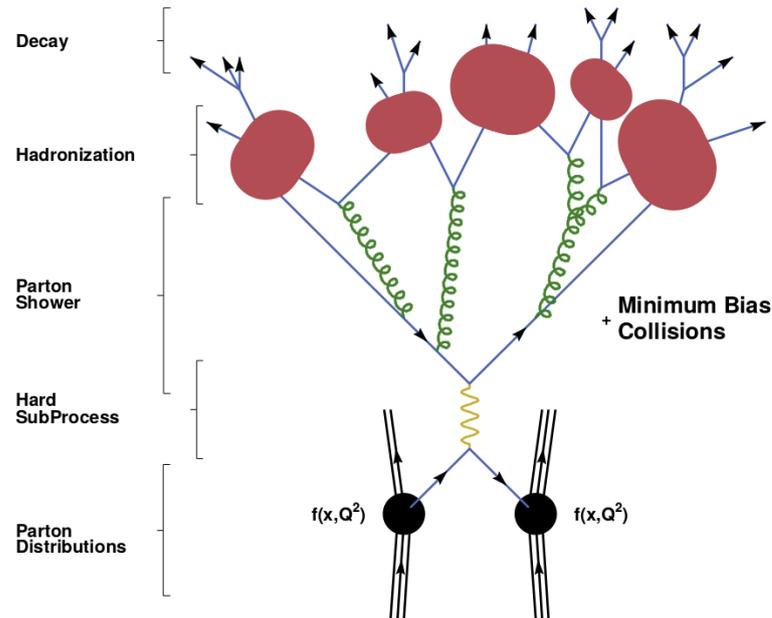


Figure 3.1.: Overview of the event generation process. The diagram shows the different steps of the event generation that can be done sequentially due to the factorization. The parton distribution functions provide the probabilities to colliding partons with a certain fraction of the proton momentum. The hard subprocess describes the actual particle interactions that are the subject of study. During the generation of the hard subprocess the matrix element and therefore the production cross section is calculated. The parton shower is responsible for the generation of gluon radiation. In the hadronization step, all color-charged particles are recombined to form neutral bound states, which subsequently decay in other particles. As last step, the influence of pileup is modeled by overlaying the events with collisions, based on minimum bias collisions. For clarity the underlying event, the influence of the proton remnants, is not depicted in this diagram. Adapted from [64].

section is based on References [64,65], if not otherwise stated.

3.1.1. Hard Scattering

The model of two colliding protons is not sufficient at energies of the LHC as actually the partons inside the proton take part in the deep-inelastic scatterings. Unfortunately, the momentum information of the proton constituents at the time of impact is inaccessible.

However, sets of parton distribution functions (PDFs) provide the probability to find a specific parton with a certain longitudinal momentum fraction x at a resolution scale Q^2 inside the proton. The evolution of the PDFs cannot be calculated entirely perturbatively and thus the DGLAP QCD evolution equation [66–68] is employed.

Further input for the PDFs is provided by many deep-inelastic scattering experiments, partly already long shut down, such as HERA at DESY or the Tevatron at Fermilab. Sets of PDFs are provided by several groups, the most famous and widely used are the sets of the CTEQ [69], the

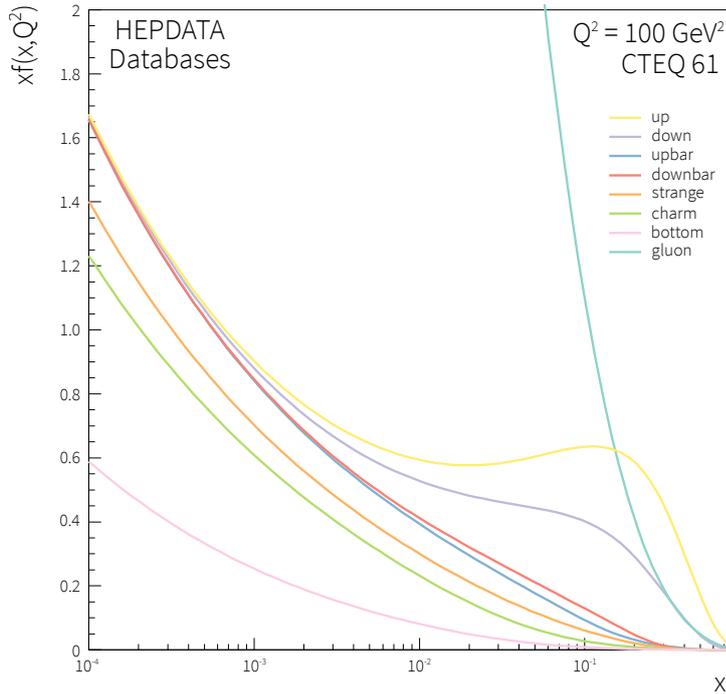


Figure 3.2.: The distribution shows the parton distribution function for up-, down-, antiup-, antidown-, strange-, charm- and bottom quark and the gluon as provided by the CTEQ 6.1 PDF set at $Q^2 = 100 \text{ GeV}^2$. The template for this diagram has been produced with the help of Reference [73].

MSTW [70] and the NNPDF [71] collaborations. The most important sets of PDFs are combined in the LHAPDF package [72]. An exemplary PDF set can be seen in Figure 3.2.

By utilizing quantum field theory the matrix element (ME) and thus the cross section of a certain process can be calculated by considering all possible Feynman diagrams which contribute to the process. As these interactions are usually high-energetic and hence the strong coupling constant α_s is small, perturbation theory can be applied. The ME calculation takes possible interference effects between diagrams with same initial and final state particles into account and covers the decay of particles carrying spin information. These calculations are performed by matrix element generators, which are described later in this chapter.

The accurate calculation of the production cross section for a given processes requires the consideration of radiation corrections. Due to the relative size of α_s over the electromagnetic coupling, QCD radiations play the dominant role.

The generation of physics processes at leading order (LO) does not include additional radiations. Cross sections at LO are usually a good reference point, but do not reproduce the actual physics. The simulation of next-to-leading-order (NLO) processes which contain one hard emission or one loop in the Feynman diagrams can be directly done in the event generators, but at the cost of increased computing time due to the increased number of Feynman graphs that need to be

considered.

The direct generation of events with a given number of radiations, might seem like a solution for this problem, but creates problems that are addressed in the next section.

3.1.2. Parton Shower

The parton shower is responsible for the simulation of gluon radiation stemming from accelerated color charges at lower energy levels. All gluons eventually decay into quark-antiquark pairs or, as they are color-charged themselves, emit further gluons. These processes repeat themselves, consequently creating a whole shower of particles. With each radiation the energy of the particles decreases, therefore the energy eventually enters a domain where α_s can no longer be assumed to be small causing perturbation theory to fail. Hence, this simulation step relies on approximation methods which are tuned in a way to reproduce results observed in data events.

Successive splittings of partons into two new particles are parametrized with the Altarelli-Parisi splitting functions [68], whereas the Sudakov factors [74] are used in the calculation of the probabilities that a particle does not radiate off another particle.

Radiations can be categorized as either initial state radiation (ISR) or final state radiation (FSR), depending on whether the radiation took part before or after the hard interaction.

When simulating additional radiations for processes at NLO accuracy the potential danger of double counting has to be considered. The first radiation added by the parton shower to a LO process would overlap with the born-level diagram of a NLO process. The second radiation at LO overlaps with the first parton shower radiation of a NLO process and with the born-level diagram at next-to-next-to-leading-order (NNLO), and so forth. This problem can be solved by the application of matching algorithms. The most prominent algorithms at the time are the CKKW [75], the MLM [76] and the FxFx method [77].

b Quark Initiated Processes

The production of single top quarks or other processes that involve b quarks in the production process share a peculiarity: the mass of the b quark is higher than that of the proton. Two viable generation approaches exist to deal with this issue, the *four-flavor scheme* (4F) and the *five-flavor scheme* (5F) [78]. The 4F scheme argues that b quarks that take part in the production can only stem from a high energetic gluon splitting into $b\bar{b}$. Thus, the b PDF is set to zero, as the b quark cannot be part of the proton. The 5F scheme, on the contrary, treats the b quark as a massless particle with a non-zero probability to be found inside the proton. These approaches are inherently different, as the second b quark which does not take part in the production process but has to be present in both approaches is generated at LO by the ME generator in the 4F scheme ($2 \rightarrow 3$ process) and by the parton shower in the 5F scheme ($2 \rightarrow 2$ process). Therefore, the 4F scheme allows for a more accurate description of this second b quark, yet introduces possible infinities for large momentum transfers that would cause perturbation theory to break down. This effect is mitigated in the 5F scheme by absorbing these infinities into the b PDF via the DGLAP evolution equations. It is found that processes produced with

the 4F scheme describe the actual measured process best and therefore the 4F scheme is widely used in this thesis.

3.1.3. Hadronization and Decay of Unstable Particles

The parton shower only simulates the individually produced partons, but color confinement plays a dominant role at these lower energies, where partons cannot be treated as quasi-free particles anymore. The color-charged particles have to be recombined forming colorless bound states. This procedure of producing hadrons, known as hadronization, is tuned to mirror the physics results actually measured in earlier experiments.

The most important hadronization model is the Lund string model [79]. The color-carrying gluons are treated as self-attracting field lines which subsequently form flux tubes between colored particles. If the energy of a tube reaches a critical point a new quark-antiquark pair with new flux tubes is produced, thereby depleting the energy of the tubes in total.

Most of the baryons and mesons produced in the hadronization step are unstable and hence the short-lived resonances are decayed in the detector. This decay is simulated based on the known branching ratios of these resonances (see Reference [13]).

3.1.4. Underlying Event and Pileup

The hard scattering of two partons does not leave the rest of the proton unaffected. The colored remnants of the proton are also part of the hadronization step and their effect on the hadronization, known as underlying event, has to be evaluated as well, as they are still color-connected to the initial partons.

As concluding step pileup events (PU) need to be added to give an accurate emulation of the actual collisions. Pileup events can be grouped in two categories:

In-time pileup Collisions of protons that happen in the same bunch crossing, but have a lower vertex quality than the main primary vertex of the event.

Out-of-time pileup Collisions of protons from former or latter bunch crossings which are still visible in the detector due to a time delay in the detector components.

3.1.5. Monte Carlo Generators

The previously described steps of the event generation are usually performed by specialized software packages.

In the following section all generators which were involved in the production of the simulation samples used in this analysis are described.

MadEvent and MadGraph

MADEVENT is a tree-level event generating software based on the matrix element generator MADGRAPH, commonly only referred to as MADGRAPH [80, 81]. For a given process MADGRAPH calculates the amplitudes for all contributing Feynman diagrams. This information can

3. Event Generation, Simulation & Reconstruction

then be evaluated by `MADEvent`, which allows the user to calculate cross sections or decay widths and produce unweighted events. `MADGraph` does not provide a parton shower and is therefore often interfaced with `PYTHIA`.

Implemented in the `MADGraph` framework is the `MADSPIN` [82] tool, which is used for the decay of particles carrying non-zero spin.

`MadGraph5_aMC@NLO`

`MADGRAPH5_AMC@NLO` [83] is an automated next-to-leading order event generator including an optional parton shower step. It can calculate cross sections with full QCD corrections for a user given process, generate the hard process and is able to consistently match processes with radiations at tree-level with possible radiations of the parton shower.

`MADGRAPH5_AMC@NLO` was a merger of the two event generators `MADGRAPH5` and `AMC@NLO` superseding both packages and is currently one of the most frequently used matrix element generators.

A unique feature of `MADGRAPH5_AMC@NLO` is the presence of negative event weights arising from the usage of counterterms that are needed to smoothen the phase space transitions between matrix element and parton shower dominated parts. Negative weights reduce the effective number of events that are used to reproduce object property shapes by filling histograms, which could lead to discontinuous shapes in exotic phase spaces with a low number of events.

`Powheg`

`POWHEG` [84–86], short for **Positive Weight Hard Emission Generator**, is an NLO matrix-element generator. `POWHEG` models the hardest emission of color charged particles in an NLO process, making it necessary to interface it with a p_T -ordered parton shower or a parton shower able to veto this highest emission, as else double counting of this highest-energetic emission would occur. This feature makes `PYTHIA` a natural match for the `POWHEG` event generator.

`Pythia`

`PYTHIA` is a parton shower based multipurpose event generator. The calculation of the hard matrix element is performed at leading order. The parton shower uses a p_T -ordered emission algorithm [87] and the Lund model for the hadronization.

Highly optimized parameter sets, known as tunes, are necessary to reproduce actual collisions to a high degree of precision. The underlying event tune used in all samples for the analysis at $\sqrt{s} = 8$ TeV is the Tune Z2* [88], whereas in the analysis at $\sqrt{s} = 13$ TeV the CUETP8M1 Tune [89] is employed in all samples.

In the Run I analysis of this thesis `PYTHIA` 6.4 [90] was used. During the first Long Shutdown the new version `PYTHIA` 8.2 [91] was released and became the new standard in CMS for parton showers.

As LO matrix-element generation is often not sufficient anymore, `PYTHIA` is often solely used for the parton shower and is interfaced with other generators which take care of the matrix-element generation.

Tauola

Due to the narrow width of the τ lepton its decay can be treated separately from its production. In the Run I analysis the `TAUOLA` package [92] is applied, which only simulates the decay of the τ lepton and has hence to be interfaced with other event generators. Spin effects as well as electroweak corrections are taken into account by `TAUOLA`.

3.2. Detector Simulation

The events produced by the Monte Carlo generators are not yet comparable to the events recorded by particle physics detectors at the LHC. The simulated events need to be subjected to an accurate modeling of the CMS detector. Interactions with the detector material and the effect of the magnetic field on the particles are simulated in this step.

The CMS detector is simulated with `GEANT4` [93, 94], a toolkit for the simulation of particles passing through matter. `GEANT4` not only describes the interactions of the particles with the material budget of CMS, but also simulates the electronic signals that would be measured by all sensors inside CMS.

After the detector responses are simulated it is possible to directly compare the generated events with events actually recorded in the CMS detector.

3.3. Event Reconstruction

After the readout of millions of electronic channels, the signals need to be combined in a sensible way. A coherent statement about e. g. an electron should at best involve all electronic signals caused by an actual electron in the detector. The particle-flow algorithm [95] tries to achieve this by bundling signals throughout the different detection layers in an optimized way, identifying and reconstructing them as physical objects.

3.3.1. Particle-Flow Algorithm

The aim of the particle-flow event reconstruction algorithm is to reconstruct and identify all stable particles based on their different traces left in the detector. After the reconstruction the individual particles can be combined into larger objects like jets or other observables that can be calculated based on the particle content of the detector, such as the missing transverse energy.

The particle-flow algorithm employs an advanced tracking algorithm, the *Iterative Tracking*. Hits produced by charged particles in the silicon detectors are seeded and reconstructed to tracks under strict quality criteria. Hits which can be clearly assigned to a complete track are removed from the set and the reconstruction is repeated with continually looser track seeding criteria. In the fourth and fifth iteration the constraints on the vertex are relaxed, allowing for a reconstruction of secondary charged particles.

After the reconstruction of the tracks in the silicon detector, the calorimeter clustering is performed. In the calorimetric subdetectors, *cluster seeds*, energy maxima above a certain threshold,

are located in the calorimeter cells. Starting from these seeds, other adjacent cells containing an energy above the noise threshold are added to the initial seed to form *particle-flow clusters*. Neighboring cells can be assigned to more than one seed and if so, the energy in the cell is shared among the clusters.

Tracks in the silicon detector, tracks in the muon system or particle-flow clusters are the basic building blocks of the particle-flow reconstruction and are therefore known as *elements*. Single particles give rise to different elements in the detector, thus a linking of these elements is required in order to reconstruct a particle correctly. The linker algorithm links two elements together on a trial basis and only keeps elements linked, if their quality surpasses a certain threshold. The linker tries to extrapolate tracks starting from the outermost hit in the silicon detector to either, the first two PS layers, the ECAL in depth of a typical maximum of an electromagnetic shower or the HCAL at a depth of a typical hadron shower length. If the extrapolated track is in correspondence up to a certain link distance with a calorimeter cluster the two elements are linked together. When linking tracks to ECAL clusters the trajectory of the track is extrapolated along its tangent, looking for energy deposits by bremsstrahlung.

Elements directly or indirectly linked with each other are known as *blocks*. The high granularity of the CMS detector causes most blocks to only contain one to three elements and therefore the object reconstruction is simplified. The linking algorithm, as well as the whole particle flow algorithm, was intensively validated and commissioned [96] as e. g. a broken link could lead to a reconstruction of a ghost particle and thus to an overestimation of the total energy of the collision.

In the last reconstruction step, the found blocks in the detector are interpreted as *candidates*, as they are reconstructed as actual physical objects.

3.3.2. Vertex and Track Reconstruction

The determination of the point of collision is crucial for many other reconstruction or analysis methods, be it an improvement to the track fitting, the separation of pileup events or the b tagging algorithm. A vertex is reconstructed in three steps [60, 97]: the selection of tracks, the clustering of tracks stemming from the same vertex and the fitting to obtain the actual vertex position.

During the track selection a set of strict criteria is imposed on the measured tracks in the detector. For a track to be considered, it must have more than one hit in the pixel layer and more than four hits in pixel layer and strip detector combined, the quality of the fit to the trajectory has to surpass a given threshold ($\chi^2 < 20$) and the significance of the transverse impact parameter relative to the beam spot center has to be lower than five.

Earlier in CMS's lifetime the track clustering was based on the track's z -coordinate at their closest distance to the beam spot. This was known as *gap clustering* [98]. This was superseded by a *deterministic annealing* (DA) algorithm. A detailed description of the algorithm can be found in Reference [99].

After the identification with the DA clustering the remaining tracks are fitted with an *adaptive vertex fitter* [100], which provides the best estimate of the vertex location.

The reconstructed vertices, known as *primary vertices* (PV), are ordered by their respective squared sum of the p_T of all associated tracks. The vertex with the highest sum is chosen as

main vertex of the collision. This choice also settles which tracks and signals in the detector are the ones of interest and which are seen as noise.

3.3.3. Muon Reconstruction

As first step of the particle reconstruction and identification, muons are reconstructed. If the combined momentum of a global muon agrees with the in the tracker determined momentum within three standard deviations, it gives rise to a *particle-flow muon*. Possible energy deposits in the ECAL or HCAL have been estimated with cosmic muons to be 0.5 GeV and 3 GeV, respectively, with an uncertainty of $\pm 100\%$.

3.3.4. Electron Reconstruction

The reconstruction of electron candidates starts with a pre-identification stage in the silicon tracker. Electrons already interact electromagnetically in the tracker layers, thus causing them to produce characteristically short tracks and to lose energy due to bremsstrahlung in the tracker layers. The pre-identified electrons are then refit with a Gaussian-Sum Filter (GSF) [101] attempting to connect them to their corresponding ECAL cluster. The electron track is subsequently removed from the detector together with its associated ECAL cluster as well as its linked bremsstrahlung blocks. Together they form a *particle-flow electron*.

3.3.5. Photon and Hadron Reconstruction

After removing all *particle-flow electrons* and *particle-flow muons*, the remaining tracks are assigned to charged hadrons. The remaining clusters can additionally be caused by photons and neutral hadrons. By applying a more stringent track quality requirement, tracks are subsequently connected to clusters in the ECAL and HCAL. The number of neutral particles contributing to the calorimeter clusters is estimated by comparing the deposited energy in the calorimeters with the momentum of the associated tracks. As tracks can be assigned to multiple clusters, only the links to their closest clusters are kept intact.

Tracks contained in blocks give rise to a *particle-flow charged hadron* associated with the momentum and energy of the track under the assumption of it being a charged pion. The momentum of the charged hadron is reestimated, when track momentum agrees with the energy deposition in the calorimeter within measurement uncertainties. This is especially important at high energies or at high pseudorapidities due to the worsened energy resolution in these regions. If the energy of the linked clusters is significantly higher than the momentum of the charged-particle track and the excess is above the expected calorimeter energy resolution, a *particle-flow photon* and possibly a *particle-flow neutral hadron* is created. If the energy excess is larger than the total ECAL energy of the cluster, the ECAL energy is assigned to the photon and the remaining energy is credited to a neutral hadron. If the excess is lower than the ECAL energy the energy is only credited to a photon. Photons are prioritized here, as roughly 25% of jet energy is credited to photons, but neutral hadrons only deposit 3% of their energy in the ECAL.

Remaining unmatched clusters in the ECAL or HCAL also give rise to a photon or a neutral hadron, respectively.

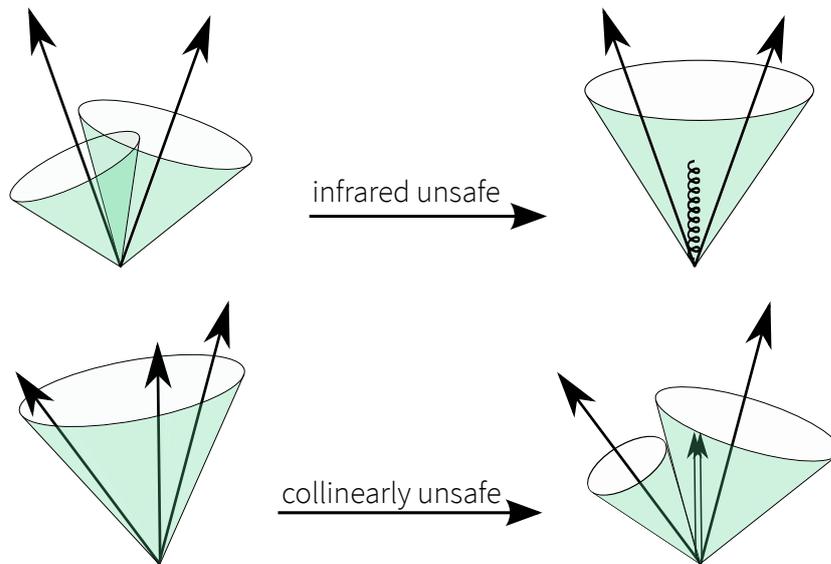


Figure 3.3.: A depiction of infrared (top) and collinear (bottom) safety violations. As shown on the top an additional radiation is able to change the outcome of an infrared *unsafe* jet clustering algorithm. The diagram on the bottom presents how a collinear splitting can alter the jet clustering, if collinear safety is not provided. Both safety requirements are satisfied by the employed anti- k_t algorithm.

3.3.6. Jet Reconstruction

Strongly interacting particles produced in the collisions lead to a shower of hadrons and their decay products. A precise study of every constituent of these showers is unfeasible and unfit. By grouping all these particles into one larger object, the so-called jet, the handling becomes much easier and almost the complete physics content is contained. The following section covers different jet clustering algorithms, as well as corrections that are applied to jets.

Jet Clustering Algorithms

An ideal clustering algorithm would exclusively group all particles that originated from the same parton into one jet. In reality jets originally stemming from different partons overlap in the detector and produce signals in the same calorimetric cells. Since the hadronic activity at the LHC surpasses everything seen in previous experiments a careful treatment of the jet clustering is necessary and two self-imposed theoretical concepts are to be satisfied: The concepts of *infrared safety* and *collinear safety*. A description of both concepts can be found in Figure 3.3. As mentioned in Section 3.1.2 the modeling of low-energetic gluon radiations and small-angle radiations are theoretically very difficult to model. The concept of infrared safety forbids that the clustering of jets changes when soft radiations are added to the jet. The concept of collinear safety implies that the clustering of a jet has to be invariant under a collinear splitting of an initial object.

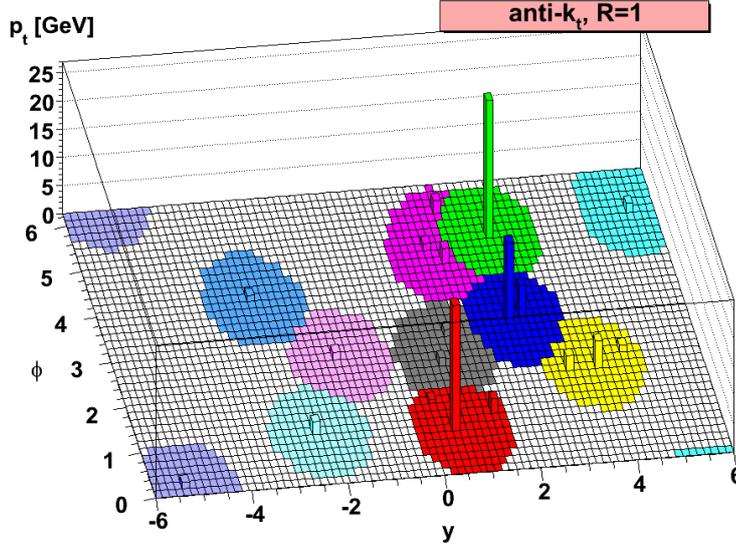


Figure 3.4.: A depiction of a jet clustering employing the anti- k_t algorithm. The diagram shows a flat projection of the HCAL with energy deposits shown as towers. The anti- k_t algorithm produces cone-like shaped jets and successfully resolves overlapping jets. Taken from [102].

anti- k_t clustering algorithm

Throughout the years several jet clustering algorithms have been developed that are infrared and collinear safe. The standard jet clustering algorithm used in CMS and also in the analyses of Chapters 5 and 6 is the anti- k_t algorithm [102].

The anti- k_t algorithm is a sequential recombination algorithm obtained by generalizing the previously existing k_t [103] and Cambridge/Aachen algorithms [104]. An exemplary jet clustering employing the anti- k_t algorithm can be found in Figure 3.4. The clustering is performed by introducing distances between two particles and the distances between particles and the beam. The distance measures are defined as

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

and

$$d_{iB} = k_{ti}^{2p} \quad ,$$

where $\Delta_{ij}^2 = (y_i - y_j)^2 + (\phi_i - \phi_j)^2$ is the distance in the detector, k_{ti} the transverse momentum, y_i the rapidity and ϕ_i the azimuthal angle of particle i . The distance to the beam is denoted with the index B and p is an algorithm specific parameter.

The distances between particles d_{ij} and the distance to the beam d_{iB} are calculated iteratively. If the smallest distance is found between two particles they are combined into a new object and the distances are recalculated. If the smallest distance is found to be between an object and the beam, the clustering ends for this jet and all particles associated to this jet are removed from

the detector. This procedure is repeated until no objects are left in the detector.

The choice of parameter p is crucial as it changes the priorities of the algorithm. By setting the parameter p to -1, one retrieves the anti- k_t algorithm. With the values 1 and 0 for p one recovers the k_t and the Cambridge/Aachen algorithm, respectively. Although the choice of a negative p seems unintuitive, the anti- k_t algorithm produces cone-shaped jets and is infrared and collinear safe. The choice of a negative p causes soft particles to preferentially cluster with harder particles instead of other soft particles, as d_{ij} between a hard and a soft particle is dominated by the momentum of the hard particle.

The parameter R can be chosen freely, as it is a measure for the maximum cone radius. In the analyses in this thesis size parameters of $R = 0.5$ and $R = 0.4$ are chosen producing the so-called *AK5-jets* and *AK4-jets*.

Calculating all distances when running the algorithm can take up a lot of resources as the number of calculations for an ensemble with N particles would scale with N^3 [105]. The implementation of the algorithm in the FASTJET package [106] however, is able to reduce the complex clustering to a two-dimensional nearest neighbor problem, therefore reducing the scaling factor to only $N \ln N$. This allows for a fast application of the anti- k_t algorithm in the high-particle multiplicity environment of the LHC.

Jet Energy Corrections

The study of QCD interactions in the detector becomes feasible after the abstraction of thousands of hadrons into jets. However, this comes at first at the loss of generalization, as the CMS detector is not homogeneously instrumented and the same particles would result in different jets in different regions of the detector. In order to compare the reconstructed jets to predictions and to make them comparable to each other a set of corrections has to be applied. The complete correction chain can be seen in Figure 3.5.

The CMS collaboration uses a factorized model to apply jet energy corrections (JEC) [107–109]. Different correction levels are applied in a fixed order, where at each level the jet four-vector is multiplied by a scaling factor based on jet properties. All correction levels applied in this thesis are explained in the following paragraphs.

Pileup Correction Pileup events (see Section 3.1.4) pollute the detector environment and distort jet energy measurements. The *Level-1 corrections* (L1) aim to remove the surplus of measured energy which is dependent on the number of concurrent events in the detector. The offset corrections are determined with a simulated QCD dijet sample with and without overlaid pileup events. The corrections are derived as function of the energy density ρ , jet area, jet pseudorapidity and jet transverse momentum.

The dedicated L1 corrections for data events account for residual differences compared to the detector simulation as a function of the pseudorapidity. They are determined with a random cone method in zero-bias events, the events with minimal trigger restrictions [110].

MC-truth Corrections In order to infer the original parton energy from the jet energy and to compare jets with each other, the detector response in different pseudorapidity regions and for different transverse momenta has to be equalized. In Run I these corrections were facilitated in two steps, the *L2Relative* correction, which adjusted the detector response to

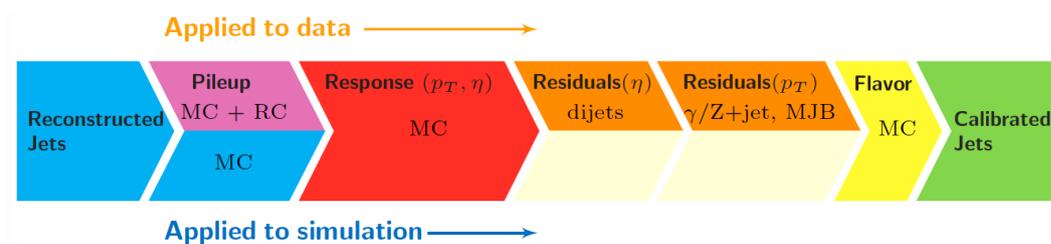


Figure 3.5.: Overview of the different jet energy correction levels. The Level-1 corrections negate the influence of pileup jets in the energy measurements and are determined in MC simulations. Additional data correction factors are derived with a random cone method (RC). Level-2 and Level-3 corrections improve the varying detector response. Further, Level-2 and Level-3 residual corrections fix remaining deviations between data and simulation. The Level-5 corrections would account for flavor-dependent variations in the detector response, but are not applied in the analyses of this thesis. The diagram is taken from [108].

be a constant function of the pseudorapidity, and the *L3Absolute* correction, which corrected for the transverse momentum dependence of the response.

In Run II these corrections are combined into a single step, although the historical nomenclature might suggest otherwise. In Figure 3.7(a) the simulated jet response for different pseudorapidity regions at $\sqrt{s} = 8$ TeV can be found.

The corrections are derived in a simulated QCD dijet sample by comparing the reconstructed jet p_T to the true values of the MC simulation.

Residual Corrections After the aforementioned corrections small deviations are still apparent between the MC simulation and data. The *Level-2 residual corrections* (L2Residuals) and *Level-3 residual corrections* (L3Residuals) are applied to data to account for these differences.

The L2Residuals are determined in recorded dijet events and correct for a pseudorapidity-dependence of the jets in recorded events. The corrections are derived by comparing a jet to a reference jet in the central barrel region with similar transverse momentum.

The L3Residuals correct the absolute scale of the jet energy. The corrections are determined in $Z(\rightarrow \mu\mu)+\text{jet}$, $Z(\rightarrow ee)+\text{jet}$, photon+jet and multijet events.

b Tagging of Jets

Quarks described in the standard model appear in six different flavors, as explained in Section 1.1. Whereas the lifetime of the top quark is too short to take part in the hadronization, bottom quarks pose an important handle to identify processes with heavy flavor content.

The decay of bottom quarks is suppressed, as they cannot decay inside of their generation due to the high top quark mass and hence have to decay into quarks of lower generations. This suppression leads to a longer lifetime of the B mesons of roughly $\tau = 1.6$ ps and therefore allows for the creation of displaced vertices with respect to the primary vertex, so-called secondary vertices. An illustration of the production of a secondary vertex can be seen in Figure 3.6. The analyses of this thesis rely heavily on the identification of jets originating from bottom quarks,

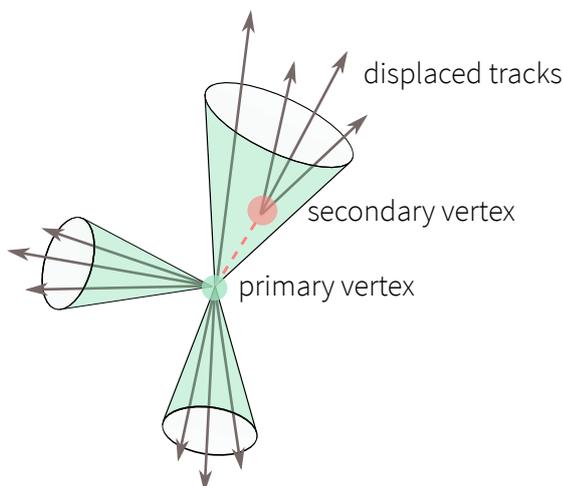


Figure 3.6.: The illustration shows the production of a secondary vertex due to the longer lifetime of a produced B meson. Multiple tracks inside of one jet are interpolated to a displaced vertex, the so called secondary vertex.

a method known as b tagging [111, 112].

Different b tagging algorithms are available in the CMS software framework and the *Combined Secondary Vertex Tagger* (CSV) is used in the analyses of this thesis. The CSV algorithm combines information about the impact parameter significance, the secondary vertex and jet kinematics to a single discriminant. In Run II an updated version of the CSV algorithm is used (CSVv2), which uses a neural network instead of a likelihood ratio as multivariate method to combine the information. The information on the secondary vertices is obtained with the Inclusive Vertex Finder algorithm [113].

Working points with fixed mistagging efficiencies, the probabilities to falsely identify a jet stemming from a light quark or a gluon as a b jet, are made available in CMS. The complete distribution of the CSVv2 discriminator for different processes can be found in Figure 3.7(b). Different tagging efficiencies in simulation and data can lead to large normalization differences therefore corrections have to be employed. In the Run I analysis described in Chapter 5 scale factors are applied which correct the (mis-)tagging efficiencies at given working points.

In the Run II analysis in Chapter 6 the complete distribution of the discriminator of the CSVv2 algorithm was reweighted, to achieve an agreement between simulated and measured events. This way the shape of the distribution can be utilized making the analysis independent of working points. A more thorough description on both corrections is given in the respective chapters.

3.3.7. Missing Transverse Energy

The bosons of the weak force produced in the LHC collisions can decay into the strictly weakly-interacting neutrinos. With an extremely low probability for interactions within the detector, neutrinos leave without trace.

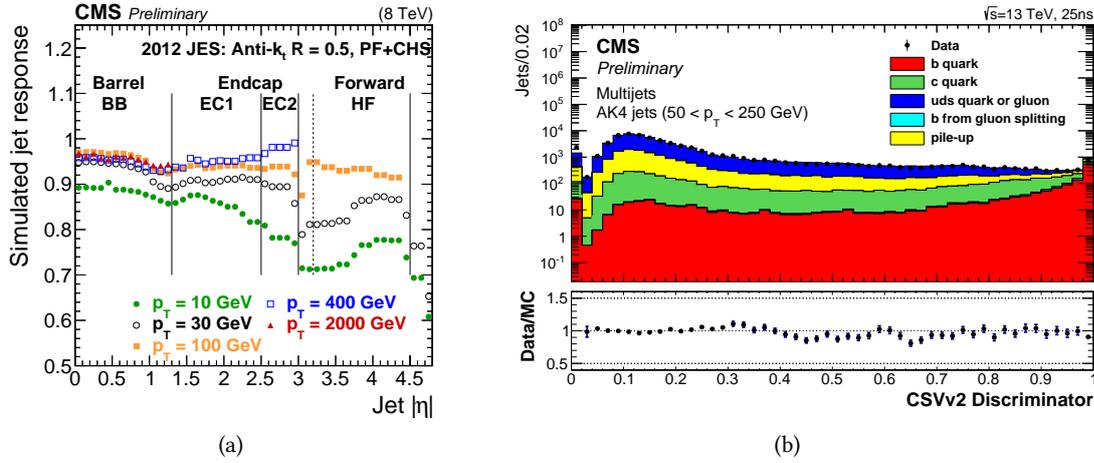


Figure 3.7.: In diagram (a) the simulated detector response in the different detector regions can be seen at $\sqrt{s} = 8$ TeV. Especially in the forward region, like the HF, and for low- p_T jets the detector shows a lower response, which is mostly corrected in the L2L3 corrections. In diagram (b) the distribution of the CSVv2 discriminator at $\sqrt{s} = 13$ TeV can be seen. The higher the values of the discriminator the higher the probability for a jet to be stemming from a b quark and the lower the mistagging efficiency. Tested with a multijet sample, the discriminator shows good agreement. The plots are taken from [109] and [114], respectively.

The circumstance that the partons of the hard subprocess carry in sum only longitudinal momentum can be exploited. Neutrinos render their share of the total transverse momentum undetectable, therefore causing an imbalance in the total transverse momentum.

The missing transverse energy (MET) could also be attributed to detector effects, e. g. particles leaving the detector undetected through the beam pipe, which have to be negated. Other analyses search for large amounts of missing transverse energy as a foot print of particles predicted by physics beyond the standard model, e. g. neutralinos predicted by super symmetry or possible dark matter particles, but these particles play no role in the scope in this analysis.

The *particle-flow missing transverse energy* is calculated as negative sum over the transverse momenta of all *particle-flow candidates*:

$$\vec{E}_T^{\text{raw}} = \sum_i (E_i \sin \theta_i \cos \phi_i \hat{x} + E_i \sin \theta_i \sin \phi_i \hat{y}) \quad ,$$

where \hat{x} and \hat{y} are the unit vectors in their respective axes. This quantity has to be corrected for aforementioned detector effects or distortion due to pileup pollution. An in-depth description of these corrections can be found in the References [115, 116]. Unlike other corrections, the MET corrections are not factors to apply to the raw MET but summands.

The *Type-I corrections* propagate the L2Relative and L3Absolute jet energy corrections to MET. The vector sum of all particles which could be clustered into jets is replaced by the vector sum of the particles exceeding a certain p_T threshold, but with applied jet energy corrections.

The *Type-II corrections* cover the consistent treatment the L1 jet energy corrections negating

3. Event Generation, Simulation & Reconstruction

pileup effects. Additionally, the Type-II corrections take the particles into account that fell below the threshold of the Type-I corrections and particles which were not clustered into jets at all. The norm of the corrected missing transverse energy $\vec{E}_T^{\text{type1p2}}$, abbreviated as E_T , is identified as the transverse momentum of the neutrino.

4. Multivariate & Statistical Methods

Simple counting experiment constituted the majority of analyses performed at particle colliders in the last decades, but they have been superseded by today's analyses employing multivariate techniques, which aid in bringing even faint signals in overwhelming background environments into prominence. By exploiting the correlation of a multitude of observables multivariate methods can classify objects much more effectively.

The analyses described in this thesis employ two different types of multivariate methods: Artificial Neural Networks (ANN) and Boosted Decision Trees (BDT). In the first part of this chapter the basics of these analysis methods are described. The second part will cover the tools of statistical inference utilized in the course of this thesis and the intricacies of the fitting procedure and the limit calculation are described in detail.

4.1. Multivariate Methodology

The classification of objects into two or more categories is one of the most abundant problems in science. By transferring the decision finding process to machines, a new level of optimization and reproducibility can be achieved.

Via plenty of different *machine learning* algorithms it is possible for the trained machine to predict the outcome of a certain target variable with the help of a set of input variables. The multivariate analysis tools employed in this thesis are provided by the TMVA software package [117] implemented within the ROOT framework [118].

4.1.1. Artificial Neural Networks

Artificial Neural Networks are a machine learning class inspired by biological neural networks as they are found in nature, most prominently in the human brain. A set of interconnected neurons produces a certain response for a given set of input variables, thereby mapping the n -dimensional space of n input variables x_1, x_2, \dots, x_n into a one-dimensional space. *Supervised learning* methods are used to train the neural networks of this thesis, which means that the networks are provided with a training data set, where the true target value is known and the network can learn the correlations of the input variables with each other and with the target. This enables the network to make a prediction for the target value in events where the true value is unknown.

The ANNs used in this thesis are classified as *multi-layer perceptrons*, as their neurons are organized in layers, such that only connections from neurons of one layer to the next are allowed. The ANNs consist of a first input layer with one node per input variable, a last output layer, which holds the final neural net estimator of the ANN, and an arbitrary number of hidden layers in between. In this thesis only ANNs with one hidden layer are employed and a schematic

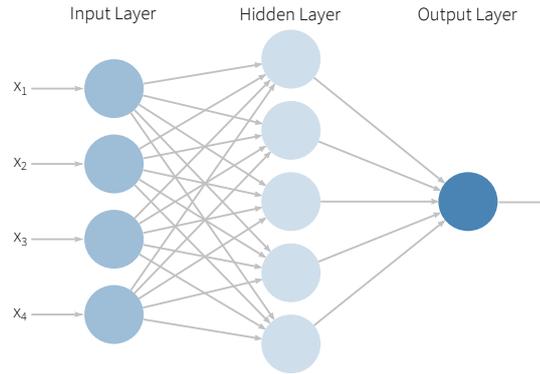


Figure 4.1.: Illustration of a neural network with four input variables and one hidden layer.

illustration of such a network can be found in Figure 4.1.

Each node is provided with a weighted n -dimensional input, which is mapped by a *synapse function* κ onto one dimension and subsequently converted by the *activation function* α to its output value. For this thesis a simple sum synapse function

$$\kappa : (y_1, \dots, y_n | w_0, \dots, w_n) \rightarrow w_0 + \sum_{i=1}^n w_i y_i$$

is chosen, which maps the n inputs y_i with their corresponding weights w_i onto a single value. As activation function the hyperbolic tangent is applied:

$$\alpha : x \rightarrow \frac{e^x - e^{-x}}{e^x + e^{-x}}.$$

The training of the neural networks is performed via *back propagation*, where the weights assigned to each interconnection are adjusted such, giving more significance to more separating variables and thereby optimizing the classification performance. By adjusting the weights of the connections an *error function* E , which is a measure for the rate in which the network classifies events correctly, given by

$$E(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{w}) = \sum_{a=1}^N E_a(\mathbf{x}_a | \mathbf{w}) = \sum_{a=1}^N \frac{1}{2} (y_{\text{ANN}, a} - \hat{y}_a)^2,$$

can be minimized, where N is number of events in the training, \mathbf{w} is the set of weights of the network, \mathbf{x}_a are the input values of event a , \hat{y}_a the true target value of event a and $y_{\text{ANN}, a}$ is the output of the neural network for event a . The algorithm used for the minimization of the error function in this thesis is the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm [119–122], an algorithm that uses the Hessian matrix of second derivatives for the weight adaption. A detailed description can be found in Reference [117].

4.1.2. Decision Trees

Naive Decision Trees (DT) [123] consist of consecutive binary questions based on a set of input variables. Depending on the previous answer a different decision has to be taken on the next

node. After a maximal number of nodes, a final verdict is reached depending on the majority of event class in the occupied endpoint, the *leaf*.

In the training of Decision Trees the best possible selection criterion for a node is defined by maximizing the *separation gain* S between two consecutive nodes. In order to determine the separation gain, the purity of each leaf has to be determined as

$$P = \frac{\sum_s w_s}{\sum_s w_s + \sum_b w_b} ,$$

where $\sum_{s,b} w_{s,b}$ are the summed weights of the signal and background events, respectively¹. The purity allows for the determination of the Gini index G , a measure of statistical dispersion, as

$$G = \left(\sum_{i=1}^n w_i \right) \cdot P \cdot (1 - P),$$

where w_i are the normalized weights for each event. The best possible question for the node can then be determined by maximizing S , given by

$$S = G_{\text{father node}} - G_{\text{child I}} - G_{\text{child II}}.$$

As the Gini index has a maximum at $P = 0.5$ of $G = 0.25$ (for unweighted events) the maximal separation gain would also be $S = 0.25$ for a maximally disperse starting ensemble ($G_{\text{father}} = 0.25$) and two perfectly separated child nodes ($G_{\text{child}} = 0$). Based on the purity after the training, leafs are either labeled as signal or background leaves.

By using a multitude of slightly altered decision trees and averaging over the predicted outcome the robustness of the procedure can be ensured, as a badly trained tree due to aberrations does not have enough power to sway the decision into the wrong direction. The average of all classifier outputs is later used as a final discrimination variable, labeled as "BDT output".

Adaptive Boosting

By using the *boosting* method the full potential of decision trees can be exploited. By subsequently training a multitude of trees, where in each iteration falsely categorized events get assigned a higher weight, the focus of the training is shifted to classify these previously misidentified events correctly.

In this thesis the *AdaBoost* [124] algorithm is utilized. Misclassified events are assigned a boost weight α_i for each tree i which is dependent on the misclassification rate Γ_i of tree i in the previous training iteration:

$$\alpha_i = \frac{1 - \Gamma_i}{\Gamma_i}.$$

Correctly classified events get reweighted such that the sum of weights remains constant. The BDT output of a boosted ensemble of N decision trees is determined as

$$y_{\text{Boost}}(\mathbf{x}) = \frac{1}{N} \cdot \sum_i^N \ln \alpha_i \cdot h_i(\mathbf{x}),$$

¹In an unboosted DT and in the first iteration of the training of a BDT the weights for each event are 1, so $\sum_{s,b} w_{s,b}$ corresponds to the total number of signal and background events, respectively.

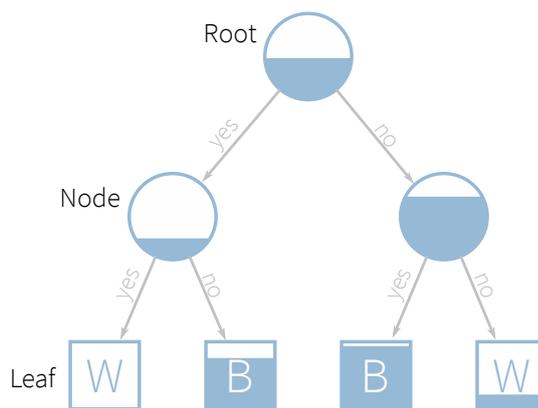


Figure 4.2.: Illustration of a decision tree trying to separate an initial sample which contains events that shall be classified in one of two categories (white/blue) and is maximally disperse at the root. At each node a binary question is posed which is answered based on the properties of the specific event. At a given depth the purity of the contained events determines the final leaf as leaf of either of the two categories.

where $h(x)$ is the binary outcome of tree i , mathematically realized as $h(x) = 1$ for events ending on a signal leaf and $h(x) = -1$ for events ending on a background leaf.

4.1.3. Overtraining

With suboptimal parameter settings multivariate methods are prone to learning statistical fluctuation of the training data set. This effect, known as *overtraining*, seems to improve the classification power of the method during the training, but evaluating it on an independent test sample shows a diminished potency, as the learned fluctuations are not present anymore. An illustration of the issue of overtraining can be seen in Figure 4.3. In order to test for overtraining the training sample is split into two parts, a separate training and a testing sample. The MVA is trained on the training sample and subsequently checked for a similar performance on the testing sample. Both MVA output distributions should match, which is tested with a Kolmogorov-Smirnov test (KS test), which provides the probability that two distributions originate from the same mother distribution.

In order to avoid a bias in the subsequent analysis training and testing sample are discarded, as they have both been used to optimize the MVA method. The analysis uses a third sample, the evaluation sample, with events that have never been in contact with the MVA method.

4.1.4. Variable Ranking

When optimizing MVA methods it is interesting to see which variables are the most important in the training of the classification. For neural networks this variable importance I for a certain variable i is determined by

$$I_i = \bar{x}_i^2 \sum_{j=1}^{n_h} \left(w_{ij}^{(1)} \right)^2, \quad i = 1, \dots, n_{\text{var}},$$

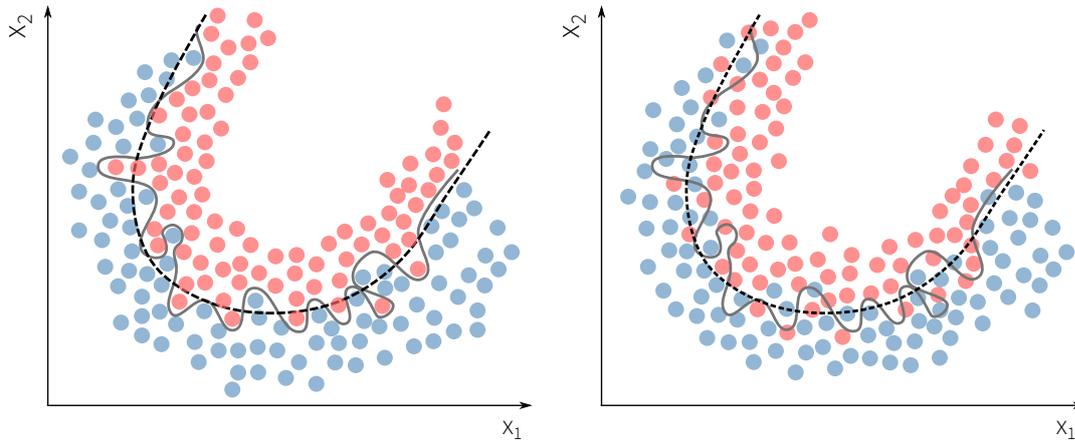


Figure 4.3.: Illustration of the effect of overtraining on an independent testing sample. The distributions show the correlations of variable x_1 and x_2 for two classes of events (red and blue) for a training sample (left) and a testing sample (right). An extremely overtrained method would result in a separation line that perfectly separates the two classes (solid, gray). The application of this separation line to the testing set shows a much worsened performance (right). A robust, non-overtrained MVA would pick a separation line (dashed, black) that performs equally well for training and testing sample.

where w_{ij} are the weights of the connections from input node i to the nodes of the hidden layer, n_h is the number of hidden nodes in the hidden layer and \bar{x}_i^2 is the arithmetic mean of the variable.

The ranking of BDT input variables is determined by their number of occurrences in the tree nodes. This count is weighted with its squared separation gain and the number of events in the node.

4.2. Statistical Inference

In order to lend substance to performed measurements, a coherent statistical inference is obligatory. This chapter shortly discusses the applied methods during this thesis. The information for this section is obtained from References [125–128], if not noted otherwise.

4.2.1. Maximum Likelihood Parameter Estimation

The determination of parameters that are in best agreement with a set of measurements, poses an often faced problem in science. For a given, fully known probability density $f(x|a)$, where a is a set of unknown parameters of the density, \hat{a} is the best estimation for the reproduction of a set of measurements x_1, x_2, \dots, x_n .

The basis for the maximum-likelihood estimation (MLE) is the likelihood function, given by

$$L(a) = \prod_{i=1}^n f(x_i|a).$$

The set of parameters a which minimize the likelihood function are found to be the best estimators. Possible difficulties with extremely low likelihoods for high number of measurements can be avoided by utilizing the monotony of the logarithm. The negative log-likelihood function is given by

$$F(a) = -\ln L(a) = -\sum_{i=1}^n \ln f(x_i|a),$$

which converts the former product into a sum, mediating the problem of tiny likelihood values. The negative sign is a convention chosen for historical reasons. By minimizing $F(a)$ the set of best estimators can be obtained. The likelihood for binned distributions that are the core of the analyses of this thesis, is given by the product of Poissonian probabilities for each bin. The likelihood is given by

$$L(\text{data}|\mu, \theta) = \prod_i \frac{(\mu s_i(\theta) + b_i(\theta))^{n_i}}{n_i!} e^{-\mu s_i(\theta) - b_i(\theta)},$$

where s_i and b_i are the expected signal and background yields for bin i , respectively, n_i is the observed number of events in bin i and signal strength modifier $\mu = \sigma/\sigma_{\text{SM}}$. The expected yields s_i and b_i are dependent on a set of nuisance parameters θ which is used to incorporate the analysis uncertainties into the likelihood function.

4.2.2. Nuisance Parameter Treatment

Each source of uncertainty in the course of an analysis introduces an additional nuisance parameter θ_i . Two types of nuisance parameters are considered in this thesis: rate uncertainties and shape uncertainties. Rate uncertainties are bin-independent, but can be dependent on process p , and scale the affected template by a constant factor. Shape uncertainties are bin-dependent and do not only alter the normalization, but can also change the complete shape of a template. The nuisance parameters can be incorporated into the likelihood with an additional Gaussian distribution. An uncertainty on the mean number of background events b can be modeled with

$$\pi_\beta(\beta) = \frac{1}{\sqrt{2\pi}\sigma_\beta} \cdot \exp\left[-\frac{(\beta - \beta_0)^2}{2\sigma_\beta^2}\right],$$

for $\beta = \ln b$, its mean value β_0 and its standard deviation σ_β . The relative uncertainty is identified as $\sigma_{\text{rel}} = e^{\sigma_\beta} - 1$. Each considered nuisance parameter adds an additional Gaussian² function, which is multiplied with the likelihood.

Due to their bin-dependence the treatment of shape uncertainties is more complex. An often posed problem is the existence of a nominal template and two varied templates, which correspond to an up- and downward shift of an uncertainty by one standard deviation. To be able to access also templates at other but these three discrete values, a *template morphing* method is utilized. More details can be found in Reference [129].

²Statistical uncertainties which arise due to the use of Monte Carlo simulation samples with limited event numbers are incorporated with the use of Poissonian distributions.

4.2.3. Exclusion Limits

The discovery or exclusion of new physical processes from the statistical point of view boils down to the testing of hypotheses. Two hypotheses, the null hypothesis H_0 and the alternative hypothesis H_1 , are defined. For the exclusion of a certain physics model, as it is performed in this analysis, H_0 states the existence of a specific signal ($\mu = 1$) and H_1 states that only the background is present ($\mu = 0$). The existence of this process can hence be excluded by rejecting H_0 in favor of H_1 . As a quantitative measure for the likeliness of a hypothesis a *test statistic* q is introduced. The *Neyman-Pearson lemma* [130] states that the likelihood-ratio test is the most powerful test and for LHC purposes the profile likelihood ratio

$$q_\mu = -2 \ln \frac{L(\mu, \hat{\theta})}{L(\hat{\mu}, \hat{\theta})}$$

is used. Here $\hat{\theta}$ is the set of nuisance parameter values that maximizes L for the given signal strength modifier μ . The parameters $\hat{\mu}$ and $\hat{\theta}$ in the denominator globally maximize the likelihood.

For given signal and background models the probability density function (p.d.f.) $f(q_\mu | \mu, \theta)$ of the test statistic can be sampled with the use of pseudo experiment for the range of μ of interest. Two exemplary p.d.f. of test statistics for the exclusion of a signal hypothesis can be found in Figure 4.4.

The outcome of a complete analysis can be condensed into a single observed value of q , which can be used to calculate the *p-values* of the observation:

$$p_\mu = \int_{q_{\mu, \text{obs}}}^{\infty} f(\tilde{q}_\mu | \mu) dq_\mu$$

$$1 - p_b = \int_{q_{\mu, \text{obs}}}^{\infty} f(\tilde{q}_\mu | 0) dq_{\tilde{\mu}}$$

The higher the value of $q_{\mu, \text{obs}}$ the more background-like is the observation. The p_μ -value states the probability that a value of $q_{\mu, \text{obs}}$ or higher would be obtained, if a signal with a signal strength modifier μ is present. So if $p_\mu < 5\%$, the signal hypothesis (for a specific μ) could in principle be rejected at 95% Confidence Level (C.L.). However, for small expected signals the two p.d.f. of the test statistics q_μ and q_0 almost overlap completely making it possible that downwards fluctuations of the background would lead to a rejection of the background hypothesis or that a signal hypothesis is rejected, although the experiment is not sensitive to it. To protect against these dangers, a modified *p-value* is calculated as

$$p'_\mu = \frac{p_\mu}{1 - p_b} \equiv CL_s,$$

which is also known as CL_s limit [131]. By scanning over μ and stating the μ^{up} for which $p'_{\mu^{\text{up}}} = 5\%$ an upper limit of μ^{up} can be set. If $\mu^{\text{up}} < 1$, the signal model can be excluded at 95% C.L.

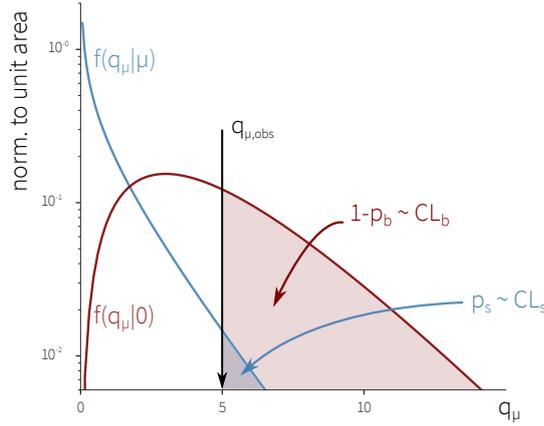


Figure 4.4.: Illustration of p.d.f.s of two test statistics for the signal+background hypothesis $f(q_\mu|\mu)$ (blue) and the background-only hypothesis $f(q_\mu|0)$ (red). For an observed value of the profile likelihood $q_{\mu,\text{obs}}$ the respective p -values can be calculated by integrating the area under the p.d.f.s from $q_{\mu,\text{obs}}$ to infinity. The ratio of the two areas serves as the signal confidence level CL_S .

In order to find deviations from the expectation it is common to also quote *expected upper limits*. Without the need for data one can perform thousands of pseudo experiments based on the background-only hypothesis, where for each pseudo experiment one value of $\mu_{\text{toy}}^{\text{up}}$ is retrieved. The median of the distribution of these $\mu_{\text{toy}}^{\text{up}}$ values serves as the expected upper limit. The outer boundaries of the distribution that include 68% and 95% of the area centered around the median, serve as one and two standard deviation values, respectively, and are often visualized by green and yellow uncertainty bands around the expectation.

4.2.4. Asymptotic Limits

The construction of the test statistics and the calculation of the expected upper limits with its uncertainty bands is very computing expensive. The discovery of the asymptotic formulae for likelihood-based tests [132] allows for an extremely fast approximation of the limit setting procedure described above. By employing the theorems of Wilks [133] and Wald [134] the p.d.f.s of the test statistic can be approximately calculated instead of the cumbersome construction with pseudo experiments.

Another aspect of the limit calculation with the asymptotic method is the usage of artificial *Asimov data set*. This data set serves as a representative data set for the $\mu_{\text{toy}}^{\text{up}}$ values, which can be used to calculate expected limit and the respective standard deviations in only a single iteration. The Asimov data set is defined such that the true nuisance parameter values are recovered when it is used to evaluate the best estimations.

The asymptotic method is found to yield slightly lower upper limits, an effect that is enhanced in regions with a low number of events.

In this thesis the statistical inference is largely performed with the COMBINE package [135] provided by CMS, which itself employs the ROOFIT toolkit [136].

5. Search for Associated Production of Single Top Quark and Higgs Boson at $\sqrt{s} = 8$ TeV

In this chapter a complete presentation of the search for a Higgs boson produced in association with a single top quark (tHq) at a center-of-mass energy of $\sqrt{s} = 8$ TeV is given. With a branching ratio of 58.1% Higgs bosons decay into a pair of bottom- and anti-bottom quark. A search in this decay channel allows for the exploitation of the most produced Higgs bosons, but puts the analysis in a very demanding high-jet-multiplicity environment. By requiring a leptonic decay of the top quark, only about 19% of the tHq events are available to study, but the lepton is essential as a handle to suppress the QCD multi-jet background drastically and it poses as a valuable trigger object.

The analysis of tHq events allows for a probe of the Yukawa coupling of the Higgs boson to the top quark. A flipped sign of the coupling with respect to its expected value would increase the cross section of the process significantly. More information about the tHq process, its attributes and the underlying theoretical model is provided in Section 1.3. The goal of this analysis is the exclusion of the presence of the tHq process under the assumption of an anomalous coupling of $C_t = -1$.

The first part of this chapter gives an overview about the complete analysis strategy, the characteristics of the signal and background processes and the used physics object definitions. Following is the description of the event selection, as well as details on the two employed event reconstructions. In the last part of this chapter the final classification of events, the thorough statistical analysis of the process and the final exclusion limits are presented.

This is the first time a search for tHq, $H \rightarrow b\bar{b}$ has been performed. The complete analysis was a collaborative effort within CMS with a strong influence from an analysis group at KIT including myself. For the course of this thesis the analysis has been re-performed step by step. The analysis has been made public by CMS as a Physics Analysis Summary [137] and was part of a combination effort of multiple tHq analyses studying different Higgs boson decay channels. A paper on the combination of all different tHq analyses has been accepted for publication in the Journal of High Energy Physics (JHEP) [138] and will be published soon. The results of the combination are also shown in the last section of this chapter.

5.1. Analysis Strategy

The Standard Model of particle physics predicts a very small cross section for the tHq production of $\sigma_{\text{SM}} = 18.3$ fb at $\sqrt{s} = 8$ TeV. Early projection studies showed that a search for this process would require much more recorded data to become sensitive to this process [139]. The analysis is hence optimized for a signal process under the assumption of a negative Yukawa coupling of

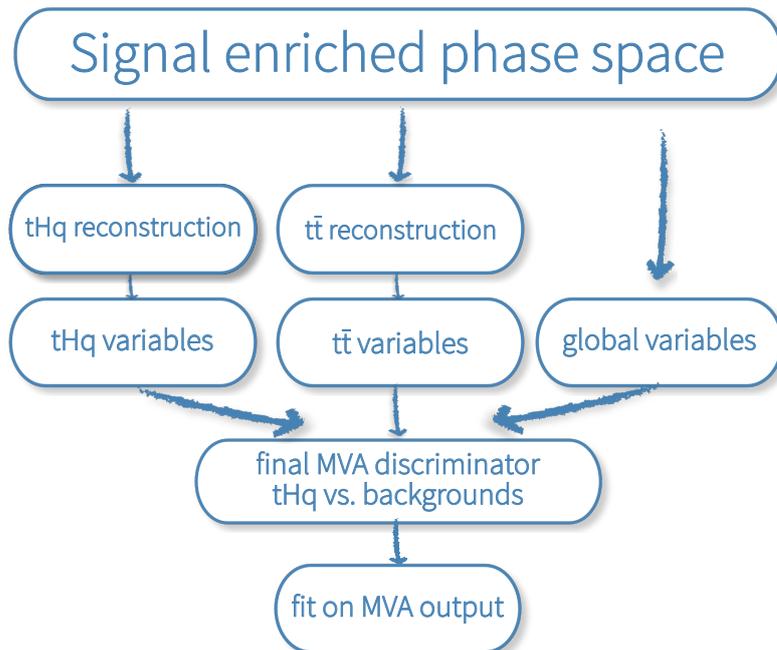


Figure 5.1.: A schematic overview of the analysis workflow. A set of optimized selection cuts ensures a favorable signal to background ratio and puts the analysis in a signal enriched phase space. The final classification uses variables out of three categories: variables based on tHq reconstruction, where jets in the event are assigned to quarks of the tHq final state, variables based on the $t\bar{t}$ reconstruction, where jets are assigned to the quarks in the final state of a semi-leptonic $t\bar{t}$ event and variables that are based on neither. These variables are used in a multivariate method which classifies events to be either more signal- or background-like. The final limit is then derived by performing a likelihood fit in the neural network output distribution and calculating CL_S limits.

the Higgs boson to the top quark $C_t = -1$, which would increase the cross section to $\sigma_{C_t=-1} = 234.8 \text{ fb}$ due to interference effects.

An overview of the analysis strategy can be seen in Figure 5.1. Firstly, an optimized set of selection criteria is applied, thereby enhancing the signal over background ratio in the defined signal regions.

The main differences between the signal process and some background processes can be traced back to the behavior of the different particles in the final states. Combined with the high-jet-multiplicity environment, a well-performing and reliable assignment of jets to the individual partons becomes essential. Two jet assignment algorithms employing the same multivariate methods but different event content hypotheses are used in this analysis. Based on the two obtained jet assignments a set of discriminating variables can be extracted. These variables are used in a multivariate event classification. As a final step in this analysis a fit of the simulation to the actual data in the classification discriminator distribution is performed. After a complex statistical investigation a statement on the nature of the anomalous tHq process can be made.

5.2. The tHq Process

The tHq process has several distinct features that are essential in order to separate this process from background processes.

The naming schema is derived from the final partons that are available directly after the production. Like in the t -channel single top production, a top quark is produced singly accompanied by a light quark, denoted simply as q . This light quark is produced predominantly in the forward region of the detector. The main interest of this analysis however, lies on the produced Higgs boson. It can be either emitted by the exchanged W boson or directly by the singly produced top quark.

Due to their high masses and thus low lifetimes, the top quark and Higgs boson decay directly after their production in the detector. Exploiting its high branching ratio, the Higgs boson is required to decay into a $b\bar{b}$ quark pair. Due to the $|V_{tb}|$ element of the CKM matrix being very close to one¹, the top quark almost always decays into a bottom quark and a W boson. In this analysis the W boson is required to decay leptonically. Tau leptons are not reconstructed separately, only their leptonic decays into either electron or muon are considered in this analysis. Furthermore, the b quark in the initial state calls for a special discussion. As bottom quarks are too heavy to be part of the colliding protons, the bottom quark has to emanate from a gluon splitting. The corresponding anti-quark also leaves a trace in the detector, however its momentum is regularly too low and hence falls outside the detector acceptance.

In summary, the signal process of this thesis, the tHq process, is characterized by four b quarks, one muon or electron, one neutrino and one light flavored quark produced in forward direction. Other decay chains, such as a hadronically decaying top quark or the decay of the W boson into τ leptons, are not considered in this analysis, as this would further complicate the analysis with an even higher jet multiplicity, without a good trigger object, or additional neutrinos in the final state.

The Feynman diagram of the full process can be seen in Figure 5.2. Other production mechanisms, like tHW (tW -channel production) or tHb (s -channel production), are not considered in this analysis.

5.3. Background Processes

The low cross section of the signal process calls for a thorough understanding of all involved background processes. The main background of this analysis is the top quark pair production, which is subjected to specialized treatments to separate $t\bar{t}$ from tHq production. All considered backgrounds are described in the following sections and the main backgrounds are depicted in Figure 5.3.

Top Quark Pair Production

By far the dominating background of this analysis is the production of a top quark pair. The pair is produced via fusion of two gluons or the annihilation of a quark-antiquark pair.

¹Recent measurements of the matrix element are still compatible with $|V_{tb}| = 1$ within the uncertainties [140, 141].

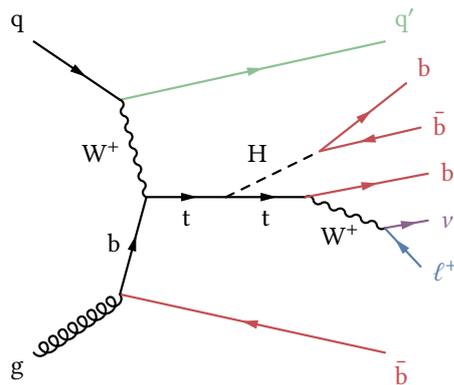


Figure 5.2.: Relevant Feynman diagram of tHq production mechanisms. The diagram shows the case, where the Higgs boson is emitted by the single top quark in the event. The Feynman diagram where the Higgs boson couples to the W boson (see Section 1.3) is omitted, as the final states are identical. The indistinguishable final state for both production mechanism is characterized by four bottom quarks (red), one charged lepton (blue), one neutrino (violet) and one light quark (green).

Based on the decay of the individual W bosons the $t\bar{t}$ process can be sorted into three categories: the full-hadronic decay mode, where both W bosons decay into a quark-antiquark pair, the full-leptonic decay mode, where both W bosons decay into a charged lepton and its corresponding neutrino and the semi-leptonic decay mode, where one W boson decays hadronically and the other leptonically.

As the signal process is expected to produce exactly one charged lepton, the semi-leptonic $t\bar{t}$ decay is the decay channel dominating the background composition, but the full-leptonic decay mode also contributes due to undetected leptons. The full-hadronic decay only plays a negligible role.

At Born-level the semi-leptonic $t\bar{t}$ production includes two b quarks, one charged lepton and one neutrino in the final state. The emission of gluons and their subsequent decay into a pair of quarks leads to additional jets in the detector. Based on the original quark flavors of the additional jets, the $t\bar{t}$ production is separated into different categories. Especially the emission of gluons which subsequently decay into a pair of bottom quarks is a process ($t\bar{t}+b\bar{b}$) theoretically not understood in its entirety, but poses as important process to this analysis, making a tailored treatment for this process necessary. Also jets from light quarks ($t\bar{t}+\text{light}$) mistakenly identified as b jets can mimic the signal process. The analysis is consequently also susceptible to the $t\bar{t}$ production with additional jets emerging from charm quarks ($t\bar{t}+c\bar{c}$), due to their increased mistagging probability.

The decay products of the $t\bar{t}$ production are expected to appear more centrally in the detector, thus making the light forward jet of the tHq process an important property to separate the signal from the $t\bar{t}$ process.

A Feynman diagram of the $t\bar{t}$ production can be found in Figure 5.3(a).

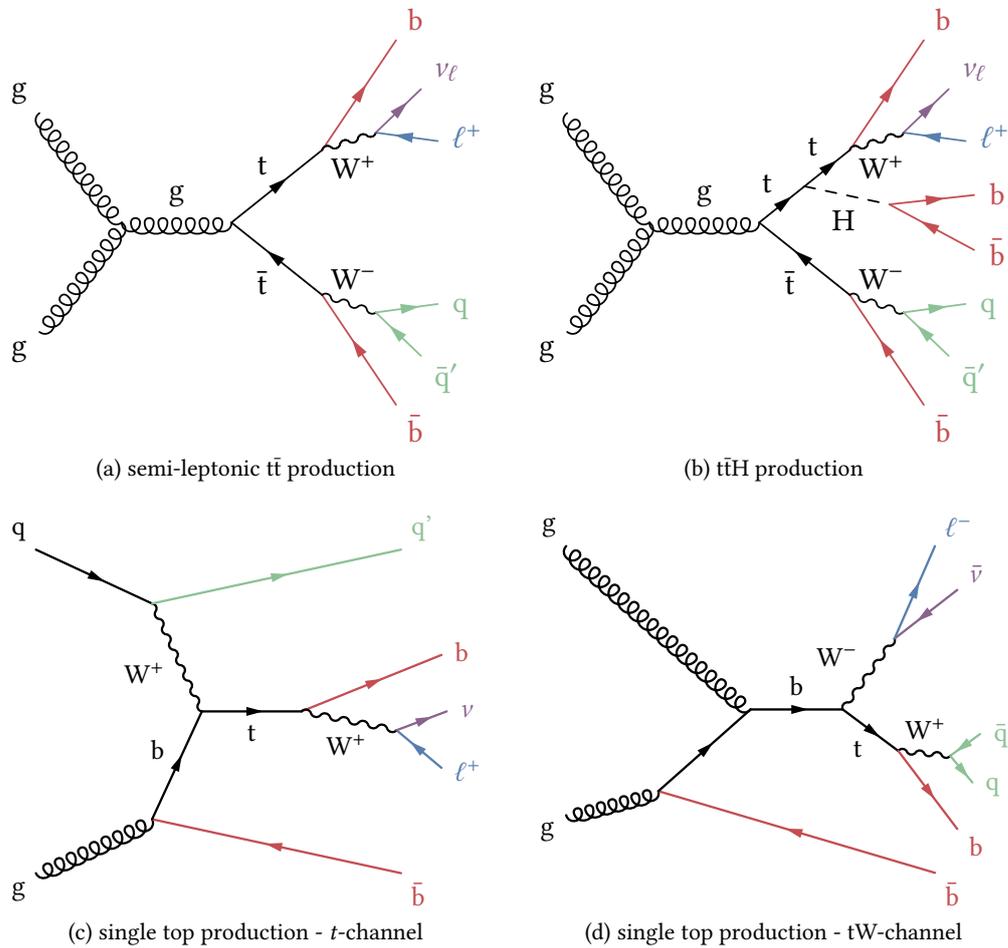


Figure 5.3.: Feynman diagrams of various background processes. In (a) the most prominent semi-leptonic $t\bar{t}$ process can be seen. Its final state consists of two b quarks (red), one charged lepton (blue), one neutrino (violet) and two light jets (green).

The top quarks in the $t\bar{t}$ event can emit a Higgs boson, which can decay into a pair of b quarks. This $t\bar{t}H$ ($H \rightarrow b\bar{b}$) background can be seen in (b).

The mother process of signal tHq production is the single top quark production. Analogous to the signal process, the single top quark can be produced in three production channels: the t -channel (c), the tW -channel (d) and the s -channel (not depicted). The single top diagrams are shown in the four-flavor scheme, where b quarks are not part of the proton and an initial gluon splitting into $b\bar{b}$ is required leading to an additional b quark in the final state.

Single Top Production

The production of single top quarks also is hard to separate from the signal process. The t -channel production process, basically the mother process of the tHq process missing the associated Higgs boson, offers similar characteristics as the signal process. Especially, the emission of a gluon decaying into a $b\bar{b}$ pair leads to a final state with exactly the same particles than the signal process. Due to its significantly lower cross section compared to $t\bar{t}$ production the single top production is not subject to any specialized treatment but still appears in phase spaces where the signal is expected.

The tW -channel and s -channel production mechanism also contribute to the background, but to a much lower degree. The corresponding Feynman diagrams for t -channel and tW -channel production can be found in Figures 5.3(c) and 5.3(d).

$t\bar{t}H$ Production

The associated production of a Higgs boson with a top quark pair also has to be considered as background process. Although CMS and ATLAS have not yet measured the $t\bar{t}H$ production with over 5σ significance and possible deviations from the SM are still within the realms of possibility [142, 143], this process is taken into account as predicted by the SM. With a semi-leptonically decaying $t\bar{t}$ pair and a Higgs boson decaying into $b\bar{b}$ this process would also produce four b quarks, one charged lepton and one neutrino. The presence of a Higgs boson also exacerbates the separation from the signal process, as variables such as the invariant mass of the two bottom quarks from the Higgs boson decay show the same behavior.

The additionally produced top quark, which predominantly decays hadronically, leads to a higher expected jet multiplicity in the $t\bar{t}H$ final state. Also the two top quarks and the Higgs boson are assumed to be more central than the particles of the tHq production, therefore making the light forward quark the most valuable asset of this analysis for the discrimination against $t\bar{t}H$ production. The corresponding Feynman diagram to this process can be found in Figure 5.3(b).

Minor Backgrounds

Three different processes are considered in the analysis but due to their low impact on the analysis they are grouped under the label "Misc".

Diboson Production The pair production of W and Z bosons pose a minor background to this analysis. A small cross section of WW , WZ and ZZ pair production combined with low branching ratios of final states that would fake a tHq signature cause the impact of this process to be insignificant.

W/Z+jets Production Another considered background is the production of a single W or Z boson produced in association with jets. Although unlikely to produce many b -tagged jets, its comparably huge production cross section makes W +jets a background that is able to produce a small number of events imitating the tHq signature. The Z +jets production contributes negligibly to the background.

QCD Production

The by far most frequent processes at the LHC are multijet events via the strong interaction summarized as QCD production. Able to fake almost every signature in high jet environments, the QCD production is a background that often calls for a special treatment. Leptons can be produced during the decay of B- or D hadrons, but are often found to be non-isolated. The introduction of a missing transverse energy selection criterion is able to reduce this background to a negligible level. This has been demonstrated with the ABCD method [144] and a QCD contribution of below 1% has been found, therefore justifying its negligence.

5.4. Datasets

This analysis is performed with the complete dataset recorded by the CMS detector at $\sqrt{s} = 8$ TeV in 2012. As the studied final state contains exactly one charged lepton the `SingleElectron` and `SingleMuon` datasets are used. After the 22Jan2013 reprocessing and the selection of luminosity sections deemed usable based on the golden JSON file [145], which is provided by the CMS collaboration, this resulted in an integrated luminosity of 19.7 fb^{-1} . The standard single lepton trigger paths `HLT_E1e27_WP80_v*` and `HLT_IsoMu24_eta2p1_v*` have been utilized in this analysis.

The Monte Carlo simulation samples used in this analysis have been produced during the `Summer12_DR53X` production campaign. The simulation sample of the associated single top quark and Higgs boson production have been generated employing the `MADGRAPH 5.1` event generator. It has been generated in the four-flavor scheme and only the t -channel process is considered. The Higgs boson is forced to decay into a pair of $b\bar{b}$ and the top quark has to decay leptonically $t \rightarrow b\ell\nu$, where ℓ is either an electron, a muon or a tau, which are subsequently brought to decay with the `TAUOLA` package and are required to decay leptonically into either electron or muon and the respective neutrinos. A signal sample generated in the four-flavor scheme has shown to reproduce the dynamics better than a sample produced in the five-flavor scheme, however the cross section has been calculated employing the five-flavor scheme. Hence, the four-flavor scheme sample is scaled to the cross section calculated with the five-flavor scheme. The cross section is subsequently reduced after considering the $H \rightarrow b\bar{b}$ and the $t \rightarrow b\ell\nu$ branching ratios as well as the selection efficiency of the loose selection criteria applied in the generation step.

The single top samples have been simulated with the `POWHEG` package. Each of the production channels is produced separately and they are subsequently scaled to their approximate NNLO cross section. Analogous to the signal sample the top quark is forced to decay leptonically.

The `PYTHIA 6.4` package was used for the generation of diboson and QCD samples. The `WW` sample was scaled to its NLO cross section prediction and the `WZ`, `ZZ` and `QCD` samples are normalized to their leading order cross sections.

The $t\bar{t}$ +jets and the V +jets samples have also been generated with `MADGRAPH 5.1` and their templates are normalized to NNLO cross sections.

The parton shower for all of the samples mentioned above is provided by `PYTHIA 6.4`. A complete table with all used data and simulation samples can be found in Appendix A.1 and A.2.

5.4.1. Heavy Flavor Splitting

As already mentioned in Section 5.3 the simulated $t\bar{t}$ samples are subjected to a special treatment considering the production of additional heavy flavor jets. Both, CMS and ATLAS, have measured different ratios of $t\bar{t}+b\bar{b}$ to $t\bar{t}$ +light than what is seen in the samples generated by MADGRAPH and PYTHIA 6.4 [146,147]. To account for a mismodeling of the ratios between $t\bar{t}$ with additional heavy flavor jets to $t\bar{t}$ with additional light jets, the samples have been separated based on the parton content of the jets. Every event in the $t\bar{t}$ simulation sample was sorted into one of nine categories, which have been subsequently condensed into only four categories as they only show negligible differences in their distribution shapes. The remaining categories are:

$t\bar{t}+b\bar{b}$ The event contains at least two additional jets which could be matched to a b quark each.

$t\bar{t}+b$ The event contains one additional jet which could be matched to one or more b quarks. This might be due to a collinear gluon splitting, which resulted in one jet containing both b quarks.

$t\bar{t}+c\bar{c}$ The event contains at least one additional jet which contains one or more c quarks.

$t\bar{t}$ +light All other events not satisfying the criteria mentioned above.

This procedure has been adopted from earlier CMS analyses described in Reference [148, 149].

5.5. Physics Objects

This analysis uses objects reconstructed with the particle-flow algorithm explained in Section 3.3. For the further analysis the CMS Software (CMSSW) version 5.3.11 is used. More details on the used physics objects can be found in Reference [150].

5.5.1. Primary Vertices

The vertex with largest p_T -sum of the collection of all primary vertices obtained by the particle-flow algorithm has to satisfy a set of quality requirements. It has to be located inside a cylinder around the beam axis with a radius of 2 cm and it must not be farther away than 24 cm from the origin in its z -coordinate. The primary vertex must at least have four degrees of freedom and must not be tagged as a *fake* vertex. After the reconstruction of all vertices the effect of pileup can be mitigated by removing every charged hadron particle-flow candidate which cannot be associated to the primary vertex. This procedure is known as charged hadron subtraction (CHS) [151].

5.5.2. Muons

Muons are reconstructed as recommended by the Muon Physics Objects Group (POG) for Run I applying the MuonID [152]. The considered particle-flow muon candidates have been global

muons with a reduced chi-squared statistic ($\chi^2/\text{n.d.o.f.}$) less than 10 each. The reconstructed muon must have initiated at least one hit in the muon chambers and muon segments in at least two muon stations. Additionally, every reconstructed muon must satisfy for the transverse impact parameter $|d_{xy}| < 2$ mm and for the longitudinal distance $|d_z| < 5$ mm. The reconstructed muon must have left at least one hit in the pixel silicon detector and at least hits in five different tracker layers.

All reconstructed muons satisfying the requirements above are then subjected to a set of kinematic cuts: the muons must satisfy $p_T > 26$ GeV, $|\eta| < 2.1$ and the isolation requirement $I_{\Delta\beta} < 0.12$, where the relative $\Delta\beta$ -corrected isolation is defined as

$$I_{\Delta\beta} = \frac{1}{p_{T,\ell}} \left(I_{\text{CH}}^\ell + \max \left(I_{\text{NH}}^\ell + I_{\text{Ph}}^\ell - 0.5 \cdot I_{\text{CH,PU}}^\ell, 0 \right) \right) . \quad (5.1)$$

Here, $I_{\text{CH}}^\ell, I_{\text{NH}}^\ell$ and I_{Ph}^ℓ are defined as the energy deposited in a cone with $\Delta R = 0.4$ around the muon track by charged hadrons, neutral hadrons and photons, respectively. $I_{\text{CH,PU}}^\ell$ is the energy deposited by charged hadron pileup candidates in the same cone.

5.5.3. Electrons

Two different requirement sets are used to reconstruct a *loose* and a *tight* electron collection. Loose electrons must satisfy $p_T > 20$ GeV, $|\eta| < 2.5$ and $I_\rho < 0.15$. The isolation I_ρ is defined as

$$I_\rho = \frac{1}{p_{T,\ell}} \left(I_{\text{CH}}^\ell + \max \left(I_{\text{NH}}^\ell + I_{\text{Ph}}^\ell - \rho A_{\text{eff}}, 0 \right) \right) , \quad (5.2)$$

where $I_{\text{CH}}^\ell, I_{\text{NH}}^\ell$ and I_{Ph}^ℓ are the deposited energies in a cone with $\Delta R = 0.3$ around the electron by charged hadrons, neutral hadrons and photons, respectively. The average angular p_T density is denoted as ρ and A_{eff} is the effective area compensating for the neutral contribution of pileup in the cone.

The tight electrons must satisfy tighter kinematic cuts with $p_T > 30$ GeV, $|\eta| < 2.5$ and $I_\rho < 0.1$. Additionally, stricter requirements on the reconstruction quality are imposed: A veto on a possible photon conversion is applied. The electrons must receive a response of the triggering MVA ID [153] larger than 0.9 and pass the trigger-emulating preselection. Also a veto on electron clusters in the gap of the ECAL ($1.4442 < |\eta| < 1.5660$) is applied.

The transverse momentum of the particle-flow algorithm is substituted by the GSF momentum of the electron.

5.5.4. Jets

This analysis uses jets clustered with the anti- k_r algorithm with a cone size of 0.5. Before the clustering charged hadrons identified as originating from pileup events as well as muons and electrons with looser isolation criteria, $I_{\Delta\beta} < 0.2$ and $I_\rho < 0.15$, respectively, are removed from the particle collection used for the clustering.

The clustered jets must pass the loose working point of the JetID algorithm [154].

Due to discrepancies seen in the transverse momentum spectrum of jets in the forward region between recorded and simulated collisions, pseudo-rapidity-dependent selection requirements are imposed. Jets in the central region with $|\eta| < 2.4$ have to satisfy $p_T > 20$ GeV and for jets in the forward region satisfying $2.4 < |\eta| < 4.7$ a selection criterion of $p_T > 40$ GeV is chosen. The aforementioned discrepancies are described in detail in Section 5.6.5.

Jets in simulation are corrected with the appropriate L1L2L3 MCtruth corrections, which are described in Section 3.3.6, and the residual corrections are applied to jets in data [155]. Also the resolution in simulation is smeared to match the resolution observed in data.

5.5.5. Missing Transverse Energy

The particle-flow missing transverse energy is used and Type-0 and Type-1 corrections (explained in Section 3.3.7) are applied. Additionally, the x - y -shift correction [115] is applied, which reduces the modulation seen in the azimuthal angle of the missing transverse energy distribution. Depending on the number of reconstructed primary vertices in the event the coordinate system is shifted in the x - y -plane, thereby reducing this unphysical modulation.

5.5.6. W Boson Reconstruction

The analysis of this chapter relies on a full reconstruction of the final states of the tHq process. The leptonically decaying W boson can be reconstructed by adding the four-vectors of the lepton and the neutrino. However, the only measurable quantity is the transverse part of the missing energy, rendering the z -component of the neutrino momentum $p_{z,\nu}$ inaccessible. By setting up the equation for the invariant mass of the W boson a quadratic equation for $p_{z,\nu}$ can be derived:

$$m_W^2 = \left(E_\ell + \sqrt{\cancel{E}_T^2 + p_{z,\nu}^2} \right)^2 - \left(\vec{p}_{T,\ell} + \vec{\cancel{E}}_T \right)^2 - (p_{z,\ell} + p_{z,\nu})^2 \quad , \quad (5.3)$$

with the lepton transverse momentum $\vec{p}_{T,\ell}$, the lepton energy E_ℓ , and the z -components of the momentum of the lepton $p_{z,\ell}$ and neutrino $p_{z,\nu}$. Assuming an on-shell production of the W boson the invariant mass of the W boson can be set to $m_W = 80.4$ GeV, thus allowing the solving of the equation for $p_{z,\nu}^\pm$. The solution for the quadratic equation can be written as

$$p_{z,\nu}^\pm = \frac{\Lambda \cdot p_{z,\ell}}{p_{T,\ell}^2} \pm \sqrt{\frac{\Lambda^2 \cdot p_{z,\ell}^2}{p_{T,\ell}^4} - \frac{E_\ell^2 \cdot \cancel{E}_T^2 - \Lambda^2}{p_{T,\ell}^2}} \quad (5.4)$$

with the abbreviation

$$\Lambda = \frac{m_W^2}{2} + \vec{p}_{T,\ell} \cdot \vec{\cancel{E}}_T \quad . \quad (5.5)$$

Depending on the sign of the discriminant in Equation (5.4) two cases arise: If the discriminant is positive, two different solutions can be found out of whom the one with the smaller absolute value is selected. If the discriminant is negative, caused by a finite \cancel{E}_T resolution in the detector, the p_x and p_y components of the neutrino momentum are varied such that the discriminant is zero. Of the two possible solutions, the one with the minimal distance between $p_{T,\nu}^\pm$ and $\vec{\cancel{E}}_T$ is chosen. This reconstruction method is adapted from Reference [156].

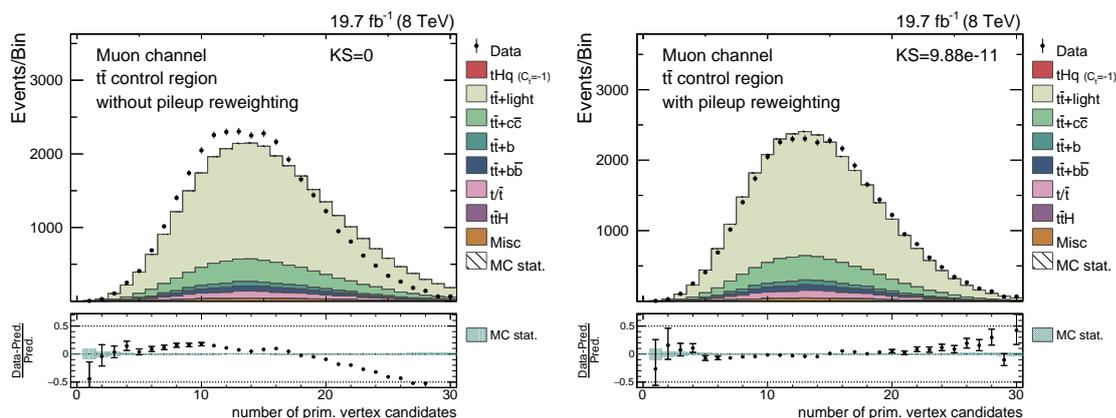


Figure 5.4.: Distribution of the number of primary vertices per event before (left) and after (right) application of the pileup reweighting. The simulation predicts more primary vertices than are observed in data. For both distributions the simulation is scaled to match the observed event numbers in data. After the reweighting a good agreement of simulation and recorded events is observed.

5.6. Monte Carlo Corrections

5.6.1. Pileup Reweighting

In the generation of the Monte Carlo samples the expected pileup scenario has to be added to the generated primary collision as explained in Section 3.1.4. The actual pileup environment differs from the simulation thus creating a discrepancy primarily visible in the distribution of the number of primary vertices (NPV) reconstructed in a collision. Events are reweighted such that the actual NPV distribution is reproduced in simulation. The situation before and after reweighting can be seen in Figure 5.4. A systematic uncertainty is assigned to this procedure and is explained further in Section 5.10.

5.6.2. Lepton Efficiency Scale Factors

Various discrepancies can be found when comparing lepton efficiencies between simulation and real data.

The calculated muon scale factors correct for different trigger efficiencies, isolation discrepancies, an imperfect tracker response and deviations of the muon ID efficiencies. The muon scale factors are parametrized as function of pseudorapidity and have been derived by utilizing a tag-and-probe technique in a Drell-Yan enriched phase space [157]. These corrections are made centrally available by the Muon POG.

The electron efficiency scale factors are correcting differences in electron ID, isolation and trigger efficiencies. The scale factors for the triggering MVA ID are employed in this analysis [153]. The scale factors are a function of electron p_T and η of the electron superclusters and are made centrally available by the E/gamma POG.

5. Search for tHq production at $\sqrt{s} = 8$ TeV

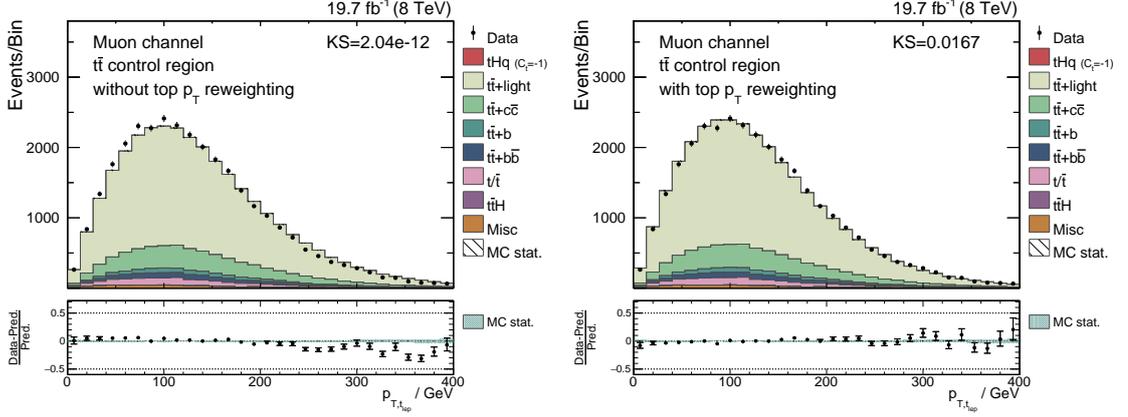


Figure 5.5.: Distribution of the transverse momentum of the reconstructed leptonically decaying top quark of the $t\bar{t}$ process before (left) and after (right) application of the top quark p_T reweighting. The simulation predicts a harder momentum spectrum of the top quarks than what is observed in data. For both distributions the simulation is scaled to match the observed event numbers in data. After the reweighting a good agreement of simulation and recorded events is observed.

5.6.3. Top Quark p_T Reweighting

The simulated $t\bar{t}$ sample generated with MADGRAPH predicts a harder p_T spectrum for top quarks than what is observed in real data. This effect is propagated to all of the top quark decay products. NNLO predictions [158] give a reasonable description of the top quark p_T . Event weights for the $t\bar{t}$ simulation sample have been derived globally to account for these seen discrepancies.

The CMS collaboration provides parametrized event weights dependent on the generator-level top quark p_T . The effect of the reweighting on the transverse momentum of the reconstructed leptonically decaying top quark can be seen in Figure 5.5.

5.6.4. b-tagging Efficiency Correction

A corrections of b-tagging efficiencies is needed, as the performance of tagging algorithms on simulated samples differs from their performance in data events. With scale factors provided by the CMS b-tagging & vertexing POG [159] the tagging efficiencies can be adjusted properly. The correction procedure applied in this analysis is adapted from [160].

The probability for an arbitrary event to get assigned a certain number of b tags can be calculated as

$$\mathcal{P}_{MC} = \prod_{i \in \text{tagged}} \epsilon_i \cdot \prod_{j \notin \text{tagged}} (1 - \epsilon_j) \quad , \quad (5.6)$$

with the simulated b-tagging efficiencies ϵ_i for jet i and the corresponding mistag efficiencies $1 - \epsilon_i$. The multiplication of tagging probabilities of all tagged jets with the mistagging probabilities of the untagged jets yields the probabilities to assign an arbitrary number of b tags to this event.

The same can be done for real data with

$$\mathcal{P}_{\text{Data}} = \prod_{i \in \text{tagged}} s_i \epsilon_i \cdot \prod_{j \notin \text{tagged}} (1 - s_j \epsilon_j) \quad , \quad (5.7)$$

where s are the p_T - and η -dependent scale factors needed to reproduce the same b-tag configurations. Every event is then assigned a weight w according to

$$w = \frac{\mathcal{P}_{\text{Data}}}{\mathcal{P}_{\text{MC}}} = \prod_{i \in \text{tagged}} s_i \cdot \prod_{j \notin \text{tagged}} \left(\frac{1 - s_j \epsilon_j}{1 - \epsilon_j} \right) \quad . \quad (5.8)$$

The b-tagging efficiencies in simulation have been determined for this analysis as they strongly depend on the event selection. The efficiencies are a function of jet transverse momentum and pseudorapidity. The efficiencies are calculated in four different jet flavor categories, separating b-flavored, c-flavored and light flavored jets as well as a combined category for gluon jets or unmatched jets.

5.6.5. Jet Pseudorapidity in the Forward Region

One of the main discriminant features of the tHq signal process is the light forward jet. Unfortunately, a severe mismodeling of the jet pseudorapidity in the forward region is evident. The mismodeling is clearly visible in ratio distributions when comparing the yields found in data and simulation.

The main features are an overestimated number of jets in simulation in the region with $2.2 \leq |\eta| \leq 2.9$ followed by an underestimation in the region with $2.9 \leq |\eta| \leq 3.5$. In the most forward region of the detector with $|\eta| > 3.5$ the predicted number of jets in simulation again exceeds the number measured. The effect is most apparent for low- p_T jets, but still visible at the higher end of the p_T spectrum, as can be seen in Figure 5.6.

A comparison of the reconstructed and the generated jet- η distributions revealed that the depression around $|\eta| \sim 2.7$ is caused by migration effects during the jet reconstruction. Jets with $|\eta_{\text{gen}}|$ around 2.5 and 3.1 are more frequently reconstructed with a $|\eta_{\text{reco}}| \sim 2.7$. This leads to an overestimation of simulation events in this region, thus the depression seen in the ratio plot. This migration in jet- η bins is limited to a region $2.4 < |\eta| < 3.2$ resulting in a constant normalization in this region. As a solution for this effect the affected region is taken as a single bin in the further analysis.

The steep descend in the ratio plot for $|\eta| > 3.1$ can be understood as an effect of the L2L3Residual corrections, which are described in Section 3.3.6. Due to low number of events in this region the corrections are derived in a large single bin with $|\eta| > 3.139$. Although the corrections cause the clearly visible slope in the ratio plot as the detector response gets lower with higher pseudorapidity values, the normalization averaged over the correction bin width in this region is still correct. Again, the solution for this discrepancy is the treatment of this region as a single bin. To ensure a coherent treatment of the pseudorapidity for the further analysis, the actual η values of the jets in the forward region are modified:

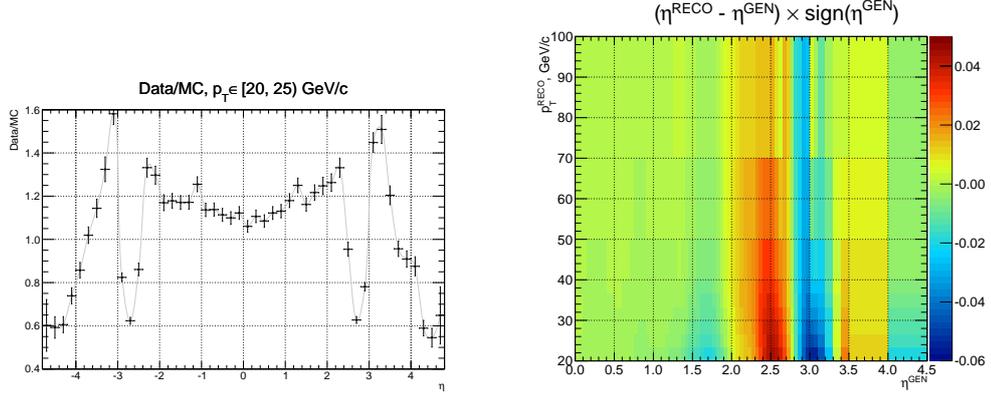


Figure 5.6.: The mismodeling observed in the jet pseudorapidity. On the left the ratio of observed to predicted event yields in bins of η for low- p_T jets is shown. An overestimated number of jets at $|\eta| \sim 2.7$ is followed by an underestimation at $|\eta| \sim 3.2$. In the most forward regions the number of predicted jets again exceeds the observed number.

On the right the jet- η migration effect during the jet reconstruction is depicted. Jets generated with $|\eta| \sim 2.5$ are reconstructed with a slightly higher jet $|\eta| \sim 2.7$. In contrast jets generated with $|\eta| \sim 3$ are reconstructed with lower values. This migration towards $|\eta| \sim 2.7$ from both sides causes an overestimation of events in that region and an underestimation in the migrated from regions. The diagrams are taken from [150].

$$\eta \rightarrow \eta' = \begin{cases} \eta & \text{if } |\eta| < 2.4, \\ 2.8 \cdot \text{sign}(\eta) & \text{if } 2.4 \leq |\eta| \leq 3.2, \\ 3.5 \cdot \text{sign}(\eta) & \text{if } 3.2 \leq |\eta|. \end{cases}$$

The transformation is applied to simulation as well as data, making the redefined pseudorapidity an effective observable. The analysis of this chapter only relies on a good agreement of simulation with the recorded data in the used variables, which is provided in these newly introduced large bins. The analysis is therefore blind to the mismodeling in the forward region. After the applied correction a slightly different normalization in the forward region in simulation is still observed, but this can be cured by requiring $p_T > 40 \text{ GeV}$ for jets with $|\eta| > 2.4$.

5.7. Event Selection

Especially in a search for a tiny signal it is essential to limit the search to a well defined phase space with an enhanced signal-to-background ratio (S/B). As explained in Section 5.2 the most protruding feature of the tHq process is the high b-jet multiplicity. With two b quarks from the Higgs boson decay, one from the top quark decay and one from the initial gluon splitting, up to four bottom quarks are expected to appear in the detector. In contrast to the decay products of the centrally produced top quark and Higgs boson, the second b quark is more likely to be produced in a forward direction. Combined with its lower expected transverse momentum the

jet originating from this b quark is less likely to be b-tagged by the CSV algorithm. Therefore, two different signal-enhanced b-tag multiplicity categories are defined, the three-tag region (3T) and the four-tag region (4T), where three or four jets with a b-tagging value larger than 0.898 are required, respectively. The considered jets must fulfill $p_T > 20$ GeV and $|\eta| < 2.4$, as information needed for b-tagging is only available in the central region covered by the silicon tracker.

The tight working point of the CSV algorithm is chosen although a lower working point would increase the number of signal events in this region, but also the $t\bar{t}$ +light background would be enhanced significantly. Even though the tight working point has a mistag probability of only 0.1% the $t\bar{t}$ +light background, which needs at least one falsely tagged jet to enter the signal region, is still the largest contribution to the background averaged over the signal regions. Changing to the medium (1% mistag probability) or the loose working point (10% mistag probability) would hence increase the contribution from $t\bar{t}$ +light roughly by factors of 10 or 100, respectively.

The leptonic decay of the top quark results in one charged lepton and its corresponding neutrino. In the case of a hadronic decay the analysis would be missing a clear object for the trigger system and a large contribution of QCD multijet production would be expected thus making a sensitive analysis, at least in the $H \rightarrow b\bar{b}$ decay channel, impossible. The analysis therefore requires exactly one tight lepton and no additional loose leptons. Here, tight and loose leptons only refer to muons or electrons. The analysis is performed in two separate categories, the electron channel and the muon channel.

As a basis for the upcoming event reconstruction, where jets are assigned to a total of four quarks, every event must offer at least four *reconstructable jets*. If a reconstructable jet is found in the central detector region with $|\eta| < 2.4$, it needs to satisfy $p_T > 30$ GeV or else a stricter cut is required with $p_T > 40$ GeV, due to the mismodeling in the forward region explained in Section 5.6.5.

To account for the light forward quark, every event is required to have at least one untagged jet. In the 3T region this requirement is redundant, as this case is already covered by the collection of reconstructable jets.

As a last requirement the measured missing transverse energy needs to exceed a certain threshold to account for the neutrino in the final state. Requirements on the transverse missing energy are chosen as $\cancel{E}_T > 35$ GeV for events with a muon and $\cancel{E}_T > 45$ GeV for events with an electron in the final state such that the vast majority of remaining QCD multijet events is suppressed, rendering this background negligible. A summary of the applied selection criteria can be found in Table 5.1.

In addition to the two signal regions a $t\bar{t}$ control region (2T) is defined which requires exactly two jets with a CSV output value above 0.898, but is in every other aspect defined in analogy to the signal regions. The control region is used to study the modeling of distributions in simulation and search for discrepancies when comparing simulation to data. A selection of control distributions can be found in Figure 5.7.

5. Search for tHq production at $\sqrt{s} = 8 \text{ TeV}$

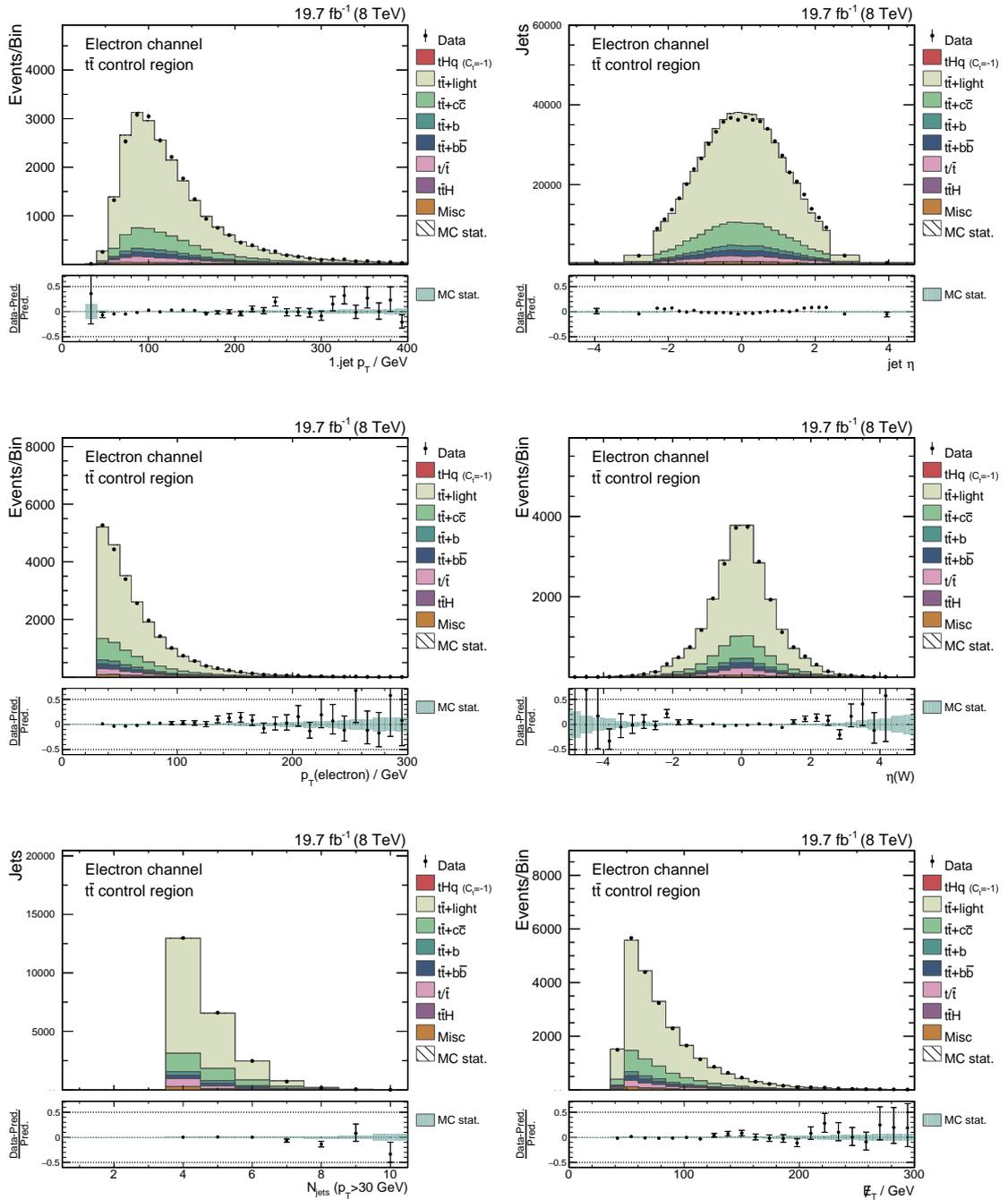


Figure 5.7.: A set of control plot in the $t\bar{t}$ control region in the electron channel are shown. The diagrams show the transverse momentum of the hardest jet of the event, the pseudorapidity of all selected jets, the transverse momentum of the electron, the pseudorapidity of the reconstructed W boson, the number of jets with a transverse momentum greater than 30 GeV and the missing transverse energy of the event. For all distributions the simulation is scaled to match the observed event numbers in data. A good agreement of simulation and data is observed. The corresponding diagrams in the muon channel can be found in the Appendix A.1.

Table 5.1.: Overview of the selection criteria applied to define the two signal regions, the three-tag region (3T) and the four-tag region (4T).

Reconstructable jets are defined as having $p_T > 30$ GeV in the central region and a more strict requirement in the forward region with $p_T > 40$ GeV. The additional jets are part of the standard jet collection with $p_T > 20$ GeV in the central region and $p_T > 40$ GeV in the forward region.

	3T region	4T region
# jets with CSV > 0.898	3	4
# tight leptons	1	1
# additional loose leptons	0	0
# reconstructable jets	≥ 4	≥ 4
# additional jets	≥ 0	≥ 1
\cancel{E}_T	$> 35/45$ GeV(μ/e)	$> 35/45$ GeV(μ/e)

5.8. Event Reconstruction

The most discriminating features of the signal process are the properties of the actual physical objects, such as the top quark or the Higgs boson. The assignment of jets measured in the detector to the expected particles of the process, however, is ambiguous. In the course of this analysis a reconstruction technique using multivariate analysis tools has been developed and is used to resolve this ambiguity. Artificial Neural Networks are trained to differentiate between correct and wrong jet assignments based on variables describing the reconstructed objects and their correlations.

A tHq event where all final state particles are correctly assigned to the jets is well separable from background events that are subjected to the same jet assignment. They naturally show a different behavior as the jets of the event get matched to hypothetical partons that are not necessarily part of the background process.

The same would hold for a reversed statement: If all partons apparent in a $t\bar{t}$ event could be correctly matched to the available jets, this event could consequently also be well separable from the signal events, where the assigned particles do not match the actual particle content. Hence, this analysis uses two similar event reconstructions that are performed in parallel but assign jets under a certain event hypothesis, a tHq reconstruction and a $t\bar{t}$ reconstruction. Both procedures are explained in detail in the following sections.

5.8.1. tHq Reconstruction

The purpose of the tHq reconstruction is the assignment of measured jets to the four main partons of the tHq process: the two b quarks of the Higgs boson decay, the b quark from the top quark decay and the light forward quark.

A neural network is trained to separate correct assignments from wrong assignments in each event, where the correct assignment is found, when each of the four quarks can be matched to a reconstructed jet within a radius of $\Delta R < 0.3$. If not all quarks can be matched, the event is discarded for the training. All other possible permutations but the correct assignment are

5. Search for tHq production at $\sqrt{s} = 8 \text{ TeV}$

Table 5.2.: Settings used in the employed neural nets for the two reconstructions in the TMVA software package. The definitions of the configuration options can be found in Reference [117].

Parameter	Value
NCycles	500
HiddenLayers	30
NeuronType	tanh
VarTransform	N
NeuronInputType	sum
EstimatorType	MSE

further seen as wrong.

For each event there are $N_{\text{jet}}!/(N_{\text{jet}} - 4)!$ possible jet assignments out of which only one is the correct hypothesis. In order to reduce the number of possible permutations a set of requirements has to be fulfilled by each considered assignment. The jets that are considered to be assigned to one of the three b quarks have to lie central in the detector to facilitate the exploitation of the CSV values of these jets. Furthermore, jets that are assigned to the light quark must not be b-tagged, hence fulfill $\text{CSV}_{\text{jet}} < 0.898$. For the actual training a randomly chosen wrong assignment is used for each correct assignment in order to ensure two equally large training samples.

For the training of the neural networks the Broyden-Fletcher-Goldfarb-Shanno (BFGS) training algorithm [119–122] as implemented in TMVA is used. The values used as training parameters can be found in Table 5.2.

The neural network is trained using twelve variables, whose description can be found in Table 5.3, sorted by their importance in the training. The most important variables are related to the number of b-tagged jets assigned to the b quark of the top quark decay and to the two b quarks of the Higgs boson decay. Also a very important variable is the absolute value of the pseudorapidity of the jet assigned to the light quark. The distributions of the six most important variables for wrong and correct assignments can be found in Figure 5.8. The remaining variables can be found in the Appendix A.2. The linear correlations of these variables for correct and wrong assignments can be seen in Figure 5.9.

A fifth of all available tHq simulation events of the 3T and 4T region are exclusively used for the training of the neural network and are subsequently discarded in order to avoid a training bias. Out of these events a smaller subset is used as an independent testing sample. The response of the neural network for the training and test sample can be found in Figure 5.10. The neural network is able to separate well between a set of correct and wrong assignments and the net performs equally well on the test sample, indicating that no overtraining has occurred.

After the successful training the neural network is employed to choose the best possible assignment in all simulation and data events. Each event is subjected to the neural network and for each of the allowed jet assignments of this event the physical objects are reconstructed. The neural network assigns an output value to each jet assignment based on the twelve variables

Table 5.3.: Input variables for the jet-assignment neural network under the tHq hypothesis sorted by their importance in the training. Instead of the transverse momenta variables the logarithm of these variables is used, as the neural net can process the information of narrower distributions better.

Variable	Description
tagged jet (b_t)	boolean variable indicating if the jet assigned to the b quark of the top quark decay is b tagged
$\eta(\text{light jet})$	absolute pseudorapidity of the light forward jet
# b tags of Higgs jets	number of b tags assigned to the jets of the Higgs boson decay
$\log m(b_t + l)$	invariant mass of the jet assigned to the b quark from the decay of the top quark and the charged lepton
$\log m(H)$	invariant mass of the reconstructed Higgs boson
$\log \min(p_T(\text{Higgs jets}))$	lower transverse momentum of the two jets assigned to the Higgs boson decay products
$\Delta R(\text{Higgs jets})$	ΔR between the two jets from the decay of the Higgs boson
$\max \eta(\text{Higgs jets})$	higher pseudorapidity of the two jets assigned to the Higgs boson decay products
$\Delta R(b_t, W)$	ΔR between the jet assigned to the b quark and the W boson from the top quark decay
relative H_T	percentage of the total transverse momenta (jets, lepton, \cancel{E}_T) that falls to the reconstructed top quark, Higgs boson and light forward jet
$\Delta R(H, t)$	ΔR between reconstructed top quark and Higgs boson
$q(b_t) \cdot q(l)$	jet charge (see definition in Reference [161]) of the jet assigned to the b quark of the top quark decay multiplied by the charge of the lepton

5. Search for tHq production at $\sqrt{s} = 8$ TeV

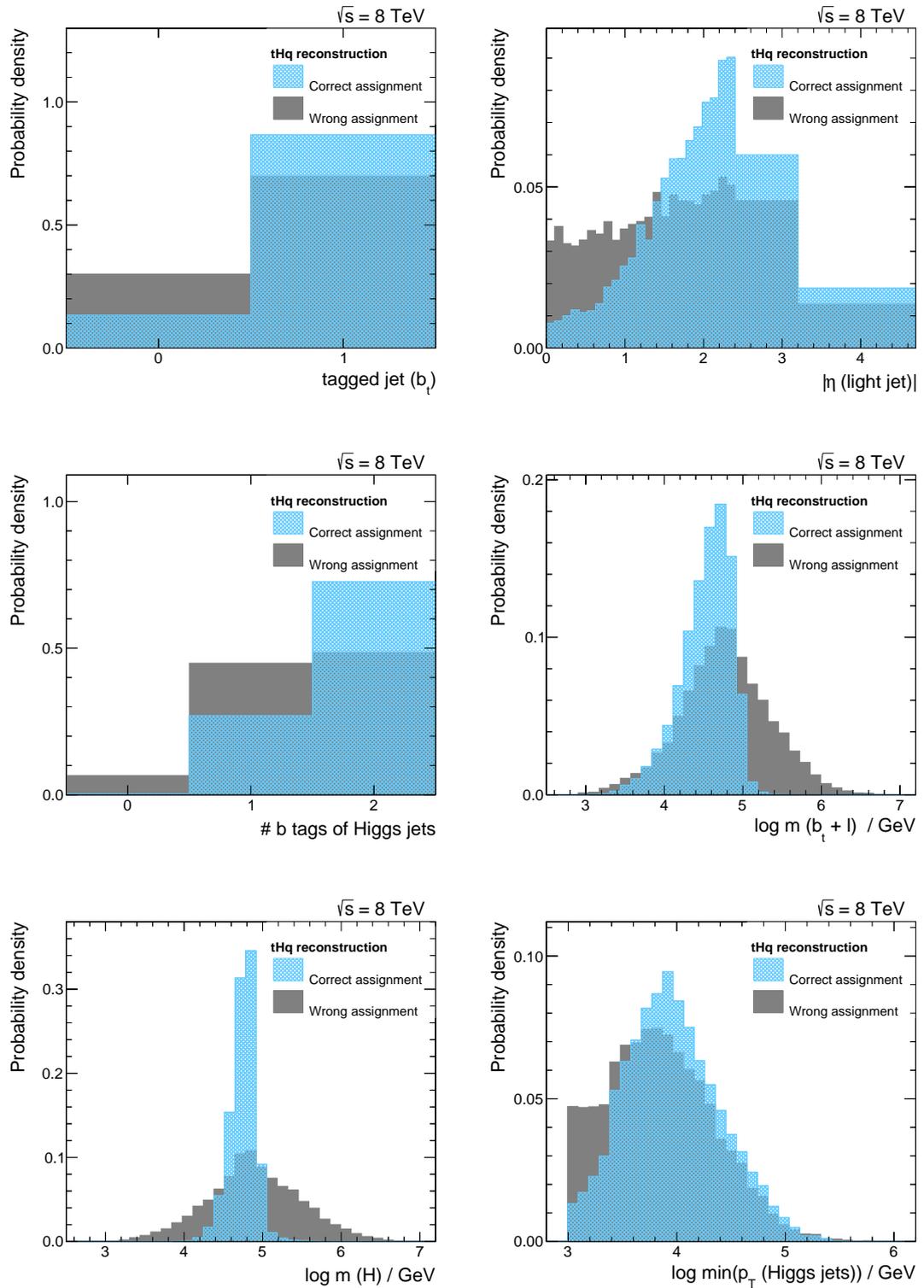


Figure 5.8.: The six most discriminating variables between correct and wrong jet assignments in the tHq reconstruction are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 5.3.

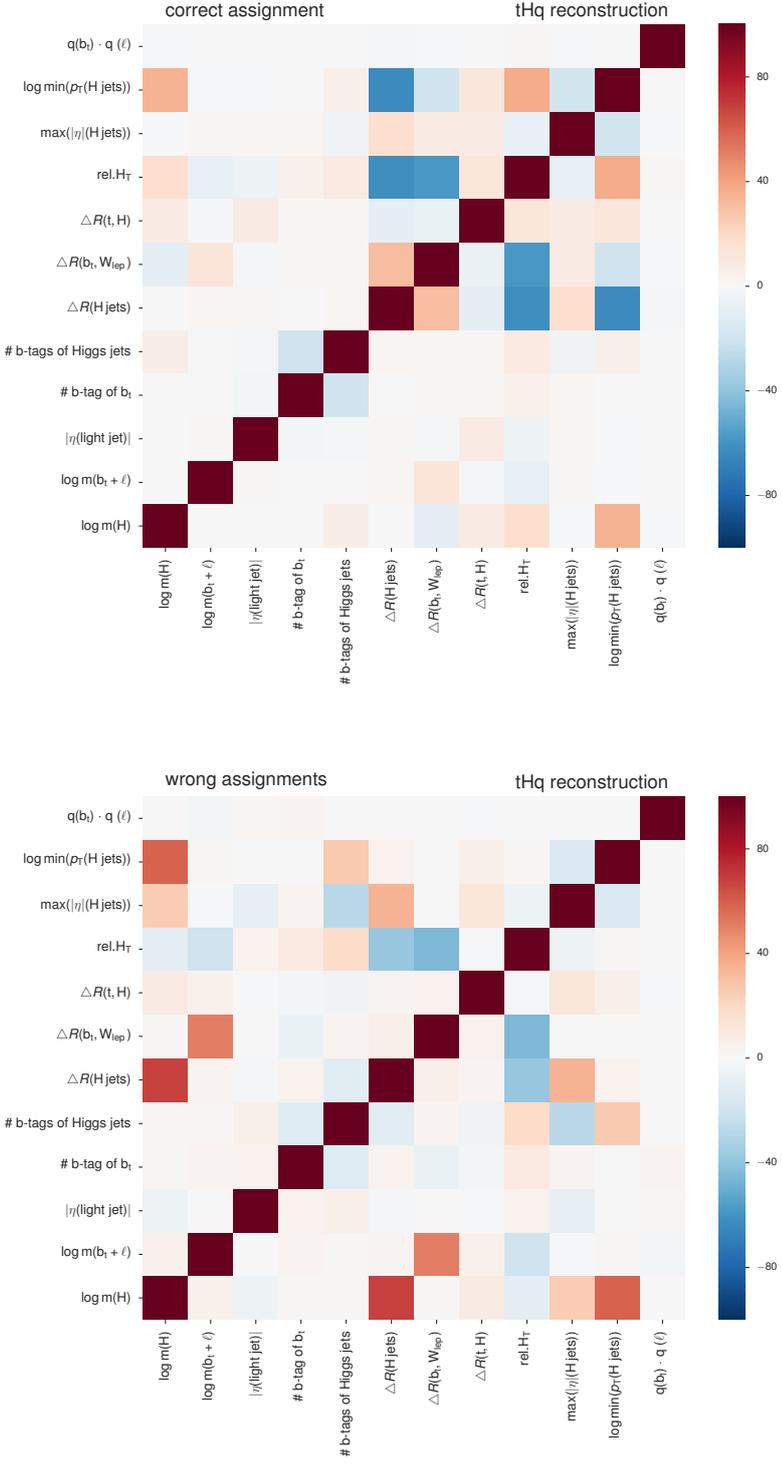


Figure 5.9.: Correlations between variables used in the tHq reconstruction for correct and wrong jet assignments. It is clearly visible that correlations are apparent between variables for correct assignments that are not visible for wrong assignments, such as the correlation between the ΔR of the two Higgs jets and the transverse momentum of the softer of the two Higgs jets, and vice versa.

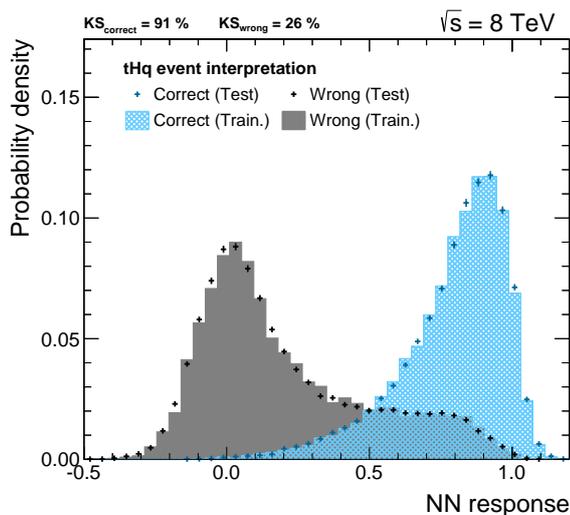


Figure 5.10.: Output values for the tHq reconstruction neural network for correct and wrong jet assignments. A clear separation between correct and wrong assignments is visible. The output of the neural net for the events used in the training are shown in the colored histograms. The training is examined with an independent set of events that were not part of the training sample. The response for this testing sample is shown as colored markers and a good agreement between the testing and training sample is seen, verified by high KS-values of the two distributions. No sign of overtraining can be found.

of the reconstructed objects. Subsequently, the assignment that is attributed with the highest output value is henceforth used for this event in the analysis. The application is done for both, simulation samples and actual data.

As further test of the reconstruction the output values of the tHq reconstruction for a randomly chosen assignment and for the assignment with the highest output value are compared between simulation and data in the $t\bar{t}$ control region. The result can be found in Figure 5.11 and a good agreement is observed for both distributions.

5.8.2. $t\bar{t}$ Reconstruction

The analysis does not only gain from a correctly reconstructed tHq event by giving each event a probability how tHq -like the event is, but also from a $t\bar{t}$ -based reconstruction. By telling how $t\bar{t}$ -like an event is, the eventually sought separation between signal and background events can be enhanced even further.

The $t\bar{t}$ reconstruction assigns the reconstructed jets to the partons of a semi-leptonic $t\bar{t}$ decay: the two light quarks of the W boson decay from the hadronically decaying top quark t_{had} , the b quark of the decay of t_{had} and the b quark of the leptonically decaying top quark t_{lep} .

In accordance with the tHq reconstruction the correct assignment is found when all four partons can be matched to a jet within a cone radius of $\Delta R < 0.3$. If no correct assignment can be found, the event is again discarded from the training. All other possible jet assignments are seen as

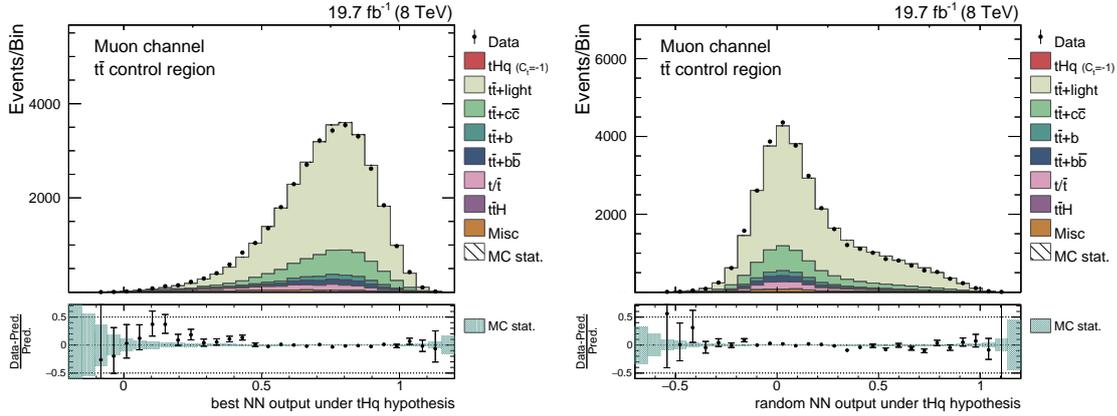


Figure 5.11.: Response of the tHq reconstruction comparing simulation to data. On the left the highest output value (chosen jet assignment) per event is shown and on the right the NN output value for a random assignment is shown. Both diagrams are shown for the muon channel in the $t\bar{t}$ control region. In both distributions a good agreement between simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all simulation corrections are applied. The corresponding distributions for the electron channel and 3T and 4T region can be found in Appendix A.3 and A.4.

wrong assignments.

To reduce the number of possible jet permutations a set of constraints is applied. Both jets that are assigned to the b quarks of the two top quark decays have to be b -tagged ($CSV > 0.898$). As b -tagging information is only available in the central region covered by the silicon tracker a pseudorapidity restriction $|\eta| < 2.4$ is implied.

A neural network with the same setting as in the tHq reconstruction is employed and thirteen different variables are used in the training. One fifth of all events in the 3T and 4T region of the semi-leptonic $t\bar{t}$ simulation sample is used as training sample and is subsequently excluded from further analysis. A list of all thirteen training variables sorted by their importance can be found in Table 5.4 and the distributions for correct and wrong assignments for the six most important variables can be found in Figure 5.12. The remaining seven variables can be found in the Appendix A.5. The linear correlations of the variables used in the reconstruction for correct and wrong assignments can be found in Figure 5.13.

The most important variables in the $t\bar{t}$ reconstruction are the invariant mass of the reconstructed W_{had} and its difference to the invariant mass to the reconstructed t_{had} . The most challenging part in this reconstruction is the jet assignment to the two light quarks of the W_{had} decay as no constraints of the possible jet candidates are posed. The response of the neural network for the training sample and a disjoint testing sample can be seen in Figure 5.14.

After the training the jet assignment is performed for all simulation and data events. For each event every allowed jet assignment is iterated over and the output of the neural network is calculated and compared. The jet assignment with the highest response value is chosen for the further analysis. A sanity check of the reconstruction technique is performed by comparing the NN response in simulation to its response for data events. This is done for the chosen

5. Search for tHq production at $\sqrt{s} = 8 \text{ TeV}$

Table 5.4.: Input variables for the jet-assignment MVA under the $t\bar{t}$ hypothesis sorted by their importance in the training. Instead of the transverse momenta variables the logarithm of these variables is used, as narrow distributions are better suited for the usage in MVA techniques than distributions with long tails.

Variable	Description
$\log m(W_{\text{had}})$	invariant mass of the two jets assigned to the W boson of t_{had}
$\log (m(t_{\text{had}}) - m(W_{\text{had}}))$	difference between the invariant masses of reconstructed t_{had} and W_{had}
$\Delta R (W_{\text{had}})$	ΔR between the two jets assigned to the W boson of t_{had}
$ \eta(t_{\text{had}}) $	absolute value of the pseudorapidity of the reconstructed t_{had}
$\log p_{\text{T}} (t_{\text{had}})$	transverse momentum of the reconstructed t_{had}
# b tags of W_{had} jets	number of b-tagged jets assigned to W_{had}
$\log p_{\text{T}} (t_{\text{lep}})$	transverse momentum of the reconstructed t_{lep}
$\Delta R (b_{t_{\text{had}}}, W_{\text{had}})$	ΔR between the reconstructed t_{had} and W_{had}
relative H_{T}	percentage of the total transverse momenta (jets, lepton, \cancel{E}_{T}) that falls to the reconstructed t_{had} and t_{lep}
$\Delta R (b_{t_{\text{lep}}}, W_{\text{lep}})$	ΔR between the reconstructed t_{lep} and W_{lep}
$\log m(b_{t_{\text{lep}}} + \ell)$	invariant mass of the jet assigned to the b quark of t_{lep} and the charged lepton
$(q(b_{t_{\text{lep}}}) - q(b_{t_{\text{had}}})) \cdot q(\ell)$	difference of the jet charges of the jets assigned to the b quarks of the top quark decays multiplied by the lepton charge
$(\sum q(W_{\text{had}} \text{ jets})) \cdot q(\ell)$	sum of the jet charges assigned to W_{had} multiplied by the lepton charge

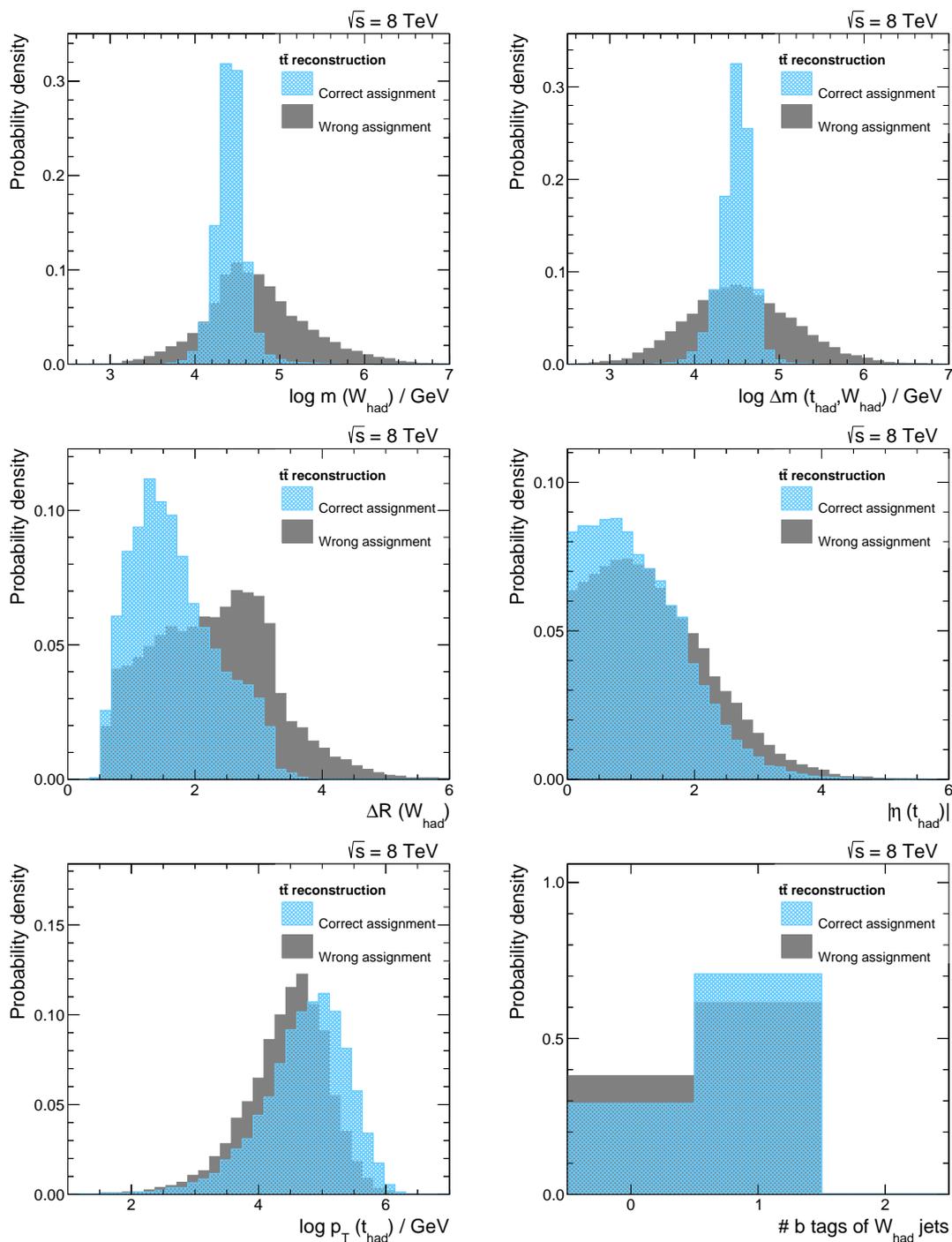


Figure 5.12.: The six most discriminating variables between correct and wrong jet assignments in the $t\bar{t}$ reconstruction are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 5.4. The remaining seven variables can be found in Appendix A.5.

5. Search for tHq production at $\sqrt{s} = 8$ TeV

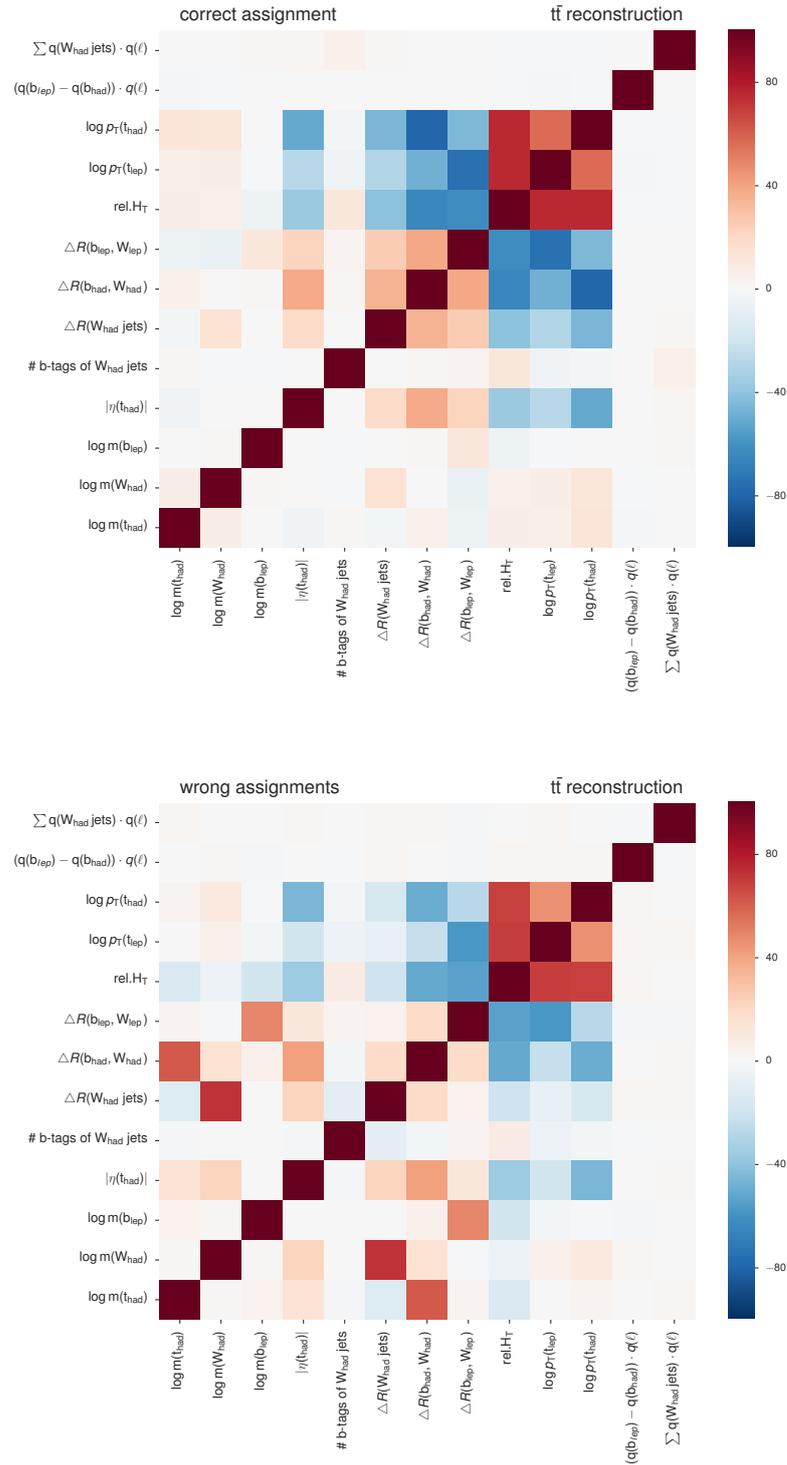


Figure 5.13.: Correlations between the variables used in the $t\bar{t}$ reconstruction for correct and wrong jet assignments. It is visible that stronger (anti-)correlations are apparent between variables for correct assignments than are for wrong assignments. Additionally, correlations are observed in wrong assignments that would be almost uncorrelated for correct assignments, such as the correlation between ΔR of b quark and W boson and the invariant mass of the reconstructed hadronically decaying top quark.

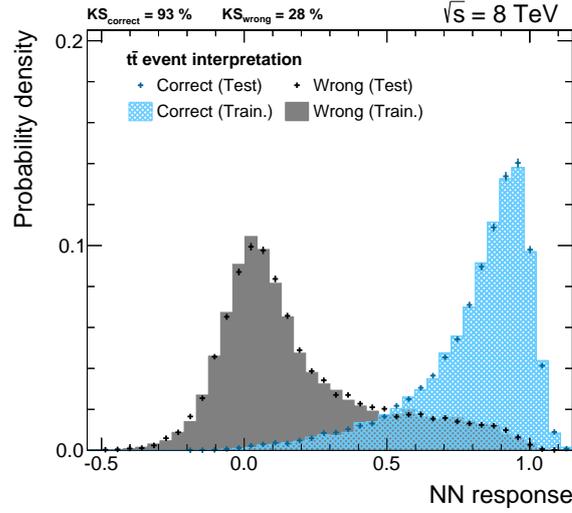


Figure 5.14.: Output values for the $t\bar{t}$ reconstruction neural network for correct and wrong jet assignments. A clear separation between correct and wrong assignments is visible. The output of the neural net for the events used in the training are shown as histograms. The training is examined with an independent set of events that were not part of the training sample. The response for this testing sample is shown as colored markers and a good agreement between the testing and training sample is seen, verified by high KS-values of the two distributions. Thus, no sign of overtraining can be found.

assignment for each event and for a random jet assignment. The comparison can be found in Figure 5.15 and a good agreement is observed in both variants.

After the application of both reconstructions, the tHq and the $t\bar{t}$ reconstruction, every event is provided with assignments to the particles of both final states. Kinematic observables of these reconstructed objects build the base for the final event classification explained in the following chapter.

5.8.3. Evaluation of Reconstruction Methods

In order to quantify the performance of the reconstruction the percentages in which the multivariate reconstruction attributes the correct jets to the single partons are calculated. As reference a simpler reconstruction method is implemented which uses a χ^2 measure to select one of the allowed jet assignments.

tHq Reconstruction Evaluation

For the tHq reconstruction the χ^2 method is implemented such that the jet assignment is performed by choosing the assignment that minimizes

$$\chi^2 = \frac{(m'_t - m_t)^2}{\sigma_t^2} + \frac{(m'_H - m_H)^2}{\sigma_H^2} ,$$

5. Search for tHq production at $\sqrt{s} = 8$ TeV

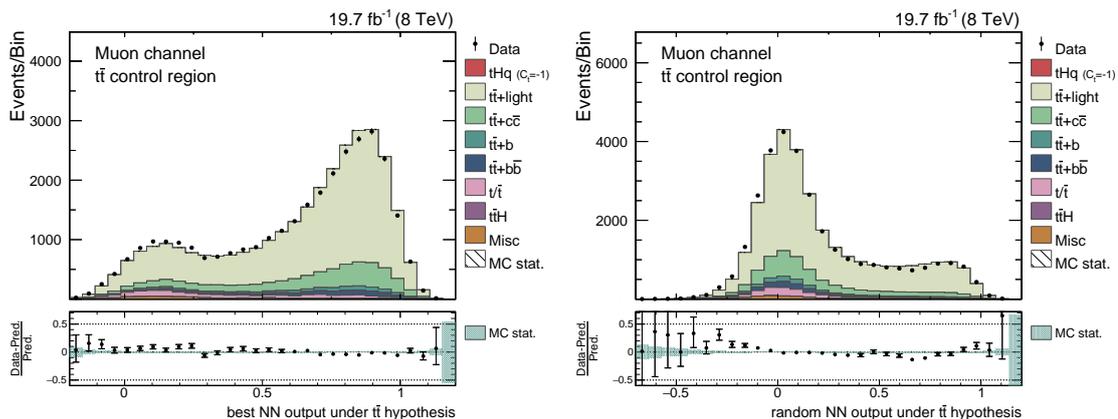


Figure 5.15.: Response of the $t\bar{t}$ reconstruction comparing simulation to data. On the left the highest output value (chosen jet assignment) per event is shown and on the right the NN output value for a random assignment is shown. Both diagrams are shown for the muon channel in the $t\bar{t}$ control region. In both distributions a good agreement between simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all simulation corrections are applied. The corresponding distributions for the electron channel and 3T and 4T region can be found in Appendix A.6 and A.7.

where the primed variables are the invariant masses of the reconstructed objects. As masses of top quark and Higgs boson $m_t = 173.21$ GeV and $m_H = 125.04$ GeV are used with corresponding mass resolutions of $\sigma_t = 30$ GeV and $\sigma_H = 15$ GeV. The jet with the highest $|\eta|$ -value is chosen as light forward jet. The comparison is performed on events where each of the four quarks in the final state can be matched to a jet within a cone radius of $\Delta R < 0.3$. As a figure of merit the percentages are calculated how often the two different reconstruction techniques successfully assign the matched jet to the individual partons. The result of this study can be seen in Figure 5.16(a).

The comparison shows that the reconstruction employing multivariate analysis technique outperforms the χ^2 reconstruction for almost every object. The bottom quark of the top quark decay as well as both bottom quarks of the Higgs boson decay are assigned correctly in 66% of all events in the 3T region. The χ^2 only reaches an assignment efficiency of 54% and 58%, respectively. The light forward jet is matched equally well (80%) in both reconstruction approaches. The efficiencies to correctly assign all four partons in the event are 47% (NN) and 37% (χ^2) in the 3T region and 39% (NN) and 28% (χ^2) in the 4T region. The reconstruction efficiencies are systematically lower for the jet assignment to the b quarks in the 4T region, as variables like the number of b-tagged jets assigned to the Higgs boson decay jets lose some of their significance due to the higher b-tag multiplicity. Although the assignment of the jet to the light quark simply based on the maximal jet pseudorapidity is as often correct as the NN method, the NN assignment is seen as superior as the $|\eta_{\text{lightjet}}|$ observable obtained by the NN reconstruction method yields a stronger separation between signal events and background events in the classification than the $|\eta_{\text{max}}|$ observable.

$t\bar{t}$ Reconstruction Evaluation

For the $t\bar{t}$ reconstruction the χ^2 reconstruction is implemented such that the jet assignment is chosen that minimizes

$$\chi^2 = \frac{(m'_{t_{lep}} - m_t)^2}{\sigma_t^2} + \frac{(m'_{t_{had}} - m_t)^2}{\sigma_t^2} + \frac{(m'_{W_{had}} - m_W)^2}{\sigma_W^2} .$$

The primed values are the invariant masses of the reconstructed objects based on the assigned jets. The value of the used invariant masses for top quarks and the W boson and their corresponding mass resolutions are $m_t = 173.21$ GeV, $m_W = 80.4$ GeV, $\sigma_t = 30$ GeV and $\sigma_W = 15$ GeV. The efficiency is calculated in a simulated semi-leptonically decaying $t\bar{t}$ sample, where to each of the assignable quarks a jet can be matched within a radius of $\Delta R < 0.3$. It is subsequently checked how often the two reconstruction methods assign the matched jet to its corresponding particle. The result of this study can be found in Figure 5.16(b).

In both, the 3T and the 4T region, the NN reconstruction shows a clear improvement over the χ^2 reconstruction. In both regions the b quark of the leptonically decaying top quark is the parton with the highest assignment success rate of 77% (3T) and 59% (4T) using the NN reconstruction method. The correct reconstruction of the complete hadronically decaying top quark is aggravated due to the two light quarks from the W boson decay. In 61% (3T) and 42% (4T) of the events the hadronically decaying W boson is reconstructed correctly, thereby limiting the overall efficiency. A completely correct reconstruction of the $t\bar{t}$ system is therefore possible in 52% (3T) and 25% (4T) of all events with the NN reconstruction. The χ^2 reconstruction leads only to a correct jet assignment in 29% (3T) and 16% (4T) of the events.

5.9. Event Classification

The overwhelming $t\bar{t}$ background makes it crucial to find variables with differing shapes for the sought signal and the sum of all backgrounds. Utilizing the different jet assignments a set of eight variables has been defined with powerful discriminating power. The set contains three variables that describe objects reconstructed under the tHq hypothesis, four variables describing objects reconstructed under the $t\bar{t}$ hypothesis and one reconstruction-independent variable. They are listed in Table 5.5.

In order to eke the most out of the given variable set, a third neural network is employed. The neural network is given two sets of events, a signal sample containing tHq events of the 3T region and a background sample containing a mixture of semi-leptonic $t\bar{t}$ events, full-leptonic $t\bar{t}$ events as well as $t\bar{t}H$ events of the 3T region. Based on the distributions of the eight discriminating variables and their correlations, the neural network produces a response value that can be interpreted as a signal-likelihood. The higher the response value the higher the signal-likeness of an event. The network is created with the same set of parameters as the reconstruction networks (see Table 5.2) and the events used for the training are subsequently discarded. The signal events are scaled to match the number of expected background events in order to have two samples of the same size.

The different shapes of the discriminating variables for background and signal events can be seen

5. Search for tHq production at $\sqrt{s} = 8 \text{ TeV}$

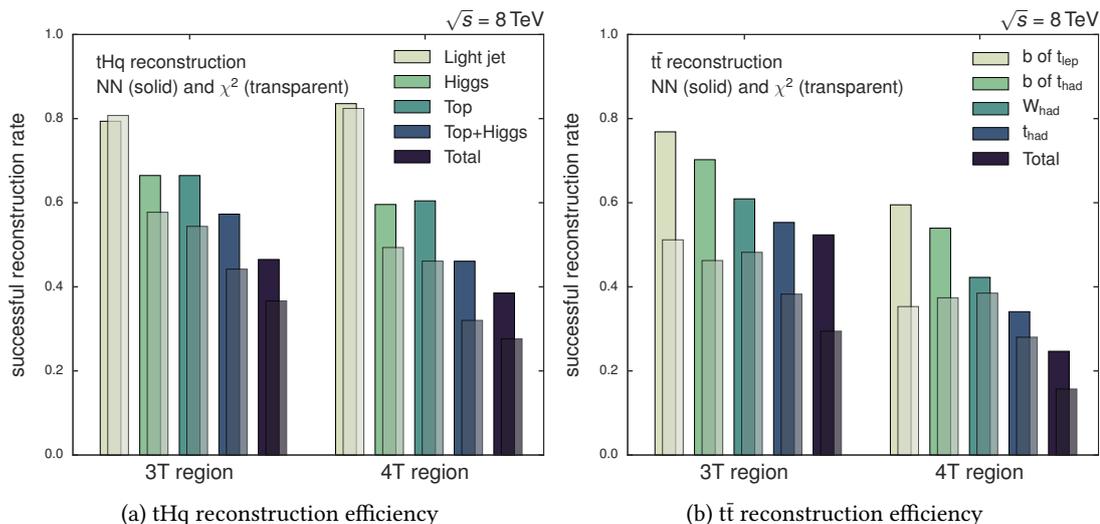


Figure 5.16.: Efficiencies of the two event reconstructions, one employing a neural network and the other using a simpler χ^2 approach, are shown in the 3T and 4T region. On the left the outcome for the tHq reconstruction and on the right for the $t\bar{t}$ reconstruction is shown. The bars show the rate a certain object or the complete event was reconstructed correctly. The solid bars represent the reconstruction of this analysis using neural networks for the assignment and the transparent bars represent the χ^2 reconstruction. The reconstruction using neural networks is clearly superior to the χ^2 method, as jets are assigned to their correct partons more frequently.

in Figure 5.17. In order to ensure that only well-modeled variables enter the neural network the variable shapes of recorded events are compared to those of simulated events. The comparisons in the $t\bar{t}$ control region and the two signal regions can be found in Figures 5.18, 5.19 and 5.20, respectively. The response of the neural network for the utilized training sample and the independent testing sample can be found in Figure 5.22 and the correlations of the used variables for signal and background events are visualized in Figure 5.21.

After the careful validation, every event is subjected to the classification neural network and is assigned a value according to its background-likeness or signal-likeness. A comparison of the classifier performance when used on simulation events and data in the $t\bar{t}$ control region can be found in Figure 5.23. The good agreement between the classifier output for simulation and data gives confidence that the procedure is working well.

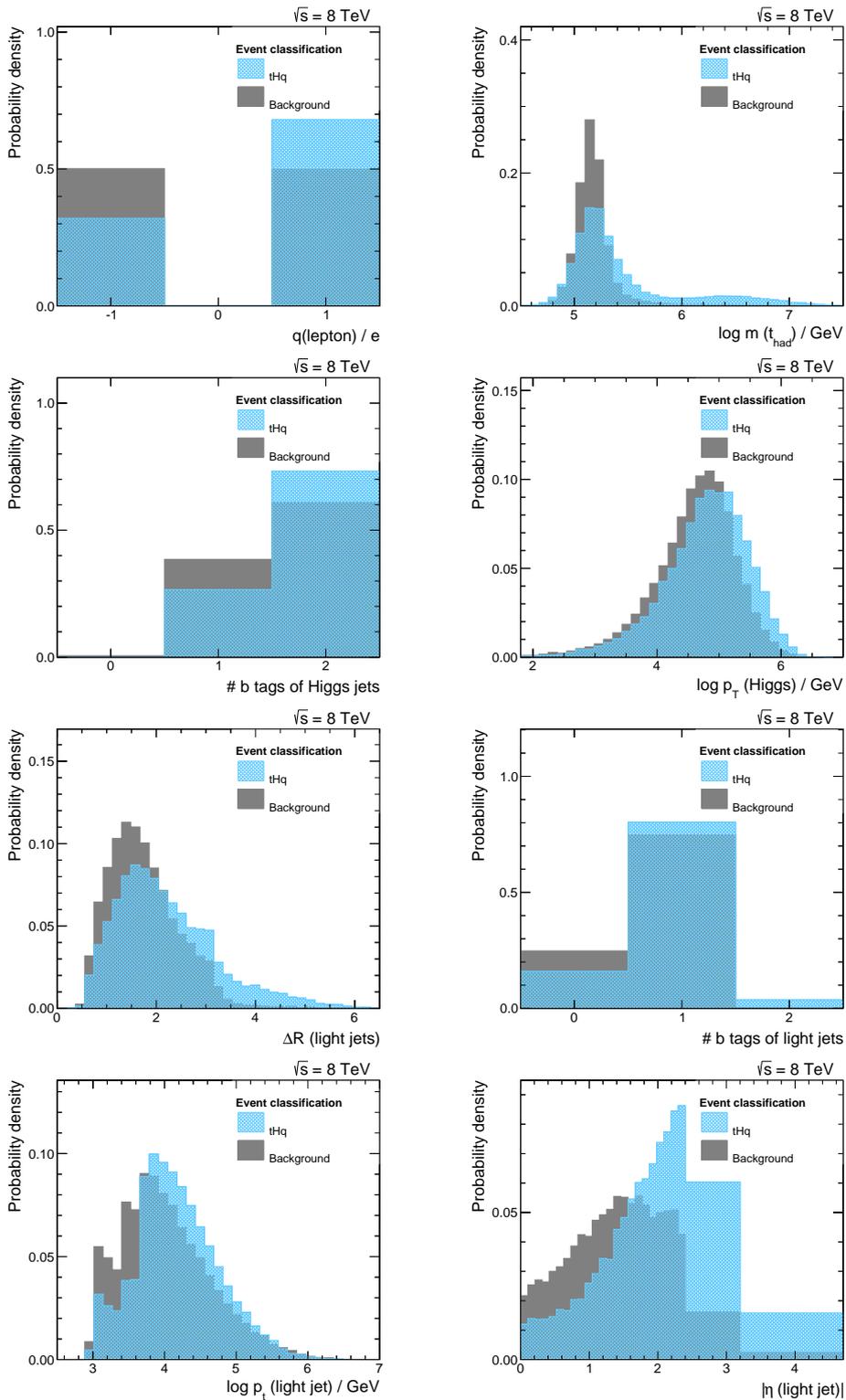


Figure 5.17.: Input Variables used in the final classification of events sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 5.5.

5. Search for tHq production at $\sqrt{s} = 8$ TeV

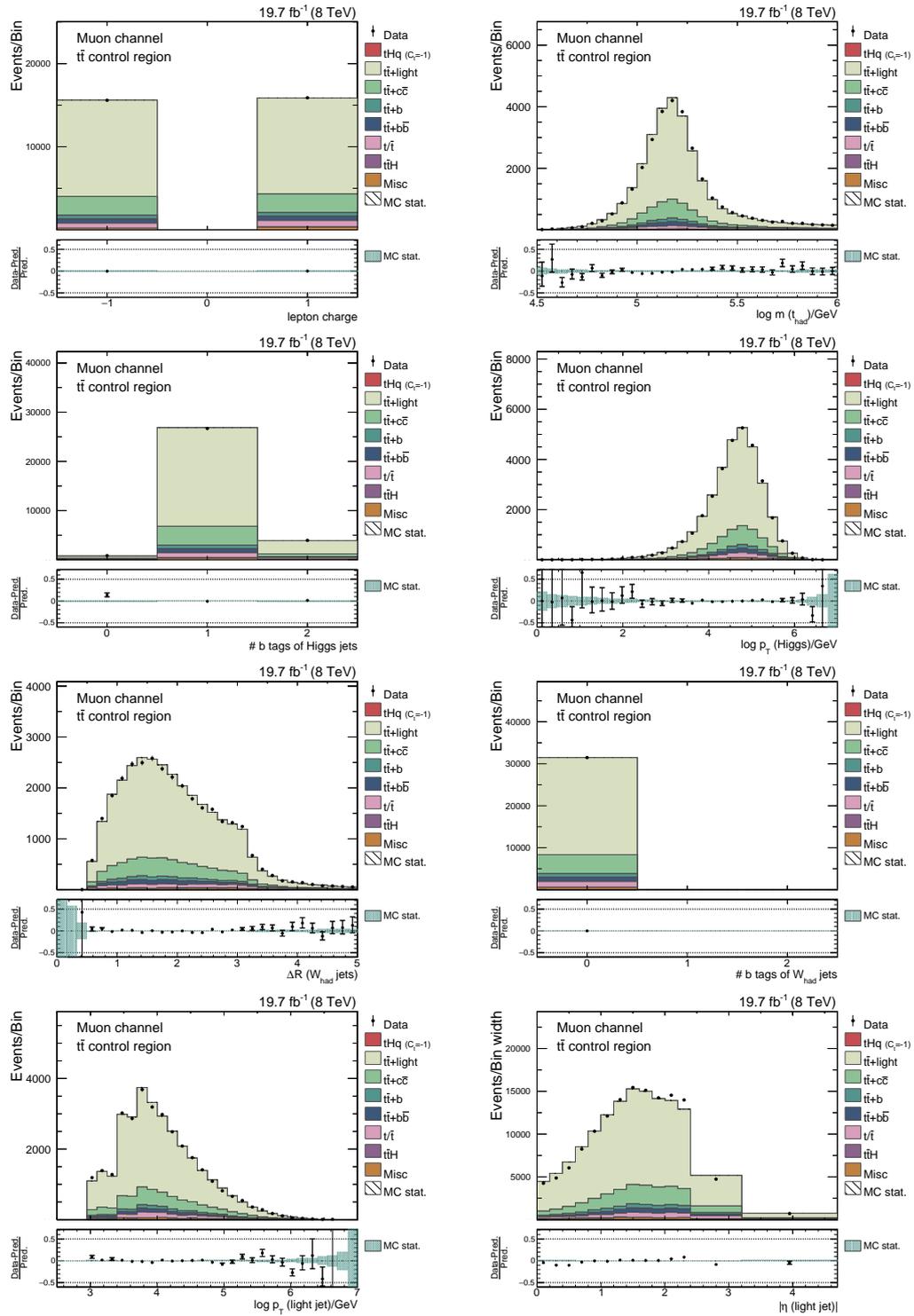


Figure 5.18.: Simulation to data comparisons for the input variables of the classification for the muon channel in the $t\bar{t}$ control region sorted by their importance in the training. A good agreement between simulation and data is observed. In both diagrams the simulation is scaled to match the event yields observed in data and all simulation corrections are applied.

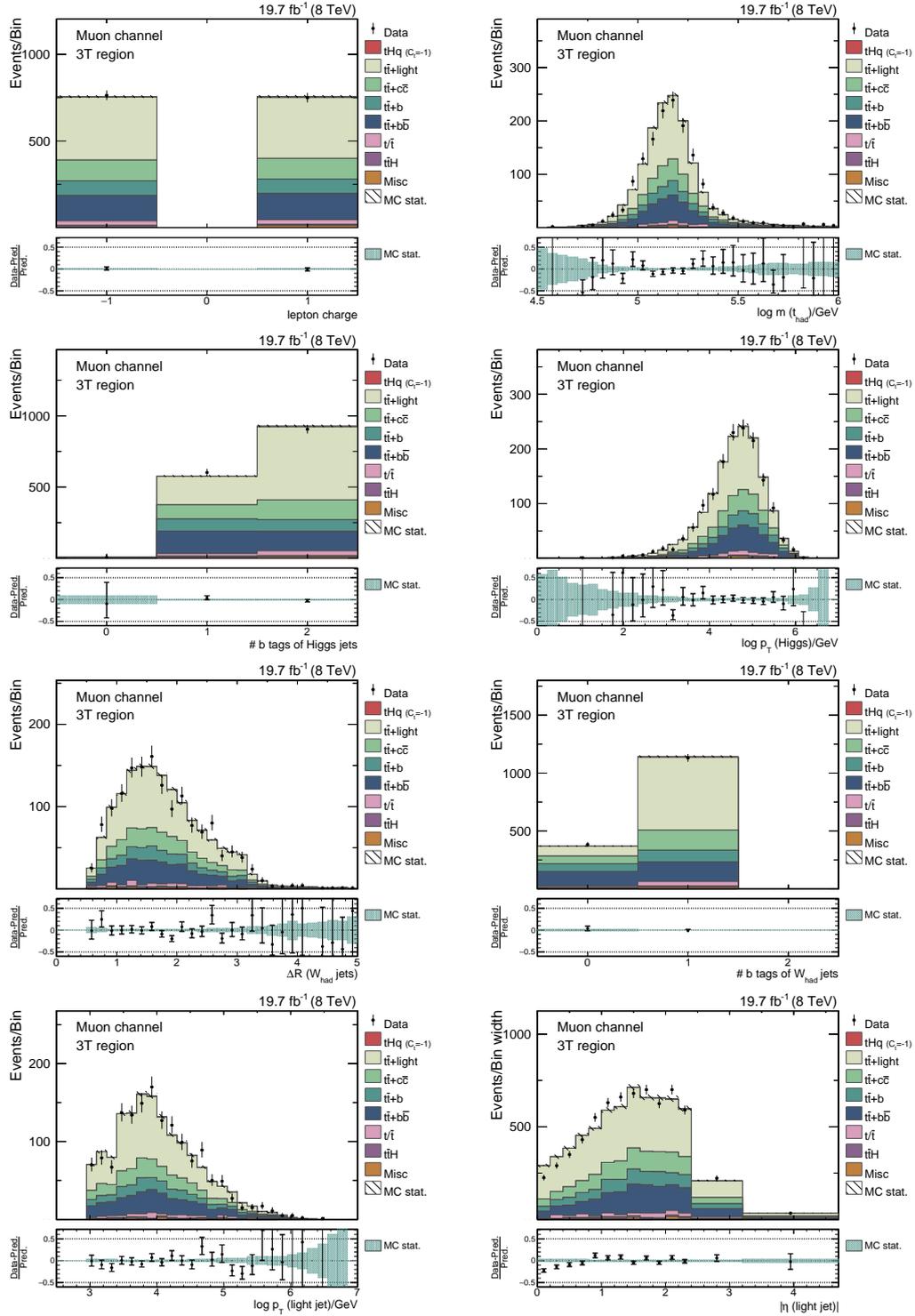


Figure 5.19.: Simulation to data comparisons for the input variables of the classification for the muon channel in the 3T region sorted by their importance in the training. A good agreement between simulation and data is observed. In both diagrams the simulation is scaled to match the event yields observed in data and all simulation corrections are applied.

5. Search for tHq production at $\sqrt{s} = 8$ TeV

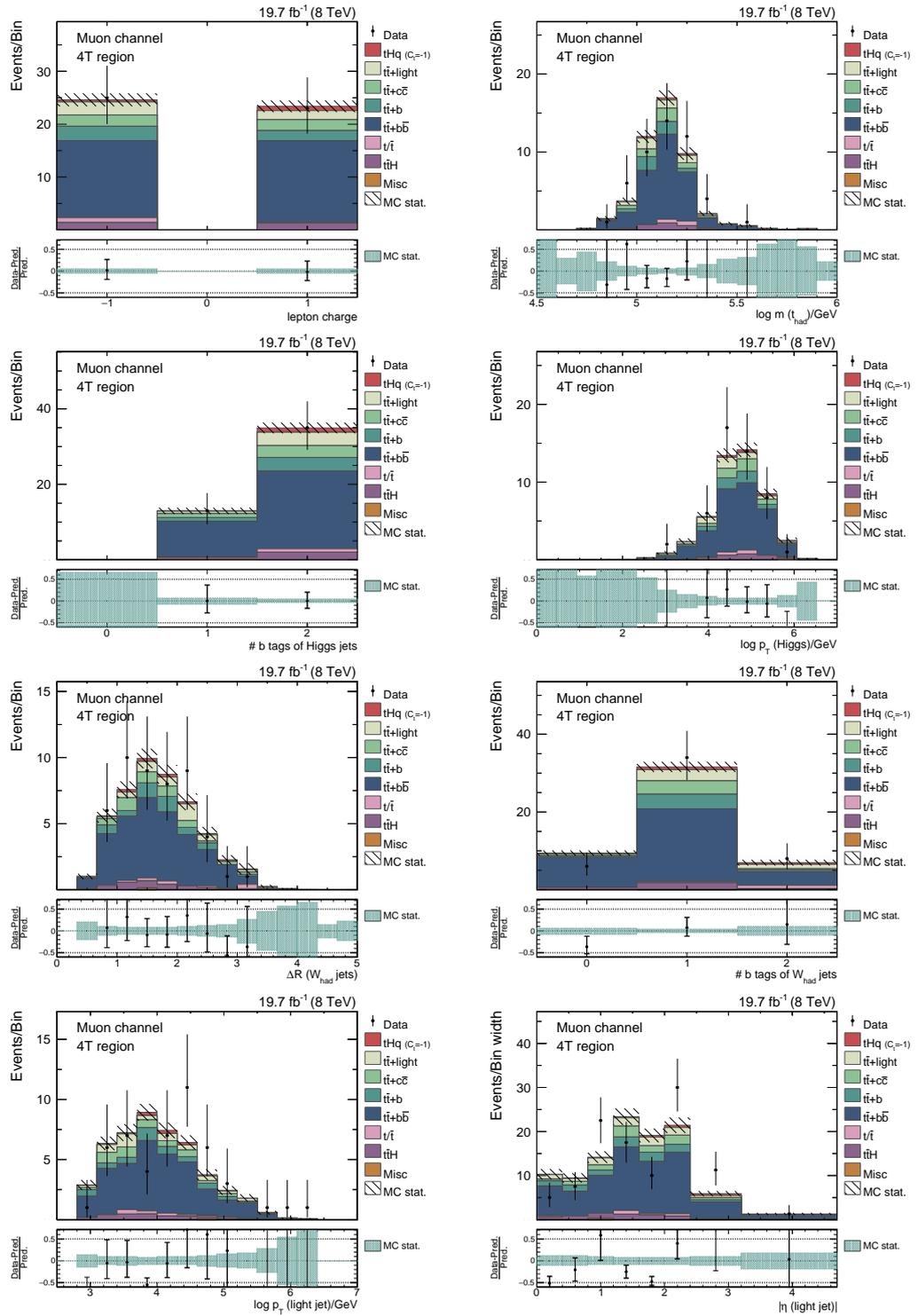


Figure 5.20.: Simulation to data comparisons for the input variables of the classification for the muon channel in the 4T region sorted by their importance in the training. A good agreement of simulation and data is observed. In both diagrams the simulation is scaled to match the event yields observed in data and all simulation corrections are applied.

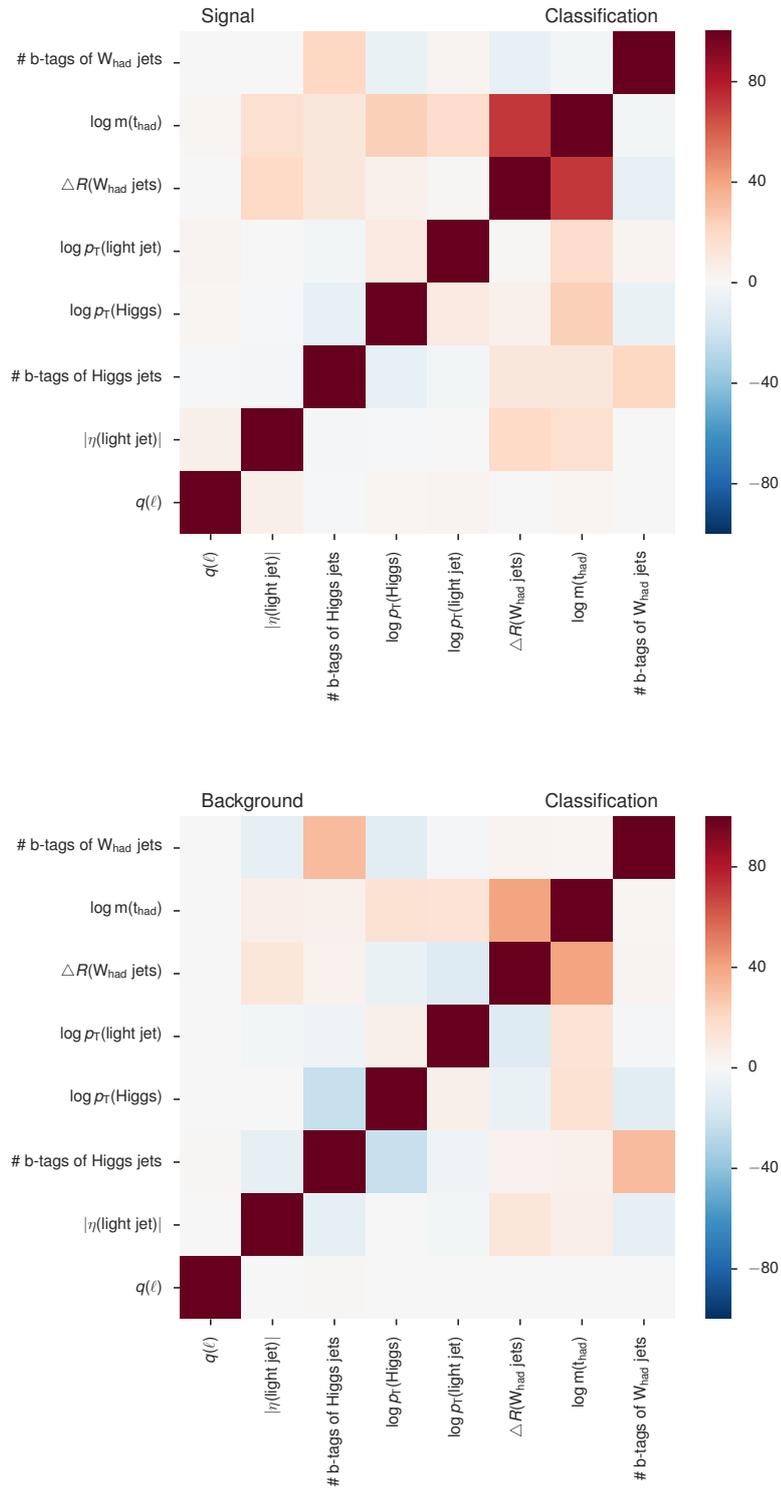


Figure 5.21.: Linear correlations between the variables used in the classification for signal and background events in the 3T region for the electron and muon channel combined.

Table 5.5.: Description of variables used in the classification and their importance ranking in the training. The variables are grouped into three categories: variables independent of any reconstruction (top), variables based on objects reconstructed under the $t\bar{t}$ hypothesis (center) and variables based on objects reconstructed under the tHq hypothesis (bottom). Instead of the transverse momenta the logarithm of these variables is used, as the neural net can process the information of narrower distributions better.

Rank	Variable	Description
1.	$q(\ell)$	electric charge of the lepton
2.	$\log m(t_{\text{had}})$	invariant mass of t_{had}
5.	$\Delta R (W_{\text{had}} \text{ jets})$	ΔR between the two light jets from the decay of W_{had}
6.	$\# \text{ b tags of } W_{\text{had}} \text{ jets}$	number of b-tagged jets of the two light jets from the t_{had} decay
3.	$\# \text{ b tags of Higgs jets}$	number of b-tagged jets of the Higgs boson decay products
4.	$\log p_{\text{T}} (\text{Higgs})$	transverse momentum of the Higgs boson
7.	$\log p_{\text{T}} (\text{light jet})$	transverse momentum of the light forward jet
8.	$ \eta(\text{light jet}) $	absolute pseudorapidity of light forward jet

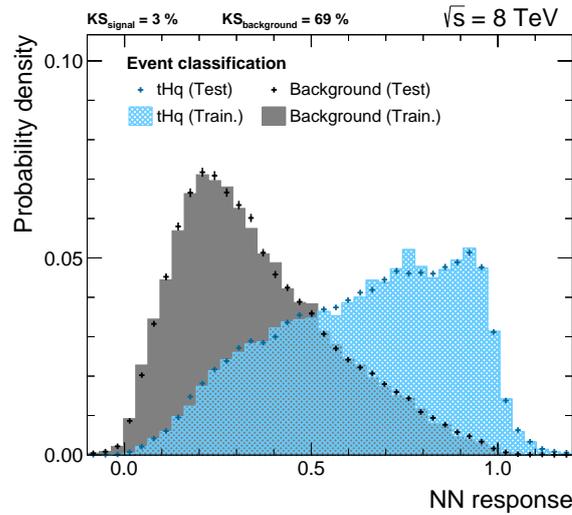


Figure 5.22.: Output values for the classification neural network for signal and background events in the 3T region for the electron and muon channel combined. A clear separation is visible. Events of the two training samples are depicted in the colored histograms. The training is examined with an independent set of events that were not part of the training sample. The response for this testing sample is shown as colored markers and a good agreement between the testing and training sample is seen, verified by high KS-values of the two distributions. No sign of overtraining can be found.

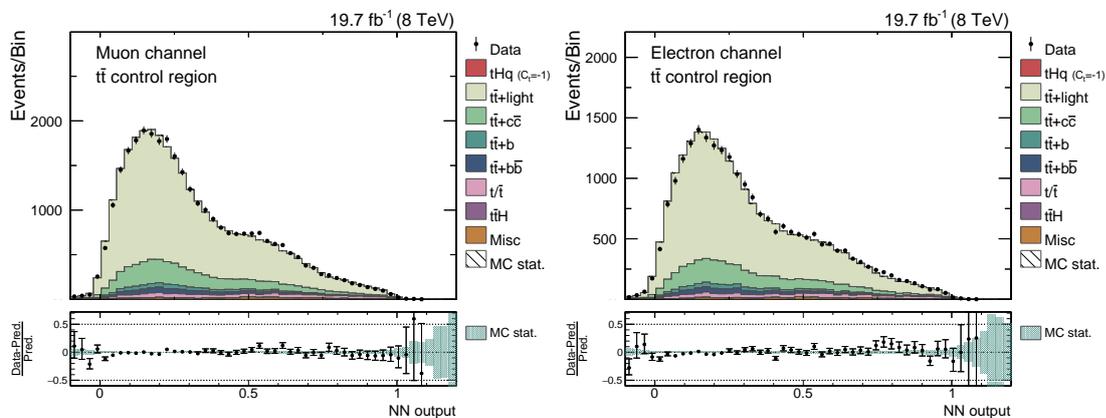


Figure 5.23.: Response of the classification comparing simulation to data in the $t\bar{t}$ control region for the muon channel (left) and the electron channel (right). In both distributions a good agreement between simulation and data is observed. In both diagrams the simulation is scaled to match the event yields observed in data and all simulation corrections are applied. The corresponding distributions for the 3T and 4T region before the fit can be found in Appendix A.8.

5.10. Systematic Uncertainties

Many different sources of uncertainties hamper measurements in high energy physics. The different uncertainties can be of statistical nature, as basically all processes are based on possibilities, or systematic nature. The systematic uncertainties can be further divided into theoretical uncertainties, such as an imperfect knowledge of the production cross sections for all processes, and experimental uncertainties, such as a limited energy resolution in the detector. All uncertainty sources contributing to this analysis are explained in the following section.

Systematic uncertainties can either occur as rate uncertainties, where a simple constant scaling factor for the affected process can be applied, or as shape uncertainties, which manipulate not only the normalization but also the shape of certain distributions. Systematic shape uncertainties can already have an impact on early stages of the analysis. A variation in e. g. the jet energy scale affects the events selected for the analyzed regions. Hence, these systematic shape uncertainties are applied to the concerned simulation samples and the complete analysis chain is reiterated with these systematically varied samples.

Experimental Uncertainties

Luminosity (rate) The estimation of the luminosity delivered by the LHC relies on experimental conditions, such as pileup and the number of protons per bunch, and the visible cross section. As these parameters are varying during the LHC run time a perfect measurement of the luminosity is impossible. The suggested uncertainty applied to the luminosity measurement of 2.6% is applied [162].

Lepton efficiencies (rate) In order to account for a disagreement between the recorded lepton efficiency and the simulated efficiency a total uncertainty of 2% is applied, independently for muons and electrons in the respective channels.

b-tagging efficiencies (shape) The corrections of the simulated b-tagging efficiencies described in Section 5.6.4 need to be varied based on the error-prone b-tagging efficiency measurements. The applied scale factors for b and c quarks are simultaneously varied within their uncertainties. These variations effectively change the values of the CSV working point resulting in different event yields. This yield variation is realized by applying a corresponding set of weights to the samples.

Independently, the scale factors for light quarks are also varied within their uncertainties, leading to an effective change of the rate of erroneously b-tagged light jets. The effect of the b-tagging uncertainty on the final NN output can be seen in Figure 5.24.

Pileup (shape) The reweighting procedure applied to simulation samples to reproduce the number of primary vertices measured in real collisions (see Section 5.6.1) also needs to be considered as an uncertainty source. The event weights calculated during this procedure are redetermined by reweighting to systematically shifted versions of the measured pileup distributions. The detailed procedure can be found in Reference [163]. The assumed variations correspond to a $\pm 6\%$ variation of the nominal total inelastic cross section of $\sigma_{pp} = 69.4 \text{ mb}$.

Missing transverse energy (shape) The energy contribution from jets with $p_T < 10$ GeV and from particle-flow candidates not clustered into jets is called *unclustered energy*. The contribution of unclustered energy to \cancel{E}_T is varied by $\pm 10\%$. Due to the recalculated \cancel{E}_T the number of events passing the selection criteria of Section 5.7 varies. This change in event numbers is again propagated to the limit calculation.

Jet energy resolution (shape) The jet energy resolution in recorded data is worse than what is simulated. This is corrected by smearing the jet energy resolution in simulation according to a recipe provided by the JetMET [164] group of CMS. These correction factors are varied within their uncertainties and the resulting samples are used to estimate the shape variations of the final classification output caused by the jet energy resolution smearing.

Jet energy scale (shape) As described in Section 3.3.6 different jet energy corrections are applied to data and simulation. A set of 16 independent systematic uncertainties each covering a different aspect of the corrections is introduced [165]. The systematic effects are then evaluated by varying the jet energies within their respective flavor-dependent uncertainties [166]. This leads to different jet p_T spectra, which already affect the analysis at the earliest stages. Simulation samples are reproduced with varied jet energy corrections and the complete analysis chain is reiterated. The effect of the JES uncertainty on the final NN output can be seen in Figure 5.24.

Theoretical Uncertainties

Cross section (PDFs/ Q^2 scale) (rate) The predicted cross sections for the considered processes are dependent on the choice of a set of parton distribution functions and the QCD scale. As shown in Table 5.6 different PDF uncertainties are applied for different processes, while treating uncertainties for the same production mechanism as fully correlated. A variation of the fixed Q^2 scale in the generation leads to a change of normalization and shape. A set of Q^2 rate systematic uncertainties is applied to cover the normalization differences. For the signal process and the dominant $t\bar{t}$ background systematically shifted samples with different utilized Q^2 scales are available. The effect of the shape variations are described at a later point.

$t\bar{t}$ +heavy flavor cross section (rate) The splitting of the $t\bar{t}$ simulation sample into sub-samples based on their heavy-flavor content (see Section 5.4.1) allows for an independent scaling of these templates. As no measurement of the normalization of the different templates in control regions has achieved a better precision than 50%, this conservative number is chosen as rate uncertainties for the $t\bar{t}+b$, $t\bar{t}+b\bar{b}$ and $t\bar{t}+c\bar{c}$ samples. These rate uncertainties effectively pose as scale factors that are determined simultaneously as all other background normalizations. By treating them independently the fit adjusts the modeled ratios of heavy-flavor content to light-flavor content to what is observed in data.

Top quark p_T reweighting (shape) When comparing the p_T spectrum of top quarks in a simulated $t\bar{t}$ sample to the measured spectrum, a softer spectrum in data is observed. The correction (see Section 5.6.3) introduces a new systematic uncertainty source. The up and

5. Search for tHq production at $\sqrt{s} = 8 \text{ TeV}$

Table 5.6.: Cross section uncertainties separated into PDF and QCD scale applied for the different processes.

Process	PDF (%)			Q^2 scale (%)
	gg	$q\bar{q}$	qg	
tHq			2.0	
t \bar{t} H	9.0			12.9
t \bar{t}	2.6			3.0
Single top			4.6	2.0
W+jets		4.8		1.3
Z+jets		4.2		1.2
Diboson				3.5

down varied t \bar{t} samples are produced by applying the correction factor twice per event or not at all, respectively. The effect of the top quark p_T reweighting on the final NN output can be found in Figure 5.24.

Q^2 scale (shape) For the signal process and the important t \bar{t} background process samples with up- and down-variations of the Q^2 scale have been produced. By repeating the complete analysis for these systematically shifted samples the effect of the Q^2 systematic uncertainty on the NN output is evaluated.

However, as these samples contain significantly less simulated events than the nominal samples, the observed event numbers in the 4T region are not sufficient for a qualified statement about shape variations. Therefore the samples in the 4T region are used to determine an additional rate systematic uncertainty for this region when comparing it to the nominal samples. The shape is considered identical to the shape of the nominal samples. The effect of the Q^2 scale systematic uncertainty on the final NN output can be seen in Figure 5.24.

Matching threshold (shape) In the generation of the t \bar{t} simulation samples a matching threshold of 40 GeV has been set, dividing the phase space into a region where emissions are generated by the matrix element generator or simulated by the parton shower. Two additional t \bar{t} simulation samples are produced by setting the matching threshold to 60 GeV and to 30 GeV, respectively. These varied samples are used to evaluate the effect of the matching threshold choice on the final NN output.

Statistical Uncertainties

Bin-by-bin uncertainties (shape) Due to the finite size of the used simulated samples an additional uncertainty has to be considered. This uncertainty is evaluated applying the “Barlow-Beeston lite” method [167, 168], which introduces a nuisance parameter in the fit for

Table 5.7.: Postfit yields in the four signal channels that are fitted simultaneously. The uncertainties include systematic and statistical uncertainties. Additionally, the sum of all expected background events is quoted and the number of observed events. The contribution from the W/Z+jets and diboson backgrounds are grouped under the label "Misc".

	3T		4T	
	Electron	Muon	Electron	Muon
$t\bar{t}$ +light	421 \pm 46	645 \pm 69	2.1 \pm 0.7	2.3 \pm 0.7
$t\bar{t}$ + $c\bar{c}$	150 \pm 67	223 \pm 106	2.7 \pm 3.7	3.6 \pm 1.5
$t\bar{t}$ +b	152 \pm 63	199 \pm 78	2.5 \pm 1.4	3.4 \pm 1.2
$t\bar{t}$ + $b\bar{b}$	245 \pm 46	346 \pm 62	22.9 \pm 3.8	33.1 \pm 6.2
Single Top	29.1 \pm 3.1	44.2 \pm 4.5	1.4 \pm 0.3	0.8 \pm 0.1
$t\bar{t}$ H	10.2 \pm 1.6	13.8 \pm 2.1	1.6 \pm 0.3	1.9 \pm 0.3
Misc	4.4 \pm 0.6	12.1 \pm 1.9	0.9 \pm 0.9	0.1 \pm 0.1
Σ Backgrounds	1012 \pm 113	1483 \pm 161	34.1 \pm 5.6	45.2 \pm 6.5
tHq	15.0 \pm 11.3	21.6 \pm 16.1	1.7 \pm 1.3	2.4 \pm 1.8
Observed	1028	1514	32	48

each bin in each sample and each region. In every introduced nuisance parameter one bin is varied within its uncertainties, while other bins are scaled such that the normalization of the complete distribution remains constant. This procedure introduces hundreds of additional nuisance parameters resulting in a significant increase of needed computing power. In order to keep the calculation complexity within reasonable bounds every bin uncertainty with a relative uncertainty below 5% is abolished.

5.11. Results

5.11.1. Fit of Final Discriminator

The statistical analysis in this thesis is performed with the `COMBINE` package. By performing a simultaneous fit in the two signal regions and for both lepton channels in the neural network output distribution an upper limit on the anomalous production of tHq with an assumed $C_t = -1$ is derived.

The provided limits are full CL_S limits at a confidence level of 95%. The event yields predicted by the simulation after the fit are provided in Table 6.8. The best agreement is found with a signal-strength factor of $\mu = 2.7^{+2.1}_{-2.0}$. The four NN distributions with full systematic and statistical uncertainty after the fit can be found in Figure 5.25.

5. Search for tHq production at $\sqrt{s} = 8$ TeV

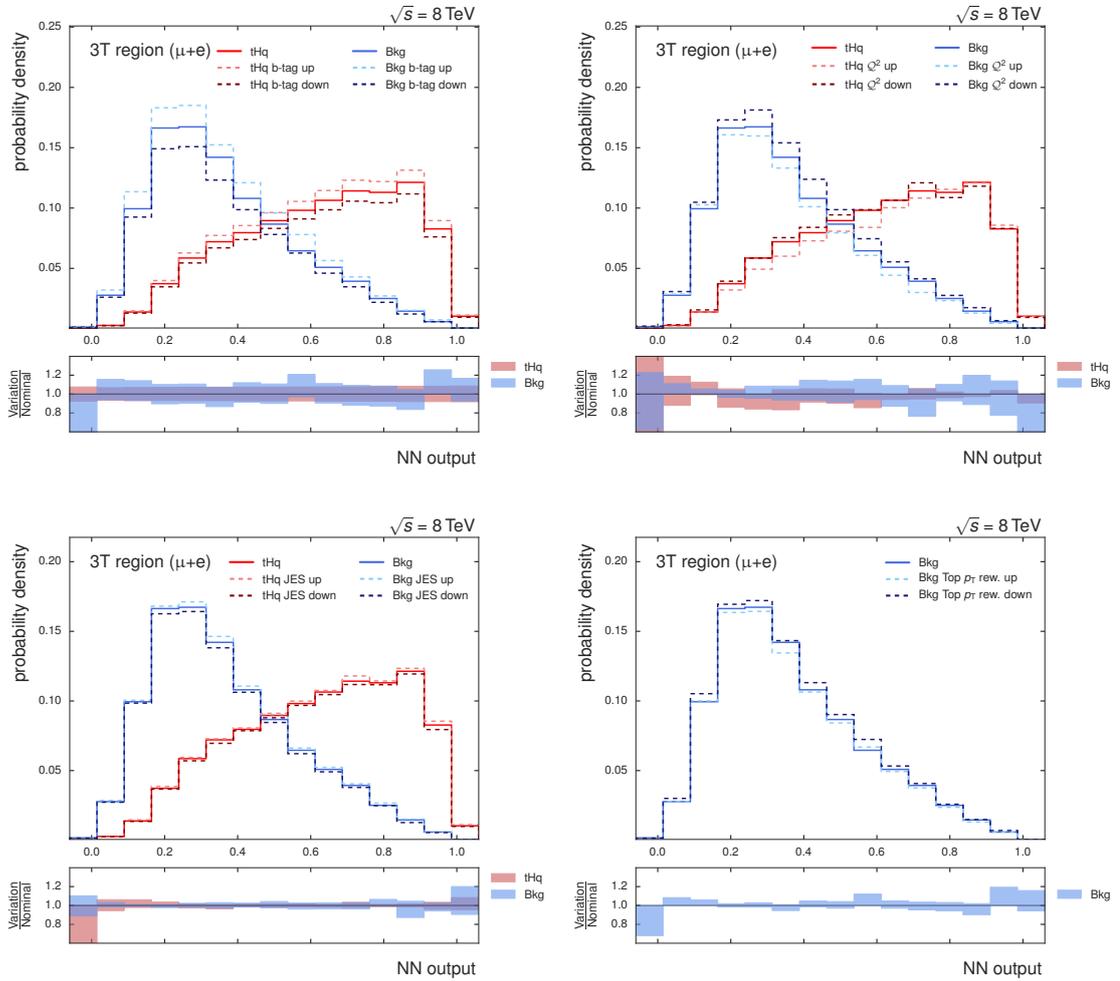


Figure 5.24.: Effect of different systematic shape uncertainties on the final neural network classifier separated into the background shape and signal shape in the 3T region. The diagrams show the effect of the b-tagging efficiency (top left), the Q^2 scale (top right), the combined jet energy scale (bottom left) and top quark p_T reweighting (bottom right) uncertainties. The effect of the 16 independent JES sources is combined by showing the envelope of all systematic variations.

The top quark p_T reweighting is only applied for the $t\bar{t}$ background, therefore only the background shape is shown.

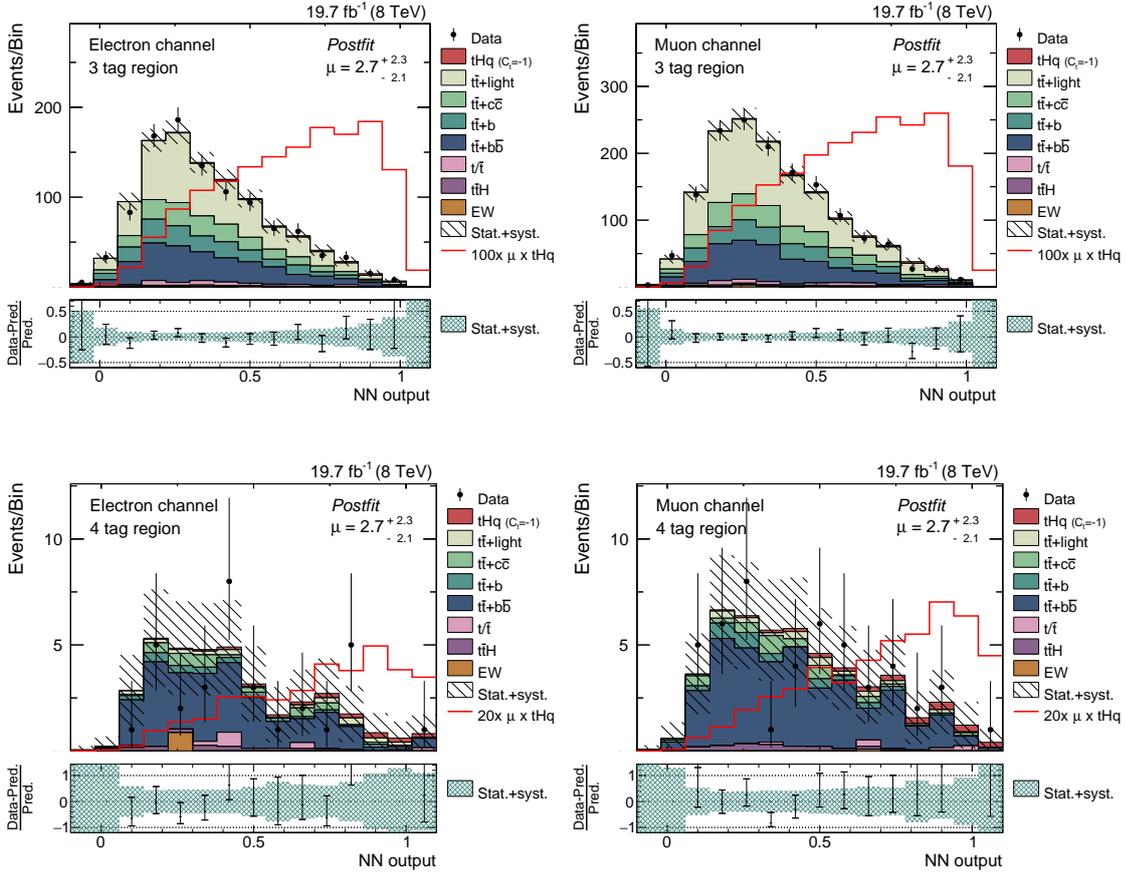


Figure 5.25.: Postfit distributions of the classifier output in both lepton channels and both signal regions. Uncertainty bands include statistical and systematic uncertainties. A good agreement in all four channels is observed, deviations from the prediction are covered by the uncertainties. The signal template is also scaled to its best-fit value of $\mu = 2.7$. The corresponding prefit distributions can be found in Appendix A.8.

5. Search for tHq production at $\sqrt{s} = 8$ TeV

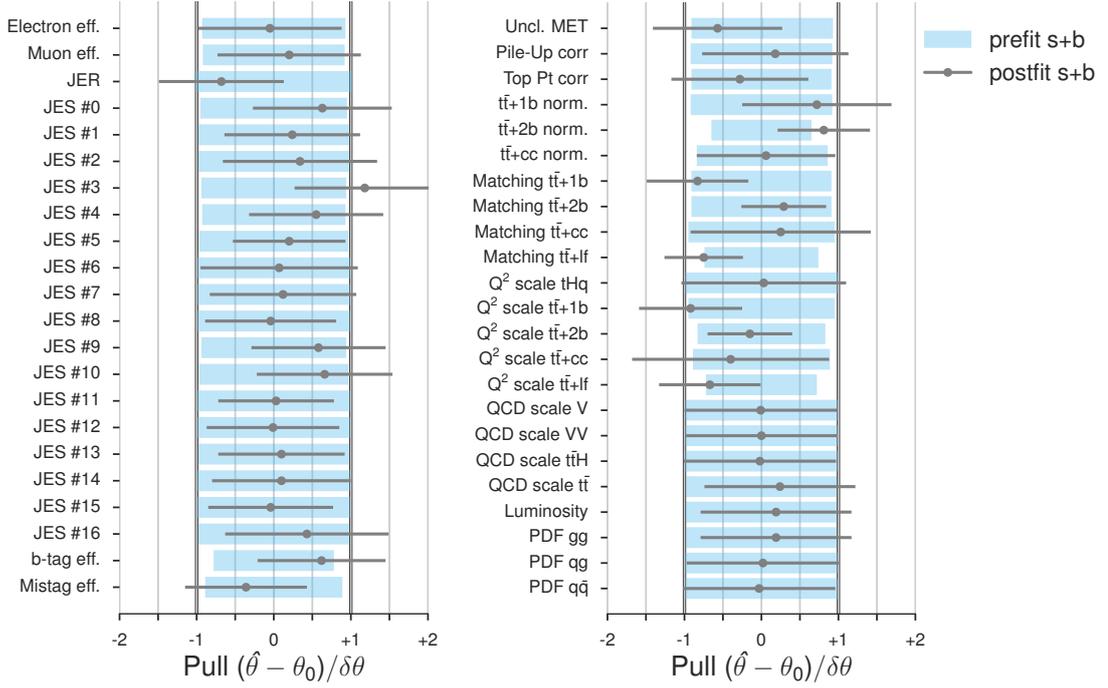


Figure 5.26.: Pre- and postfit pulls of all nuisance parameters of the analysis (bin-by-bin uncertainties excluded).

5.11.2. Analysis of Nuisance Parameters

In order to ensure a reasonable behavior of the uncertainty sources during the fitting procedure, their deviations from their initial values and the corresponding uncertainties are studied. Before the fit the nuisance parameters are all centered around their starting value θ_0 , but the fit determines the best value $\hat{\theta}$ for each parameter, hence “pulling” the mean of the likelihood away from its initial starting point. The interplay of different nuisance parameters can also lead to a constraint of the parameter uncertainty, as the fit prohibits that the nuisance parameter can exploit the complete possible uncertainty range. An unreasonable behavior of nuisance parameters would express itself in a strong pull away from the initial value or a strong constraint. The pulls of the considered nuisance parameters of the analysis but the hundreds of bin-by-bin uncertainties are shown in Figure 5.26. It is visible that only one of the JES source nuisance parameters is pulled outside its one sigma boundary and no parameter is severely constrained. Of special interest are also the nuisance parameters corresponding to the normalization of the $t\bar{t}$ events with heavy flavor content. Whereas the $t\bar{t}+c\bar{c}$ parameter is pulled only slightly from zero to 0.06, implying good modeling of the ratio of $t\bar{t}+c\bar{c}$ to $t\bar{t}$ +light, the $t\bar{t}+b$ and $t\bar{t}+b\bar{b}$ components are pulled to +0.72 and +0.81, respectively. This indicates that the assumed ratios of $t\bar{t}$ +heavy to $t\bar{t}$ +light is indeed underestimated in the simulation.

The different uncertainties introduced in Section 5.10 influence the final limit in different ways.

In order to get a grasp on the impact of the different uncertainty sources the effect of every single source on the final limit was evaluated. The following study was solely performed with asymptotic CL_S limits as the calculation of asymptotic CL_S limits needs by far less computing resources than full CL_S limits.

By performing a profile likelihood fit the optimal value $\hat{\theta}$ for each nuisance parameter is obtained. The impact of one parameter p can be calculated in two ways: In the first way all other parameters but p are fixed to their corresponding $\hat{\theta}$ values, leaving only p free to float within its uncertainties. If the limit is calculated with this nuisance set, the result reflects what the limit would be like, if p was the only uncertainty source, but conserving the information about the other parameters gained in the previous fit.

Secondly, by fixing only the parameter p to $\hat{\theta}$ and allowing all other nuisances to float, the effect of the removal of parameter p from the set of uncertainties can be evaluated. The percentage changes of the limits when taking an uncertainty as only source or when removing it from the set can be found in Figure 5.27. The nuisances are sorted by their effect after removal in descending order, therefore making visible where possible future improvements could be found. The study shows that the largest impact on the limit stems from the jet energy scale uncertainties and the Q^2 scale uncertainties. Whereas the effect of the Q^2 scale is the largest when taking it as only uncertainty its effect gets mitigated when it is removed from the whole set, as other uncertainties offset the systematic variations. The study also shows that the analysis is not clearly limited by either, experimental or theoretical uncertainties.

5.11.3. CL_S Limit

Based on the results of the fit upper CL_S limits at 95% confidence level are calculated for the anomalous tHq production. The limits are calculated in a fully frequentist way. A total of 64,000 toys have been produced at 32 different signal strength points in order to derive the final expected and observed limit of this analysis. The exact limit for the 3T and 4T region as well as the combined result can be found in Table 5.9 and a visualization can be found in Figure 5.28. This analysis is able to exclude the production of tHq with a cross section of $7.5 \times \sigma_{C_t=-1} = 1755$ fb. This observation corresponds to an excess slightly over 1 sigma above the expectations. The limit is driven by the 3T bin with an expected limit of 6.2 and an observed limit of 7.0. The 4T region has an expected limit 9.8 and an observed limit of 20.0 which corresponds to an excess slightly above 2 standard deviations.

Out of the two lepton channels the muon channel is the more sensitive one with an expected limit of 6.8 and an observed limit of 8.7. In solely the electron channel an expected limit of 8.0 and an observed limit of 11.7 have been found. It is apparent that the electron channel in the 4T region is the main culprit for the excess seen in the whole 4T region.

5.12. Combination with Other Decay Channels

This analysis was part of a combination effort of all in the CMS collaboration performed searches for the tHq production with an anomalous top quark Yukawa coupling of $C_t = -1$. This combi-

5. Search for tHq production at $\sqrt{s} = 8 \text{ TeV}$

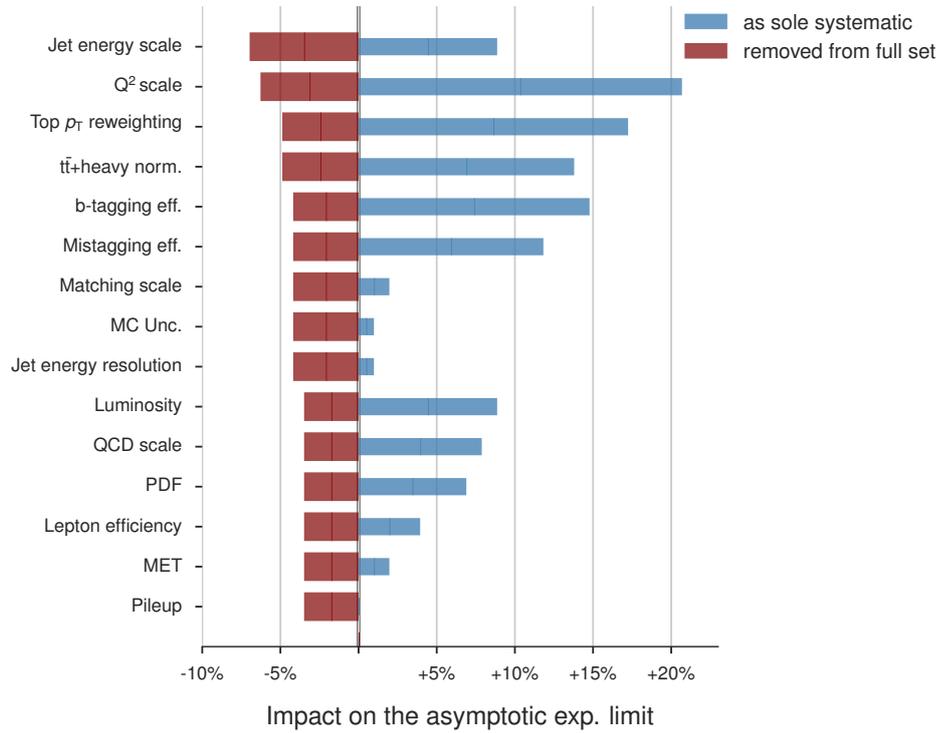


Figure 5.27.: Impact of groups of systematic uncertainties on the expected asymptotic limit. The groups of systematic uncertainties are either removed from the fit by fixing them to their postfit value, or used as single systematic by fixing all other uncertainties to their postfit values. The relative changes displayed in this diagram are calculated to the limit with all systematic uncertainties included (red bars) and to the limit, where all uncertainties are frozen to their best fit value (blue bars).

Table 5.8.: Expected and observed CL_S limits at 95% C.L. in the 3T and 4T region and their combination. Also the 68% and 95% uncertainty band values are shown. The observed limit of 7.5 corresponds to a 1.25σ upwards fluctuation. A graphical representation of the limits on the tHq production can be found in Figure 5.28.

Region	Observed Limit	Expected Limit		
		Median	$\pm 1\sigma$	$\pm 2\sigma$
3T	7.0	6.2	[4.5 , 8.8]	[3.2 , 12.0]
4T	20.0	9.8	[7.3 , 13.9]	[6.0 , 19.4]
Combination	7.5	5.0	[3.6 , 6.9]	[3.2 , 8.1]

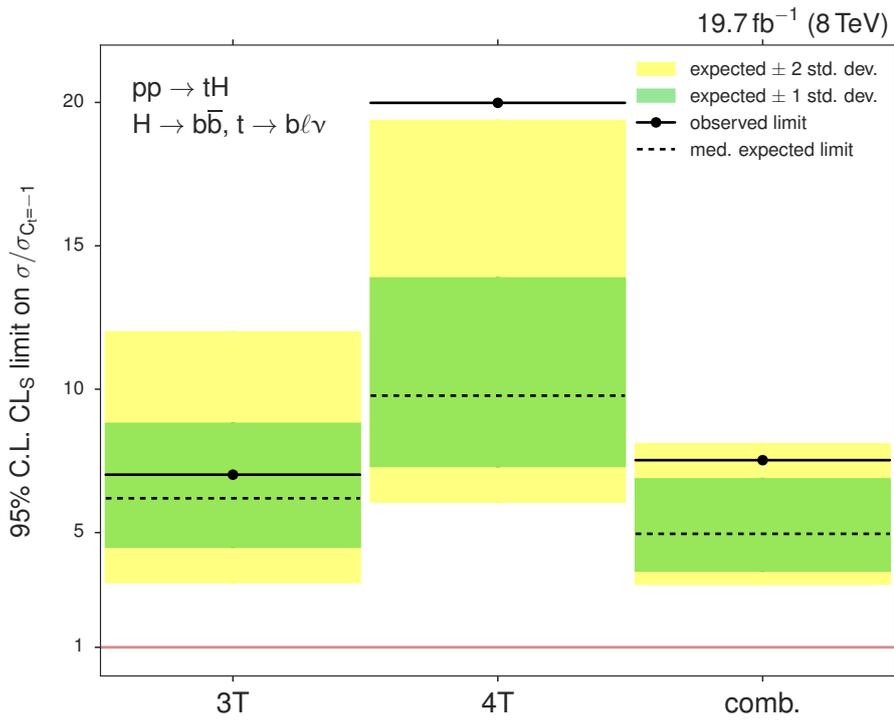


Figure 5.28.: Expected and observed CL_s limits at 95% C.L. in the 3T and 4T as well as their combination. Corresponding numbers can be found in Table 5.9.

nation was published by CMS [138] and is accepted by the journal of high energy physics for publication. The coherent treatment of systematic uncertainties spanning all involved analyses has been a major effort, partly by myself, but as a detailed description would be out of scope for this thesis the combination is only summarized in the following chapter.

The combined analyses differ in the studied Higgs boson decay products:

$H \rightarrow \gamma\gamma$ Analogous to the pioneering $H \rightarrow \gamma\gamma$ analysis for the Higgs boson discovery, the decay channel of the Higgs boson into two photons is the most sensitive analysis in the search for tHq production [169]. Although the branching ratio is very small (0.23%), the diphoton final state allows for a very good background rejection.

The decay of the Higgs boson into two photons is propagated by a loop of virtual particles, with a dominating contribution from the heaviest particles i. e. the top quark. Hence, a change in the Yukawa coupling of the top quark would also lead to an enhancement in the $H \rightarrow \gamma\gamma$ branching ratio. An increase in the $H \rightarrow \gamma\gamma$ branching ratio would also be visible in all $H \rightarrow \gamma\gamma$ analyses, but is not observed as of this writing. However, mechanisms of new physics could be able to dampen the effect of $C_t = -1$ on the branching ratio, explaining the currently measured values of the $H \rightarrow \gamma\gamma$ branching ratio.

To account for this twofold enhancement, in production and decay, the branching ratio is evacuated and treated as an independent parameter in the production, as can be seen later in the discussion of the combination.

The analysis itself is able to suppress background contributions in the signal region to a minimum by requiring two photons, one isolated lepton (muon or electron) and at least one b -tagged jet. The remaining $t\bar{t}H$ background is suppressed by employing a strict requirement on the output of a Bayesian classifier. After these selections the analysis expects 0.04 ± 0.05 background events and 0.67 signal events under the assumption of $C_t = -1$. Yet no recorded event passed the selection criteria.

Multileptons By exploiting the second highest branching ratio of the Higgs boson the $H \rightarrow WW$ analysis is searching for events with multiple leptons in the final state [170]. The three W bosons in the final state offer plenty of different possible combinations containing multiple leptons. The employed selection criteria lead to an acceptance of the $H \rightarrow \tau\tau$ decay, where one or both of the tau leptons decay leptonically.

In order to extract the limit a maximum likelihood fit of a classifier output is performed in all three channels.

$H \rightarrow \tau_{\text{had}}\tau_{\text{lep}}$ The $H \rightarrow \tau_{\text{had}}\tau_{\text{lep}}$ analysis is performed in two separate channels, each with three reconstructed leptons, the $e\mu\tau_{\text{had}}$ and $\mu\mu\tau_{\text{had}}$ channels. The limit is extracted by performing a combined maximum likelihood fit of the Fisher discriminant distribution in the two categories. This analysis is the least sensitive in the search for tHq production.

For the combination a meticulous analysis of the systematic uncertainties has been performed. Uncertainty sources which affect all analyses in the same manner, such as the uncertainty on the luminosity measurement, can be correlated and therefore the number of free parameters in the final fit can be reduced. Other uncertainties that influence only a single analysis or uncertainties that influence analyses in different ways are treated as uncorrelated. A detailed overview of

Table 5.9.: Expected and observed CL_S limits at 95% C.L. for the individual analyses that were combined with the analysis described in this chapter, and their combination. Additionally the 68% and 95% uncertainty band values are provided. No excesses are seen in either of the analyses and the observed limit of the combination of 2.8 corresponds to an one σ upwards fluctuation. The limit for the $H \rightarrow \gamma\gamma$ decay channel is calculated under the assumption that $C_t = -1$ also increases the $H \rightarrow \gamma\gamma$ branching ratio by a factor of 2.4. The limit as a function of the $H \rightarrow \gamma\gamma$ branching ratio can be found in Figure 5.29. Differences of the limits of the $b\bar{b}$ channel in the combination and in this thesis are owed to a more conservative systematic uncertainty treatment and small changes in the limit calculation procedure in the combination.

Channel	Observed Limit	Expected Limit		
		Median	$\pm 1\sigma$	$\pm 2\sigma$
$\gamma\gamma$	4.1	4.1	[3.7 , 4.2]	[3.4 , 5.3]
$b\bar{b}$	7.6	5.4	[3.8, 7.7]	[2.8, 10.7]
Multilepton	6.7	5.0	[3.6 , 7.1]	[2.9 , 10.3]
$\tau\tau$	9.8	11.4	[8.1 , 16.7]	[6.0 , 24.9]
Combination	2.8	2.0	[1.6 , 2.8]	[1.2 , 4.1]

the treatment of the systematic uncertainties can be found in Reference [138]. The limits have been derived utilizing the frequentist CL_S method.

As mentioned earlier, an anomalous $C_t = -1$ coupling would also increase the branching ratio of $H \rightarrow \gamma\gamma$, what aggravates the representation of the limit on the tHq production alone. Therefore, the limit on the production cross section of tHq is calculated for different branching ratios of $H \rightarrow \gamma\gamma$. The results of the combination for a modified $H \rightarrow \gamma\gamma$ branching ratio as well as the limits of the individual analyses can be found in Table 5.9. The median expected upper limit at 95% confidence level is 2.0 and the observed upper limit lies at 2.8, what corresponds to an upwards fluctuation of one standard deviation.

The effect of the assumption of different $H \rightarrow \gamma\gamma$ branching ratios can be found in Figure 5.29. No significant excess above the expectation can be found in the combination.

5.13. Summary

The search for associated production of a Higgs boson and a single top quark with the Higgs boson decaying into a pair of bottom quarks presented in this chapter is the first direct search for this channel.

The analysis described in this chapter of this thesis employed the full 2012 dataset recorded by the CMS detector, corresponding to 19.7 fb^{-1} of recorded data. The search is optimized under the assumption of an anomalous coupling of the Higgs boson to the top quark, represented by

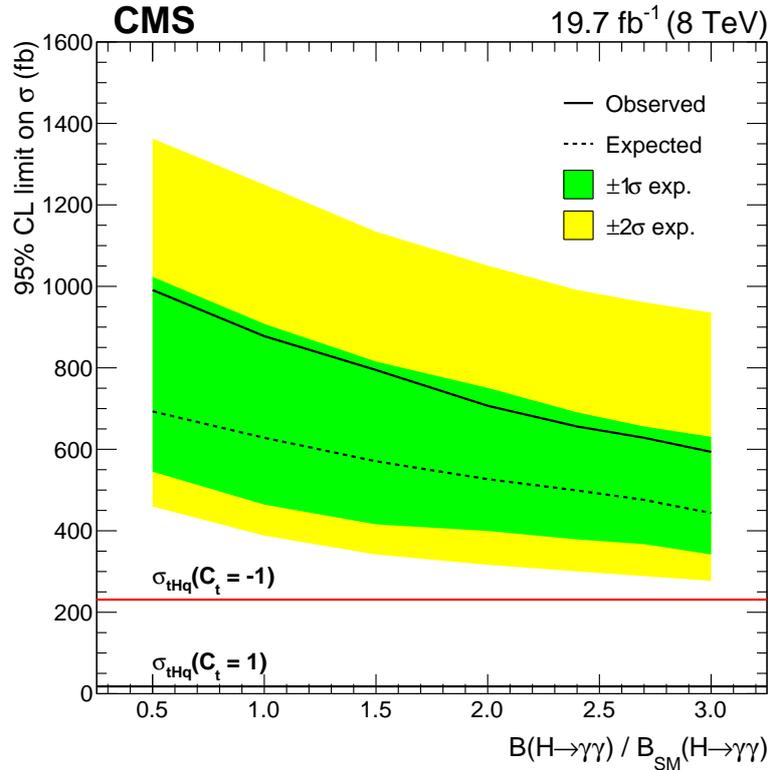


Figure 5.29.: Expected and observed CL_S limits for the tHq production at 95% C.L. provided by the combination as a function of the assumed branching ratio of $H \rightarrow \gamma\gamma$. The cross section predicted under the assumption of a flipped sign of the top-Yukawa coupling is depicted as a red line, the standard model cross section is depicted as a black line. The combination is not yet able to exclude the case of an anomalous coupling. The diagram is taken from Reference [138].

an inverted sign of the C_t scaling factor.

The analysis makes frequent use of multivariate analysis methods. Two neural networks are employed for the reconstruction of the heavy resonances present after the tHq and $t\bar{t}$ production. Jets are assigned to their corresponding partons based on a set of kinematic variables.

Another neural network is used to classify events based on the information gained by the reconstruction into either signal- or background-like.

The analysis is able to exclude a production cross section 7.5 times as large as predicted assuming an inverted sign of C_t at 95% confidence level. The observation corresponds to a 1.25σ upwards fluctuation as a limit of 5 times the prediction was expected.

The analysis was part of a combination effort to make a concluding statement about the tHq production at $\sqrt{s} = 8 \text{ TeV}$. Under the assumption that no other effects modify the $H \rightarrow \gamma\gamma$ branching ratio but the inverted top-Yukawa coupling, a limit of 2.7 times the predicted cross section assuming a negative C_t is found.

The final goal of the direct search for tHq production is the lifting of the degeneracy of the sign

of the coupling of the Higgs boson to the top quark. In global coupling fits performed by CMS and ATLAS, the case of an inverted sign of the coupling is disfavored by other analyses, but these are only indirectly sensitive to the sign of the coupling. However, the inclusion of beyond the standard model physics causes the indirect searches to lose most of their sensitivity on the sign of the coupling.

The direct searches in CMS are at the moment not part of the global coupling fits as there is a large overlap in events between the $t\bar{t}H$ and tHq analyses. Once the searches for tHq production are getting closer to exclusion of a negative Yukawa coupling of the Higgs boson to the top quark, arrangements have to be made to account for this overlap. A simple possibility to remove this overlap would be to demand at least one forward jet in the tHq analysis, whereas $t\bar{t}H$ analyses could veto these events. Another possibility would be to train a multivariate method to decide for each event in which analysis it would be better suited.

In the ATLAS collaboration there has been no direct searches for tHq production at $\sqrt{s} = 8$ TeV. However, the ATLAS $t\bar{t}H$ ($H \rightarrow \gamma\gamma$) analysis included the possible contribution of tHq simply to their signal under the assumption of an anomalous top-Yukawa coupling. This way their sensitivity to a negative C_t in the global coupling fit is increased compared to CMS, but they are not able to exclude or discover the associated production of Higgs boson and top quark either. But efforts are picking up steam in the ATLAS collaboration and direct searches for tHq production at $\sqrt{s} = 13$ TeV will be performed in the style of their CMS precursor analyses.

6. Search for Associated Production of Single Top Quarks and Higgs Bosons at $\sqrt{s} = 13$ TeV

The analysis of the tHq production mechanism at $\sqrt{s} = 8$ TeV does not allow to make a concluding statement about the nature of the top-Yukawa coupling. Run II of the LHC started at a higher center-of-mass energy of $\sqrt{s} = 13$ TeV and a total of 3.81 fb^{-1} of data has been collected in 2015 out of which 2.3 fb^{-1} can be used in the analysis of this chapter.

The increase in beam energy inevitably leads to an enhancement of production cross sections of most processes. The cross section of the standard model tHq production increases by a factor of 3.9 to $\sigma_{\text{tHq,SM}}^{13 \text{ TeV}} = 70.8 \text{ fb}$, when going from 8 TeV to 13 TeV. The $t\bar{t}$ production however, “only” increases by a factor of 3.7 to $\sigma_{t\bar{t}}^{13 \text{ TeV}} = 831.76 \text{ pb}$.

Despite the expected increase of the signal-to-background ratio at $\sqrt{s} = 13$ TeV, the standard model production of tHq is still out of reach for the complete Run II of the LHC, but the analysis searching for a flipped sign of the top-Yukawa coupling with $C_t = -1$, which predicts a cross section of $\sigma_{\text{tHq}, C_t = -1}^{13 \text{ TeV}} = 792.7 \text{ fb}$, will become sensitive to this production during this LHC era. The higher center-of-mass energy also leads to more centrally produced resonances, which makes a separation of signal events from background events more difficult.

Fighting the unfavorable circumstances new aspects of the analysis are introduced: Technical developments on the event generator side allow for an event-wise reweighting that can change the event kinematics based on specific generation parameters. This way not only the case of $C_t = -1$, but a whole range of C_t and C_V values can be investigated.

Another addition to the Run II analysis is the inclusion of the tHW process as a signal process. As a consequence of the inclusion of the tHW process the nomenclature of the sought signal process is generalized and changed to tH.

In this chapter the complete analysis searching for single top quarks produced in association with a Higgs boson at $\sqrt{s} = 13$ TeV is described. The focus of this chapter lies on the introduced changes with respect to the analysis at $\sqrt{s} = 8$ TeV.

A first study of this process at $\sqrt{s} = 13$ TeV has been documented in Reference [171], whereas this thesis constitutes an optimized and final revision of this analysis for the data collected in 2015.

In the first part of this chapter newly developed concepts of the analysis are introduced and described in detail. The second part describes the complete workflow of the analysis and upper limits on the tH production process are set. In the third chapter the sensitivity of the tH process to a possible $C\mathcal{P}$ -mixing state of the Higgs boson is exploited. An analysis is performed searching for the tH production process under different $C\mathcal{P}$ -mixing scenarios and ultimately

upper limits are set. In the last part of this chapter a summary and an outlook are given.

6.1. Analysis Developments

Two new major aspects of the analysis have been developed during the long shutdown of the LHC: an event reweighting which allows the investigation of multiple points in the C_V - C_t plane with only one generated simulation sample and the inclusion of a tHW signal sample.

The `MADGRAPH5_AMC@NLO` package provides a method which is able to reweight events of a sample generated under a certain theoretical hypothesis such that the distributions after application of the weights correspond to a different set of generation parameters. This method is employed in the signal samples of this analysis. The tHq and tHW samples have been produced under the hypothesis of a top-Yukawa coupling of $C_t = -1$ and $C_V = +1$, but can now be reweighted such that the distributions mirror the behavior of a sample generated under different C_t or C_V hypotheses. Available points for C_t include 3.0, 2.0, 1.5, 1.25, 1, 0.75, 0.5, 0.25 and zero as well as their negative counterparts. This set of C_t values is provided at three different C_V points: $C_V = +0.5$, $C_V = +1$ and $C_V = +1.5$, thus resulting in a total of 51 studied points in the C_V - C_t plane. The analysis is performed in parallel for these 51 different signal inputs. When changing the absolute values of C_t and C_V , the masses of top quark and W boson are assumed to be standard model-like as they are experimentally very well established.

The procedure also allows the storage of weights, whose application reproduces systematically shifted simulation samples, e. g. for samples produced with a different Q^2 scale.

An application of the reweighting for the signal sample can be seen in Figure 6.1, where the generated transverse momentum of the top quark and the generated pseudorapidity of the Higgs boson are depicted for different C_t points. The reweighting procedure works as expected, but it is visible that statistical fluctuations are enhanced, when weights larger than one are applied. The reweighting procedure is explained in detail in Reference [172].

The newly included signal process tHW corresponds to a tW -channel production of a single top quark in association with a Higgs boson. The hard interaction of the tHW process has been produced with the `MADGRAPH5_AMCATNLO` package at leading order. The parton shower step has been performed with `PYTHIA8`.

Interference between the tHW channel produced in the 4F scheme and $t\bar{t}H$ production would occur, as the initial and final state particles are the same ($gg \rightarrow WbWb$) for both processes. A coherent treatment of the interference would require a manual intervention when generating this sample in the 4F scheme, but would at the same time rule out the LHE reweighting for the tHW process. This technical hurdle is not overcome yet, thus the tHW process is generated in the 5F scheme at LO employing the LHE reweighting method. The cross sections are separately calculated at NLO for each of the 51 points in the C_V - C_t plane and the samples are scaled to their corresponding values. Other major changes include the migration to the medium b-tagging working point effectively increasing the number of events available for the analysis, and the usage of BDTs instead of NNs as multivariate methods for the reconstruction and classification.

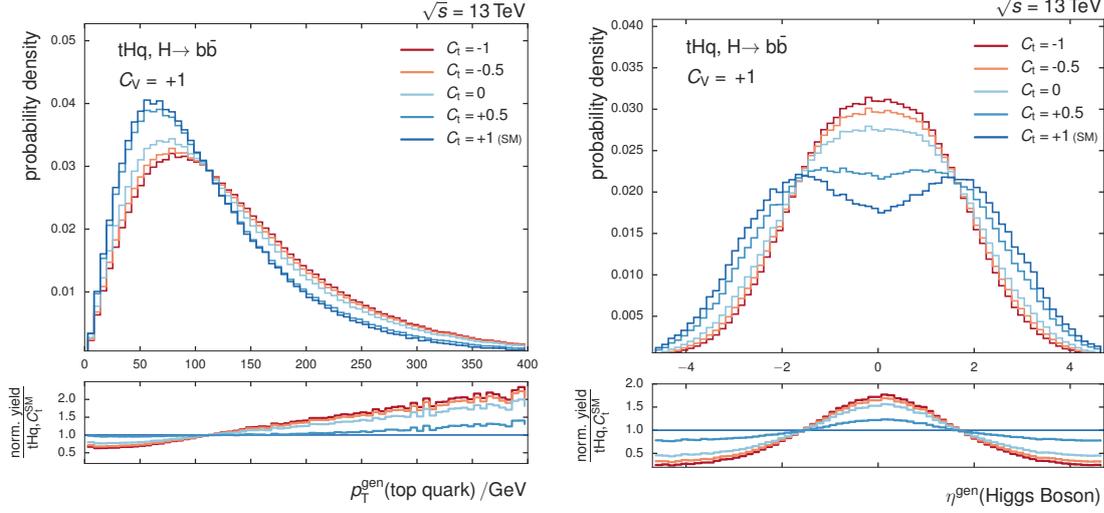


Figure 6.1.: Distributions for the generated transverse momentum of the top quark (left) and the generated pseudorapidity of the Higgs boson (right) for different values of C_t using the LHE reweighting. Higgs boson and top quark show a similar behavior, so only one distribution for each is shown. All histograms are normalized to illustrate the changing shape of the variables. Going from the SM prediction to $C_t = -1$ the p_T spectrum is shifted to higher values and the objects are produced more centrally. In the tail of the p_T spectrum a possible problem of the reweighting is visible, as statistical fluctuations are amplified when weights larger than one are applied. The distributions are produced with an inclusive tHq sample, where at least one charged lepton has been reconstructed.

6.2. Analysis Strategy

The strategy of this analysis is very similar to its predecessor analysis at $\sqrt{s} = 8$ TeV. Selection criteria which are used to optimize the expected signal-to-background ratio have been slightly altered. The analysis is extended from a search for tHq production under the hypothesis of a flipped sign of the top-Yukawa coupling to the search for tH at 51 different points in the C_V - C_t parameter plane. Especially the event reconstruction method assigning jets of an event to the partons present in the tHq final state requires a revision. As the kinematics of the signal process change for each C_V - C_t parameter point, a separate jet assignment method is trained for every coupling point. Each event is therefore provided with 51 possibly different jet assignments under the tHq event hypothesis. The assignment of jets to the partons of a semi-leptonic $t\bar{t}$ final state is unaffected by the extension to the C_V - C_t parameter plane.

As final stage in the analysis the limit is extracted from an output distribution of a multivariate classifier which separates signal from background events based on a set of input variables. Similarly to the Run I analysis, these input variables are mostly observables of reconstructed objects dependent on the event reconstructions. The 51 different tHq jet assignments lead to 51 different sets of input variables, thus necessitating a separate training of 51 final classification BDTs, one for each point in the C_V - C_t parameter plane. After the assignment of a common set of uncertainties to the different processes, an independent fit for the classification outputs is performed. Based on these classifier outputs an expected CL_S limit can be calculated for each

investigated point.

A separate analysis is performed searching for a possible $C\mathcal{P}$ -mixture state of the Higgs boson in tH production. A tH sample has been produced privately, which employs the same LHE reweighting, allowing the study of 21 different $C\mathcal{P}$ -mixing angles of the studied boson. The complete analysis strategy is copied, but reevaluated for these 21 different signal models, and upper limits are set for each of the studied angles.

6.3. Signal Processes

6.3.1. t -channel

The cross section of the SM tHq production increases from $\sigma_{\text{SM}}^{8\text{ TeV}} = 18.3$ fb at $\sqrt{s} = 8$ TeV to $\sigma_{\text{SM}}^{13\text{ TeV}} = 71.0$ fb at $\sqrt{s} = 13$ TeV. For the case of a flipped sign of the top-Yukawa coupling the cross section is increased to $\sigma_{\text{tHq}, C_t=-1}^{13\text{ TeV}} = 792.7$ fb. The remaining cross sections of the considered 51 coupling configurations are visualized in Figure 6.2. The exact values can be found in the appendix in Table B.3. The coupling point of the highest interest is still the case of $C_t = -1$ and $C_V = +1$ and will henceforth be singled out in the course of this thesis.

In retrospect to the tHq production at $\sqrt{s} = 8$ TeV the two heavy resonances are produced more centrally at an energy of $\sqrt{s} = 13$ TeV. The higher center of mass energy also leads to higher transverse momenta of all involved particles. This behavior is visible in the distributions shown in Figure 6.3. Top quark and Higgs boson both show very similar kinematics. All of the involved particles are expected to have an increased transverse momentum and are produced in a slightly more forward direction.

6.3.2. tW -channel

The tHW process is based on a conversion of an off-shell bottom quark into the much heavier top quark and an accompanying W boson in analogy to the tW -channel of single top quark production. The associated Higgs boson can again be emitted from the W boson or from the top quark causing the interference responsible for the sensitivity on the sign of C_t .

The cross section of the tHW process of $\sigma_{\text{tHW}, C_t=-1}^{13\text{ TeV}} = 147.2$ pb for a flipped sign of the top-Yukawa coupling is smaller by roughly a factor of four than that of the tHq process, as is visualized in Figure 6.2. The actual cross section values for the 51 studied points can be found in Appendix B.4.

The presence of a bottom quark in the initial state necessitates the occurrence of a second b quark, just as in the tHq production. But again, the second b -quark tends to be softer and to be produced in a more forward direction, therefore reducing its detection efficiency. Together with the b quark of the top quark decay and the two b quarks of the Higgs boson decay, the tHW production is also characterized by three or four detectable b quarks.

The final state that is contributing the most to the considered signal regions is the semi-leptonical decay of tHW . Due to the presence of two W bosons in the process an ambiguity in the decay is apparent as either of the W bosons can decay leptonically and the other can decay hadronically. By requiring both W bosons to decay leptonically a separate dilepton signal regions for tHW could be defined but this is out of the scope of this analysis and could serve as

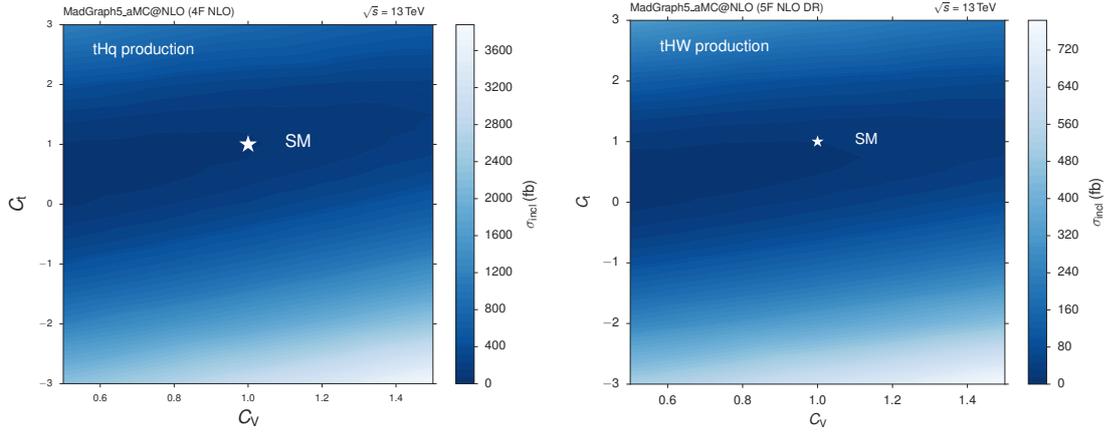


Figure 6.2.: Visualization of the cross sections of the tHq process (left) and the tHW process (right) for different points in the C_V - C_t parameter plane.

a possible analysis improvement in the future.

The final state therefore contains four b quarks, one lepton, one neutrino and two light jets. Compared to the tHq production, one additional light jet is produced, but none of these light jets are necessarily produced in the forward region of the detector. This leads to semblance of the tHW process to the $t\bar{t}$ background process.

6.4. Background Processes

The analysis is performed under consideration of the same backgrounds as in its predecessor analysis described in Section 5.3. The only newly considered background process is the $t\bar{t}W$ production. In the $t\bar{t}W$ process a W boson is predominantly emitted by one of the quarks in the initial state of the top quark pair production resulting in three W bosons, which could in principle fake the signature of a signal event, if two bosons decay hadronically and one boson leptonically. But its very small cross section causes the contribution of this background to be insignificant. In the distributions shown in this chapter the $t\bar{t}W$ process is grouped together with the Diboson and V+jets background under the label "Misc".

In principle also backgrounds like $t\bar{t}Z$, the emission of a Z boson in the $t\bar{t}$ production, or tZq , the emission of a Z boson during single top production, could contribute to the background, but low production cross section as well as a low branching ratio of $Z \rightarrow b\bar{b}$ makes them negligible.

6.5. Datasets

The data analyzed in this chapter corresponds to $\mathcal{L}_{\text{int}} = 2.3 \text{ fb}^{-1}$ collected by the CMS experiment in 2015 during the Run2015D era. The SingleElectron and SingleMuon datasets are used in this analysis after a selection of all certified luminosity sections collected in the golden json [173]. The standard single lepton trigger paths HLT_Ele27_eta2p1_WPLoose_Gsf_v*, HLT_IsoMu20_v*

6. Search for tH Production at $\sqrt{s} = 13$ TeV

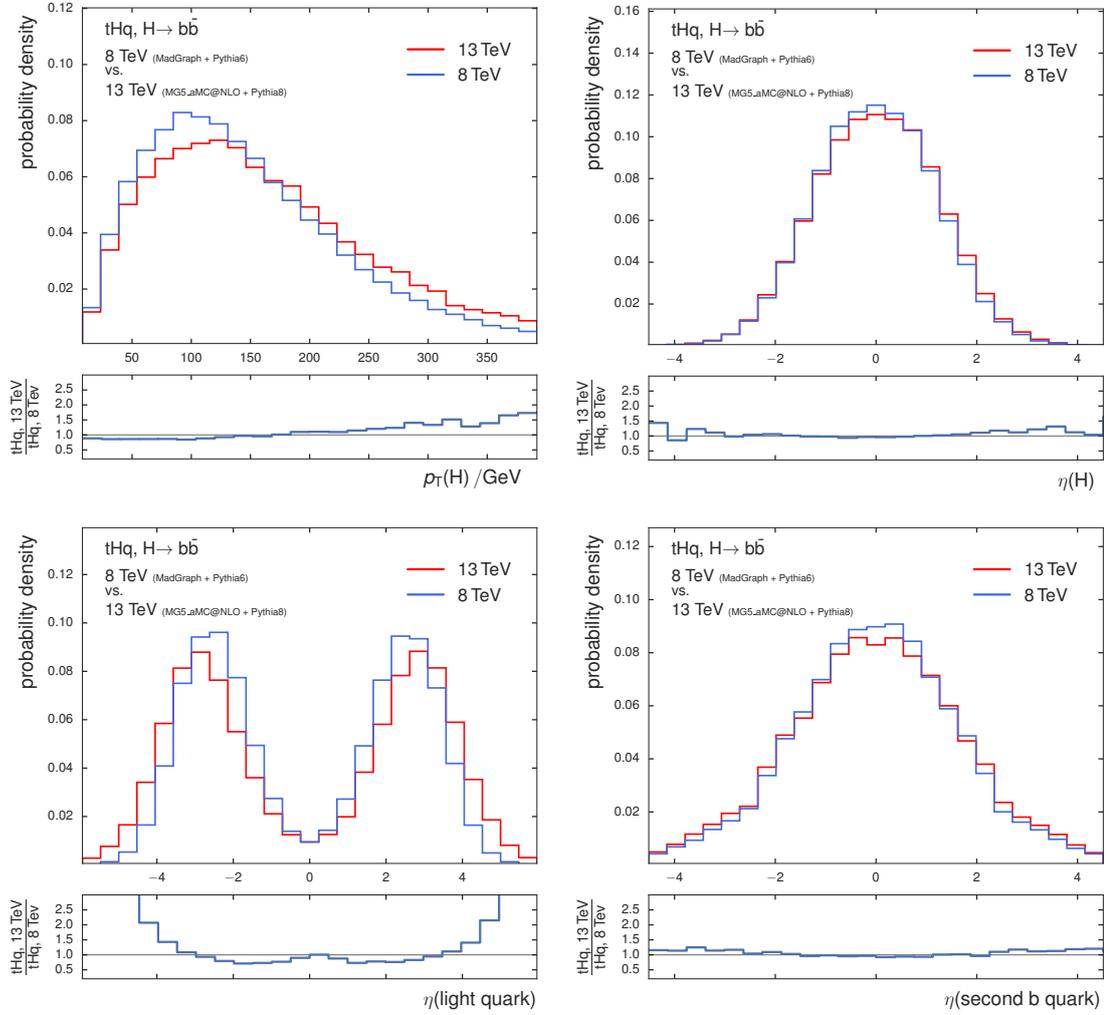


Figure 6.3.: Distributions of kinematic variables of objects of the tHq process for 8 TeV and 13 TeV. The kinematics of the top quark in the process are very similar to those of the Higgs boson, hence only one set is shown. All objects are produced with a higher transverse momentum and are produced in a more forward direction at 13 TeV. This effect is much more pronounced for the light quark in comparison to Higgs boson, top quark and second b quarks.

and `HLT_IsoTkMu20_v*` are utilized in this analysis.

The `tHq` and `tHW` signal processes are generated with `MADGRAPH5_AMCATNLO` interfaced with `PYTHIA8` as parton shower. The generation process implements the `NNPDF3.0` PDF set and a new dynamical scale $\mu = (m_T(t) + m_T(H) + m_T(b))/6$ is employed, which improves the modeling of the event kinematics for larger Q^2 scales. Additionally, the standard underlying event tune of `Pythia6` at $\sqrt{s} = 8$ TeV, `TuneZ2*`, is replaced by the `CUETP8M1` tune [89]. Both signal processes are scaled to their respective NLO cross sections. A complete table with all used data and simulation samples can be found in Appendix B.1 and B.2. When deviating from the nominal $C_V = +1$ value to $C_V = +1.5$ or $C_V = +0.5$ the cross sections of the signal samples have to be adapted due to a change of $\mathcal{BR}(H \rightarrow b\bar{b})$. This is caused by a respective increase or decrease of $\mathcal{BR}(H \rightarrow WW)$, which subsequently changes all other branching ratios such that unitarity is conserved. For $C_V = +1.5$ the cross section is scaled down by 23.3% and for $C_V = +0.5$ the cross section is scaled up by 22.24% [171]. Due to the presence of a Higgs boson, this is also done for the `tH` process, which is additionally scaled to its NLO cross section for each of the studied C_t values.

The main background, the `tH` pair production, is generated with `POWHEG` interfaced with `PYTHIA8` as parton shower. The analysis uses exclusive semi-leptonic and a full-leptonic samples for the simulation of the `tH` production, the impact of a fully hadronically decaying `tH` pair has been found to be negligible. The modeling of emissions in the `tH` process containing heavy flavored quarks, like charm or bottom quarks, still poses an impediment as the measured ratios of heavy flavored to light flavored emissions are not reproduced in the event generation. The `tH` production sample is again separated into subsamples based on the heavy flavor content of the event, but instead of the parton information, which was the basis for the splitting in the 8 TeV analysis, the splitting for this analysis is based on hadron information. This allows for an easier comparison between theoretical predictions and experimental measurements [174].

The four categories that were used in the Run I analysis are extended by splitting the former `tH+b` category into two. The categories are now defined as follows:

- tH+bH** The event contains at least two additional jets that could each be matched to one B hadron.
- tH+2b** The event contains one additional jet that could be matched to at least two B hadrons.
- tH+b** The event contains one additional jet that could be matched to exactly one B hadron.
- tH+cH** The event contains at least one additional jet that could be matched to at least one D hadron.
- tH+light** All other events not satisfying any of the criteria listed above.

6.6. Object Definitions

The objects used in this analysis are defined in a very similar way as in the analysis at $\sqrt{s} = 8$ TeV. Improved object reconstruction methods and adapted object selection suggestions, which emerged during LS1, are utilized. The selections that are applied to define the actual physical objects used in this analysis are described in the following section.

6.6.1. Primary Vertices

Primary vertices are defined as described in the Run I analysis in Section 5.5.1.

6.6.2. Muons

The muons studied in this analysis have to satisfy all requirements of the tight muonID [175], which coincides with the requirements enforced in the 8 TeV analysis. Two collections of muons are defined, the tight and the loose muon collection:

Muons passing the tight selection criteria must have a transverse momentum larger than 25 GeV and their pseudorapidity must satisfy $|\eta| < 2.4$. The relative isolation as defined in Equation 5.1 has to satisfy $I_{\Delta\beta} < 0.15$.

For muons passing the loose selection criteria the requirements on the transverse momentum and the relative isolation are relaxed. Their transverse momentum has to be larger than 20 GeV and the isolation has to satisfy $I_{\Delta\beta} < 0.20$. Otherwise, the definitions of these two collections are identical.

6.6.3. Electrons

The electrons used in this analysis have to pass the working point defined at 80% efficiency of the triggering MVA ID [176] and have to be measured outside the gap of the ECAL.

In analogy to the reconstructed muons two electron selections are defined: the transverse momentum of tight electrons needs to surpass 30 GeV and their pseudorapidity must fulfill $|\eta| < 2.1$. Their isolation, as defined in Equation 5.2, needs to be smaller than 0.15.

For loose electrons the restrictions on transverse momenta and pseudorapidity are relaxed, as they have to satisfy $p_T > 15$ and $|\eta| < 2.5$.

6.6.4. Jets

Jets used in this analysis are reconstructed with the anti- k_t algorithm with a smaller cone of $R = 0.4$ compared to the Run I analysis. Pileup candidates are removed from the jets via the CHS algorithm. The jets must pass the particle flow jet ID [177] and the individual jets are removed, if they are closer than $\Delta R = 0.4$ to a tight lepton.

Similarly to the 13 TeV analysis, jets need to satisfy $p_T > 20$ GeV to be considered in the analysis. Jets in simulation are corrected with the appropriate L1L2L3 MCtruth corrections, which are described in Section 3.3.6, and the residual corrections are applied to jets in data [155]. Also the jet energy resolution in simulation is smeared to match the resolution observed in data [178].

6.6.5. Missing Transverse Energy

The missing transverse energy provided in simulation is recorrected for possible effects caused by jet corrections. The four-vectors of the initial, uncorrected jets are added to the four-vectors of the missing transverse energy, thereby undoing the \cancel{E}_T calculation. From the resulting four-vector the corrected jets are subtracted ensuring a coherent recalculation of the missing transverse energy.

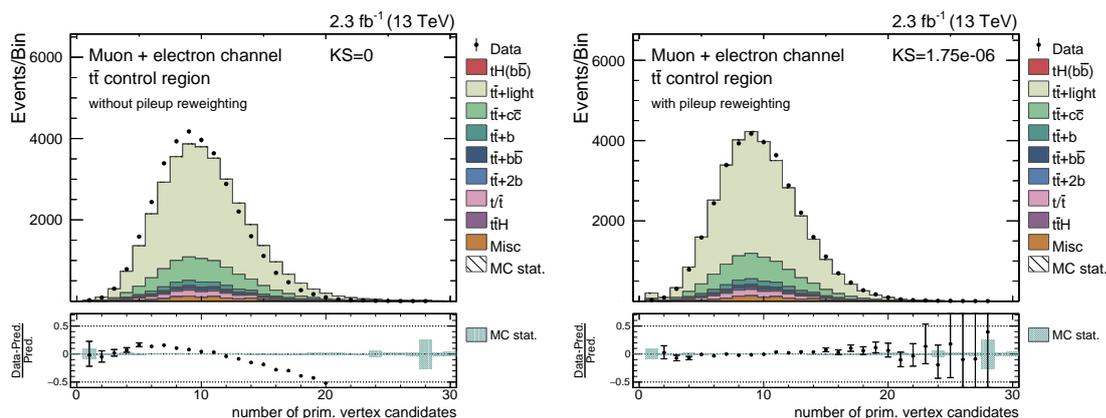


Figure 6.4.: Distribution of the number of primary vertices per event before (left) and after (right) application of the pileup reweighting. The simulation predicted slightly more primary vertices than are observed in data. For both distributions the simulation is scaled to match the observed event numbers in data. After the reweighting a good agreement of simulation to recorded events is observed.

6.6.6. W Boson Reconstruction

The reconstruction of the W boson of the leptonically decaying top quark is done in an analogous manner to that of the Run I analysis, described in Section 5.5.6.

6.7. Monte Carlo Corrections

6.7.1. Pileup Reweighting

Similarly to the procedure described in Section 5.6.1 a reweighting of the distribution of the number of reconstructed primary vertices is applied. The pileup environment seen in simulation is reweighted to match the distribution observed in data. The effect of the reweighting can be seen in Figure 6.4.

6.7.2. Lepton Efficiency Scale Factors

The correction of the muon efficiencies encompasses separate scale factors for the muon trigger, identification and isolation efficiencies. The correction factors are provided by the Muon POG [179] and are derived with tag-and-probe methods [180] at the J/Ψ or Z boson resonances.

Simulated electron efficiencies for the reconstruction as well as for the triggering MVA-ID are corrected [181]. The efficiencies are provided by the EGamma POG and have been determined on large $Z \rightarrow e^+e^-$ samples with the tag-and-probe procedure [182].

The efficiency scale factors cause a slight change in shape for some distributions and a reduction of the overall yield of $\sim 3\%$.

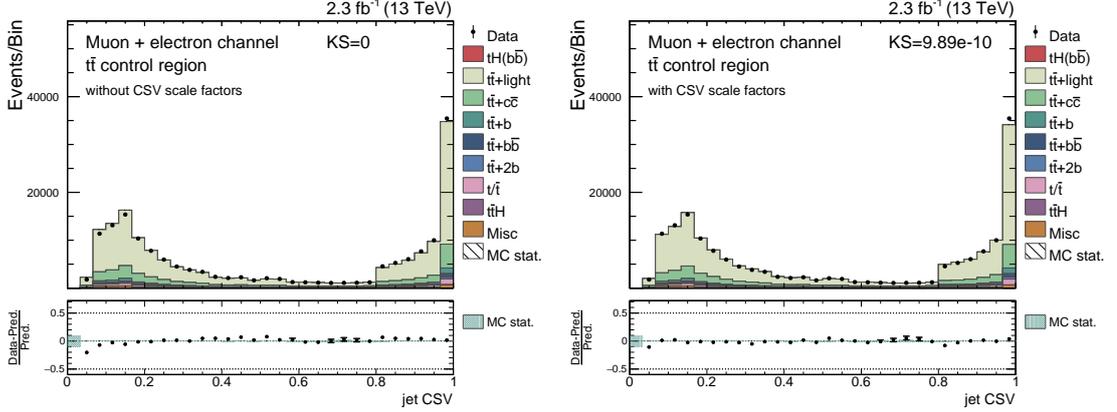


Figure 6.5.: Distribution of the CSV output for all jets per event before (left) and after (right) application of the CSV reweighting. For both distributions the simulation is scaled to match the observed event yields in data. After the reweighting an improved agreement of simulation to recorded events is observed.

6.7.3. CSV Shape Reweighting

Discrepancies between data and simulation in the output of the b-tagging algorithm are handled differently than compared to the 8 TeV analysis. Every event is assigned a weight such that the complete measured CSV distribution in data is reproduced by the simulation [183]. Weights are calculated separately for heavy and light flavored jets. The applied scale factors are functions of the jet CSV output value, the p_T and the η of the jet. The scale factors for the heavy flavored jets are derived via tag-and-probe method in a control region enriched with fully leptonically decaying top quark pairs, whereas scale factors for light flavored jets are obtained in a $Z(\ell\ell)$ control region. The scale factors to account for light flavor contamination are calculated as

$$SF_{HF}(CSV, p_T, \eta) = \frac{\text{Data} - MC_{LF}}{MC_{HF}}, \quad (6.1)$$

where for events in data the presence of one b-tagged jet is required and MC_{LF} and MC_{HF} are the simulated yields of light flavored jets and heavy flavored jets, respectively. For the calculation of the scale factors for light flavored jets, the yields in data are estimated by requiring one jet to be untagged and MC_{HF} and MC_{LF} switch places in Equation 6.1. The effect of the reweighting can be seen in Figure 6.5.

This procedure allows for the exploitation of the full information of the CSV output distribution, thus making it possible to employ a complete CSV distribution of single jets as input variables in the reconstruction and classification. More information on the scale factor estimation can be found in Reference [184].

6.7.4. Jet Pseudorapidity in the Forward Region

A similar behavior for the jet pseudorapidity in the forward region is observed than in the 8 TeV analysis when comparing simulated jets to jets in data. As in the predecessor analysis the issue is mitigated via the same measures described in Section 5.6.5.

Table 6.1.: Overview of the selection criteria applied to define the two signal regions, the three medium-tag region (3M) and the four medium-tag region (4M).

Reconstructable jets are defined as satisfying $p_T > 30$ GeV, if they are reconstructed in the central region of the detector, and satisfying $p_T > 40$ GeV, they are reconstructed in the forward region with $|\eta| > 2.4$. The additional jets are part of the standard jet collection with $p_T > 20$ GeV.

	3M region	4M region
# reconstructable jets	≥ 4	≥ 5
# rec. jets with CSV > 0.80	3	4
# tight leptons	1	1
# additional loose leptons	0	0
\cancel{E}_T	$> 35/45$ GeV(μ/e)	$> 35/45$ GeV(μ/e)

6.8. Event Selection

The event selection applied in this analysis is similar to that of the analysis at $\sqrt{s} = 8$ TeV. One major change is the migration from the tight b-tagging working point to the medium working point. The output of the CSVv2 algorithm for a jet has to be larger than 0.8 for it to be b-tagged. The tagging efficiencies for the medium working point improved w. r. t. Run I, thereby increasing the expected significance of the analysis substantially. Additionally, the migration to the medium working point mitigates a possible problem of a low event count in the MC signal samples. A high number of signal events is desirable to ensure a coherent behavior of the tHq reconstruction and the classification, as both have to be trained with a significant amount of signal events that are subsequently discarded from the analysis. Another measure to ensure a reasonably high event count is the consolidation of the muon and electron channel into a combined lepton channel.

Another small change w. r. t. the 8 TeV analysis is that only jets satisfying $p_T > 30$ GeV are considered to be b-tagged. As in the 8 TeV analysis jets in the forward region are only considered as a reconstructable jet, if they satisfy $p_T > 40$ GeV. Jets in general are part of the jet collection, if their transverse momentum is larger than 20 GeV. This allows for a simplification of the event selection by abolishing the requirement on the additional jet category, as the presence of a jet which can be assigned to the light quark is already implied by the number of required reconstructable jets.

Reconstructable jets are defined as in the 8 TeV analysis (see Section 5.7). The applied selection criteria are summarized in Table 6.1. The requirement on the missing transverse energy is kept at the same level to suppress the QCD background to a negligible level, however a future data-driven estimation of the QCD background at higher integrated luminosities could allow for a relaxation of this requirement, thereby increasing the number of signal events in the signal regions. The signal regions are defined by having either three jets passing the medium b-tagging working point (3M) or four jets passing the medium working point (4M). Additionally, a $t\bar{t}$ control region is defined, in which events are required to have exactly two b-tagged jets. This region is greatly enriched in $t\bar{t}$ events with a purity of 94% and can be used to validate a

good modeling of the $t\bar{t}$ background. A normalization difference of $\sim 10\%$ is observed between simulation and data, which is most likely caused by an interplay of a mismodeling of the jet multiplicity in $t\bar{t}$ simulation samples and the CSV reweighting procedure as the effect gets less significant for lower jet and higher b-tag multiplicities. However, this offset is covered by the applied systematic uncertainties. The mismodeling of jet multiplicities is observed in all CMS analyses at the time of writing and is attributed to an issue in the PYTHIA 8 parton shower parameters during the official MC simulation production [185]. Besides the mentioned normalization offset, no further impact on the analysis has been found. A set of control distributions, which are normalized to the observed number of events in order to facilitate shape comparisons, for general observables in the $t\bar{t}$ control region can be found in Figure 6.6. Apart from the jet multiplicity mismodeling a generally good agreement is observed.

6.9. Event Reconstruction

The event reconstruction employing MVA techniques remains the foundation of the variables used in the classification. By training the technique to learn the differences between a correct and a random jet-parton assignment the properties of reconstructed objects can be exploited in the final classifier. Two separate jet assignments are applied to each event, once under the hypothesis of it being a tH_q event and once under the hypothesis of it being a $t\bar{t}$ event. By emphasizing the features of the different processes an increased separation between signal events and the overwhelming $t\bar{t}$ background can be achieved.

In contrast to the analysis at $\sqrt{s} = 8$ TeV BDTs are used as MVA technique instead of NNs, as their performance is found to be on the same level as NNs but are trained much faster, which is especially important when optimizing the training parameters and the used employed variable set for 51 different trainings.

6.9.1. tH_q Reconstruction

The extension of the analysis to multiple points in the C_V - C_t plane makes a revision of the tH_q reconstruction necessary. As the kinematic distributions differ for various points, a generalized training for only one reference point could lead to wrong jet assignments for the tH sample, thereby diminishing the success rate of the reconstruction. To counter this effect a distinct training is performed for each of the 51 points in the C_V - C_t parameter plane. The basic structure of the reconstruction is identical to the reconstruction at $\sqrt{s} = 8$ TeV as explained in Section 5.8.1 and is not be discussed here. The set of variables used in the training is identical to ensure a good comparability among the different reconstructions.

The set of variables has been optimized with the help of the area under the ROC as performance measure and includes now fifteen variables, whose descriptions can be found in Table 6.2. The variable importance differs in the individual trainings for the 51 studied coupling points, hence an adapted ranking mechanism is used. For each training the most important variable is assigned with fourteen points and each successive rank is assigned with one point less, until the

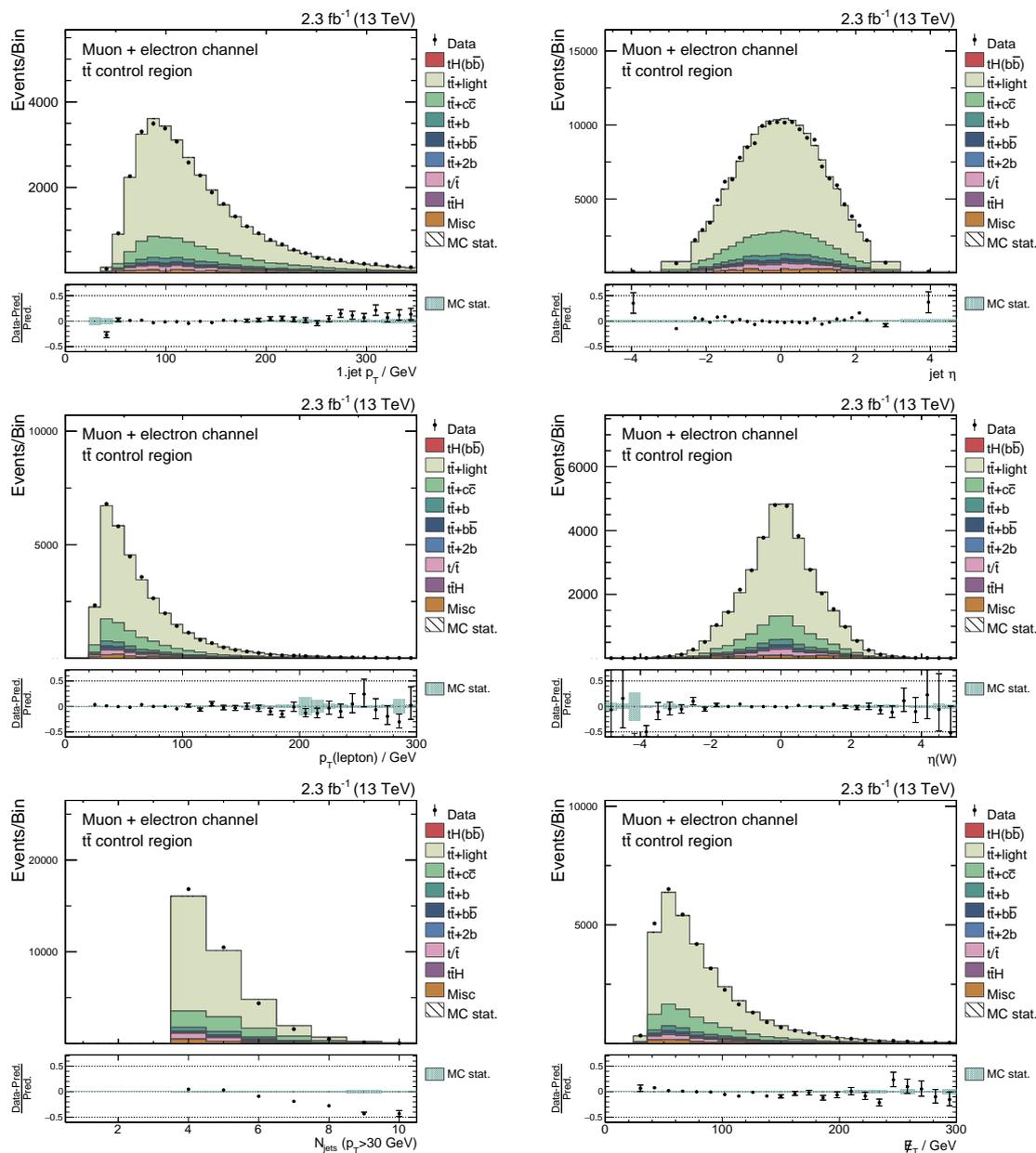


Figure 6.6.: A set of control plots in the $t\bar{t}$ control region is shown. The diagrams show the transverse momentum of the hardest jet of the event, the pseudorapidity of all jets, the transverse momentum of the electron, the pseudorapidity of the reconstructed W boson, the number of jets with $p_T > 30$ GeV and the missing transverse energy of the event. For all distributions the simulation is scaled to match the observed event yields in data. The N_{jet} distribution shows a known mismatching, based on a known issue in the parton shower settings in the official MC simulation production. The remaining variables show a good agreement between simulation and data.

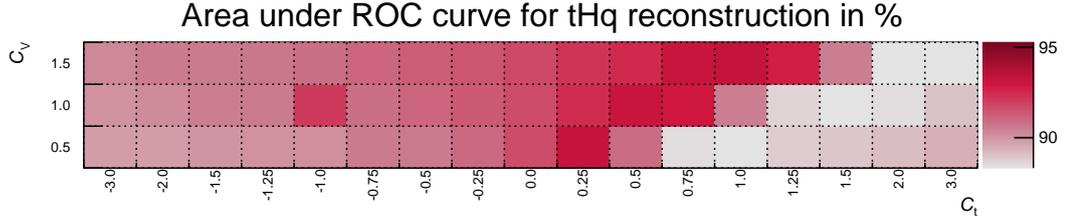


Figure 6.7.: The area under the ROC curve for all 51 tHq reconstruction trainings is shown. With values above 85% good performance is observed for all trainings. A visible ridge on the right side of the plot corresponds to the points with the lowest cross sections. Effectively lower event counts lead to a deterioration of the separation power between correct and wrong assignments.

last rank is assigned with zero points. The sum of all points is then averaged over the number of trainings, leading to a ranking system where a variable with fourteen points would be the most important variable in all trainings and a variable with zero points would always be the least important variable. The outcome of the ranking is quoted in Table 6.2. The invariant masses of Higgs boson and top quark are the most important variables overall, whereas $\cos \theta(t, \ell)$, the cosine of the angle between the top quark and the charged lepton in the lepton rest frame, is the most important novel variable. The distributions for the six most important variables for correct and wrong assignments can be found in Figure 6.9. The remaining nine variables are found in Appendix B.1 and B.2.

Out of all available simulated tHq events in the 3M region, one fifth is exclusively used for the training reconstruction, another five percent are used for testing.

The application of the LHE weights leads to a reduction of the effective event count in the training for some of the studied coupling points, which in turn increases the chance of overtraining. A reduction of the number of trees in the BDT training is applied for the coupling points that have shown signs of overtraining. The parameters used in the training are summarized in Table 6.3.

After securing that no overtraining occurred at any parameter point, the performance of the training is quantified by calculating the area under the ROC curve for the reconstruction classifier. The output for each of the 51 points can be found in Figure 6.7. Each of the 51 trainings shows a good performance, but it is apparent that the separation between correct and wrong assignments is worse for points corresponding to lower cross sections of the process, as the training is effectively performed on a smaller number of events. This effect is already visible without the reduction of the number of trees, but is amplified by that measure. The response of the reconstruction for the point of $C_t = -1$ and $C_V = +1$ and exemplary for its standard model counterpart of $C_t = +1$ and $C_V = +1$ can be seen for the training and the disjoint testing sample in Figure 6.8. After the successful training, the reconstruction BDTs are applied to all simulation and data events, where the BDT response is calculated for each allowed jet assignment. The assignment set with the highest output is selected for the studied event. A comparison of the highest output value for MC simulation and data can be seen in Figure 6.10 for the $C_t = -1$ and $C_V = +1$ case and the SM case. A good agreement between simulation and data is observed for both distributions.

Table 6.2.: Input variables for the jet-assignment BDT under the tHq hypothesis sorted by their importance in the training. Instead of the invariant masses and transverse momenta their logarithms are used, as narrow distributions are better suited for the usage in MVA techniques than distributions with long tails.

Variable	Points	Description
$\log m(\mathbf{H})$	14.00	invariant mass of the reconstructed Higgs boson
$\log m(\mathbf{t})$	11.94	invariant mass of the reconstructed top quark
$\Delta R(\text{Higgs jets})$	11.92	ΔR between the two jets from the Higgs boson decay
$\Delta R(\mathbf{b}_t, \mathbf{W})$	10.58	ΔR between jets assigned to the b quark from the top quark decay and the W boson
relative H_T	9.05	percentage of the total transverse momentum (jets, lepton, \cancel{E}_T) that falls to the b jet of the top quark, Higgs jets and light forward jet
$\cos \theta(\mathbf{t}, \ell)$	7.66	Cosine of the angle from the top quark vector to the sum vector of top quark and charged lepton in their common restframe
CSV(Higgs jet 2)	6.31	output of the CSVv2 b-tagging algorithm for the second jet assigned to the Higgs boson
CSV(\mathbf{b}_t)	6.27	output of the CSVv2 b-tagging algorithm for the jet assigned to the b quark from the top quark decay
$ \eta(\text{light jet}) - \eta(\mathbf{b}_t) $	6.00	absolute difference of pseudorapidities of the light forward jet and b quark of the top quark decay
CSV(Higgs jet 1)	5.72	output of the CSVv2 b-tagging algorithm for the first jet assigned to the Higgs boson
$ \eta(\mathbf{t}) - \eta(\mathbf{H}) $	5.56	absolute difference of the pseudorapidities of the reconstructed top quark and Higgs boson
$ \eta(\text{light jet}) $	4.54	absolute pseudorapidity of the light forward jet
$\log \min(p_T(\text{H jets}))$	3.35	lower transverse momentum of the two jets assigned to the Higgs boson decay products
$\Delta E(\text{light jet}, \mathbf{b}_t)$	1.09	jet energy difference of the light jet and the jet assigned to the b quark from the top quark decay
$ \eta(\mathbf{b}_t) $	0.94	pseudorapidity of the jet assigned to the b quark of the top quark decay

6. Search for $t\bar{H}$ Production at $\sqrt{s} = 13$ TeV

Table 6.3.: Parameter settings used for the training of the employed BDTs for the two reconstructions in the TMVA software package. A smaller number of trees is used for points that have otherwise shown signs of overtraining. These points include the following $(C_t | C_V)$ -value pairs: (1|1), (1.25|1), (1.5|1), (1.5|1.5), (2|1.5) and (0.5|0.5). After this reduction no sign of overtraining is found. The definitions of the configuration options can be found in Reference [117].

Parameter	Value
NTrees	400/150
MinNodeSize	1
MaxDepth	3
BoosteType	AdaBoost
nCuts	20
AdaBoostBeta	0.3
SeparationType	GiniIndex

6.9.2. $t\bar{t}$ Reconstruction

The reconstruction of the top quark pair is revised in comparison to the reconstruction at $\sqrt{s} = 8$ TeV. As the kinematics of the $t\bar{t}$ process do not change when varying C_V or C_t a single reconstruction is sufficient for each of the studied points in the C_V - C_t plane. Similarly to the tHq reconstruction BDTs are employed, with the same set of parameters which can be found in Table 6.3, instead of NNs. The training is performed on the semi-leptonic $t\bar{t}$ simulation sample, out of which 16% of all available events are used for the training and 4% are used for testing. The prerequisites on the allowed assignment permutations are adapted from the 8 TeV analysis such that jets that can be assigned to the b quarks of the top quark decays must satisfy the medium working point instead of the tight working point.

The variable set which is used in the training has been optimized with regard to the analysis at $\sqrt{s} = 8$ TeV, resulting in an omission of the jet charge variables. The modeling of jet charges in simulation showed small deviations when compared to data prohibiting the inclusion of related variables in the training. Additionally the variable of the invariant mass of the jet assigned to the b quark of the leptonic top quark decay and the charged lepton is substituted by the invariant mass of the reconstructed leptonically decaying top quark. The description of the eleven variables which are used in the training can be found in Table 6.4 sorted by their importance in the training. Distributions of the six most important variables for correct and wrong jet assignments can be found in Figure 6.11. Properties of the reconstructed hadronically decaying top quark, especially the invariant masses of top quark and W boson, and the distance of the objects in the detector to each other, have a large impact on the training. The training yields an extremely good AUC value of 95.5%, better than any of the tHq reconstructions. The response distribution of the training set and of the independent testing set can be found in Figure 6.12 and no sign of overtraining is observed. The reconstruction is subsequently applied

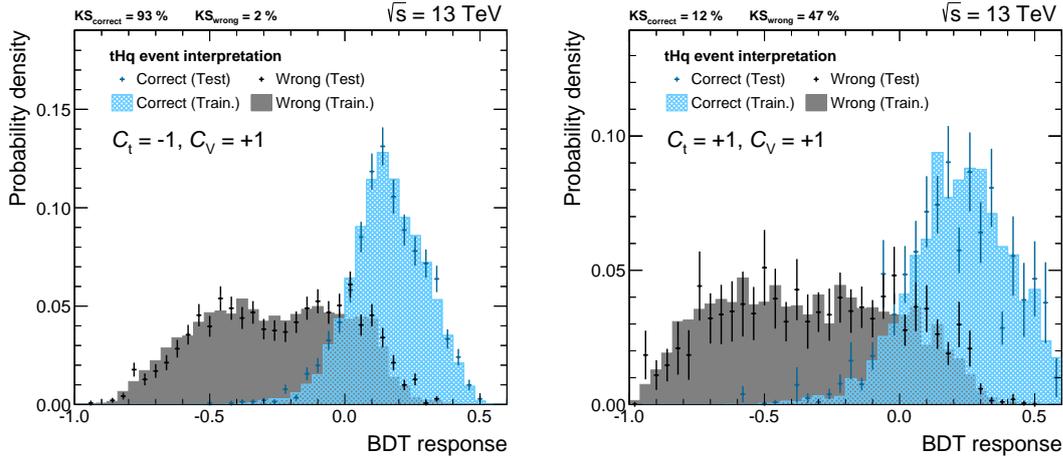


Figure 6.8.: Output values of the tHq reconstruction BDTs for correct and wrong jet assignments for the coupling case of $C_t = -1$ and $C_V = +1$ (left) and the SM prediction (right). A clear separation is visible for both trainings.

The reconstruction is examined with an independent set of events that were not part of the training sample. The events used in the training are presented as histograms, whereas the independent testing sample is represented by markers. A good agreement between the testing and training sample is observed, verified by decent KS-values of the two distributions. No sign of overtraining can be found. The larger uncertainty bars for the standard model training are a direct consequence of the reduction of the effective event count by the LHE reweighting procedure.

to all simulation and data events. The comparison of simulation and data of the highest response value per event and of a response value for a random jet assignment can be seen in Figure 6.13.

6.9.3. Evaluation of Reconstruction Methods

The reconstruction is again evaluated against a simpler reconstruction method which uses a χ^2 measure to select one of the allowed jet assignments. A detailed description can be found in Section 5.8.3. The successful assignment rates for the BDT reconstruction and the χ^2 reconstruction can be found in Figure 6.14.

tHq Reconstruction Evaluation

The tHq reconstruction yields similar results as in the 8 TeV analysis, which was not granted given the change from NNs to BDTs, the migration to the medium working point and more forward directed jets at 13 TeV. The event reconstruction improved slightly regarding the assignment of jets to the three bottom quarks in the event. The complete Higgs boson is reconstructed successfully in 66.8% (3M) and 61.3% (4M) of all events. The b quark from the top quark decay is reconstructed correctly in 67.1% (3M) and 57.2% (4M) of the studied events. The reconstruction of the light jet worsened slightly by 2% and 4% in the 3M and 4M region, respectively, when compared to the predecessor analysis. It is noticeable that the χ^2 reconstruction, which per-

6. Search for tH Production at $\sqrt{s} = 13$ TeV

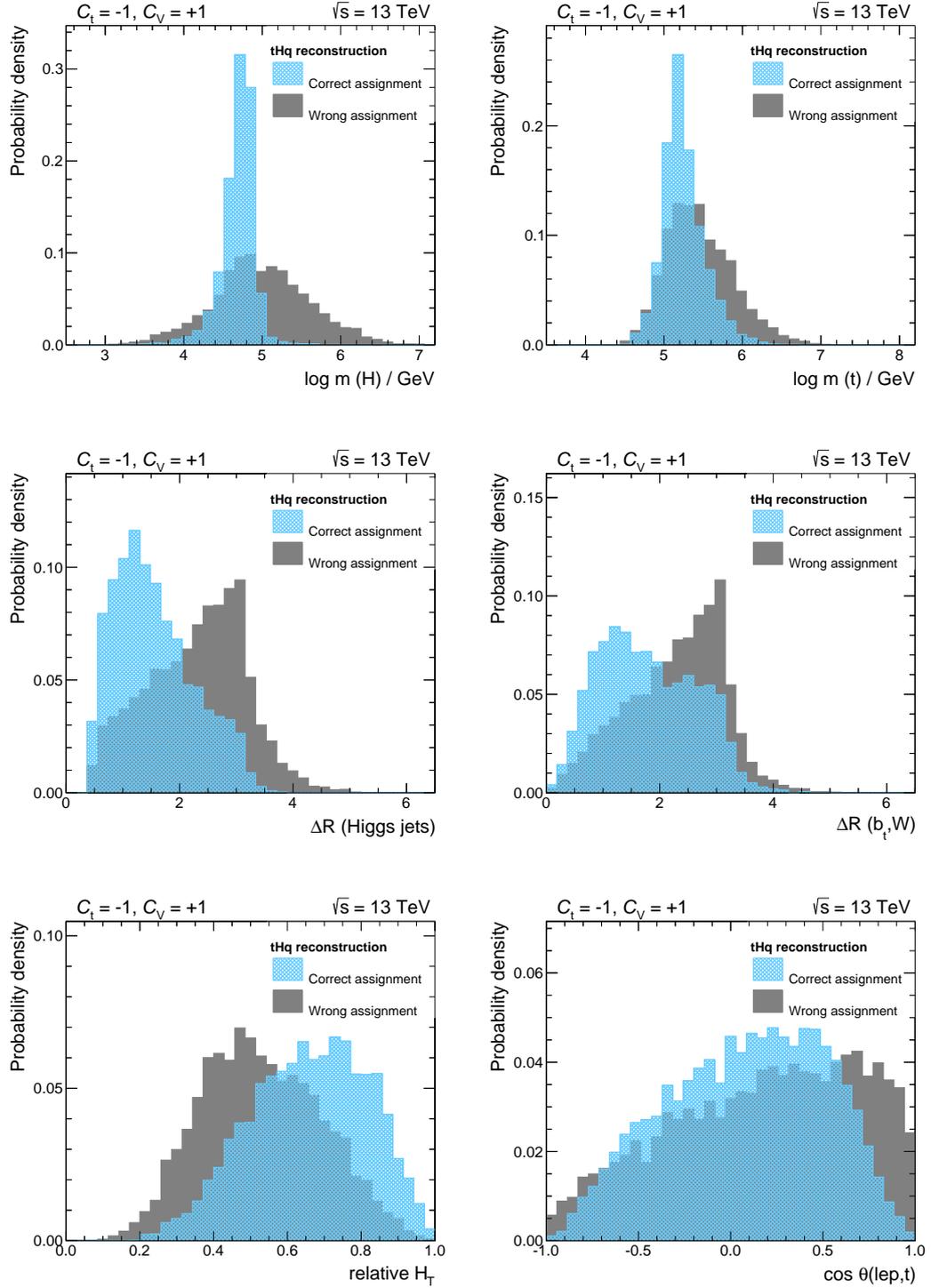


Figure 6.9.: The six most discriminating variables between correct and wrong jet assignments in the tHq reconstruction at $\sqrt{s} = 13$ TeV are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.2. The remaining variables can be found in Appendix B.1 and B.2.

Table 6.4.: Input variables for the jet-assignment BDT under the $t\bar{t}$ hypothesis sorted by their importance in the training. Instead of the transverse momenta variables the logarithm of these variables is used, as narrow distributions are better suited for the usage in MVA techniques than distributions with long tails.

Variable	Description
$\log \Delta m(t_{\text{had}}, W_{\text{had}})$	difference between the invariant masses of reconstructed t_{had} and W_{had}
$\log m(W_{\text{had}})$	invariant mass of the two jets assigned to the W boson of t_{had}
$\Delta R(b_{t_{\text{had}}}, W_{\text{had}})$	ΔR between the b quark of the reconstructed t_{had} and W_{had}
$\Delta R(W_{\text{had}})$	ΔR between the two jets assigned to the W boson of t_{had}
relative H_T	percentage of the total transverse momentum (jets, lepton, \cancel{E}_T) that falls to the reconstructed t_{had} and t_{lep}
$\Delta R(b_{t_{\text{lep}}}, W_{\text{lep}})$	ΔR between the b quark of the reconstructed t_{lep} and W_{lep}
$\log p_T(t_{\text{had}})$	transverse momentum of the reconstructed t_{had}
CSV(W_{had} jet 1)	CSV output of the first jet assigned to W_{had}
$\log m(t_{\text{lep}})$	invariant mass of the reconstructed t_{lep}
CSV(W_{had} jet 2)	CSV output of the second jet assigned to W_{had}
$\log p_T(t_{\text{lep}})$	transverse momentum of the reconstructed t_{lep}

6. Search for $t\bar{t}H$ Production at $\sqrt{s} = 13\text{ TeV}$

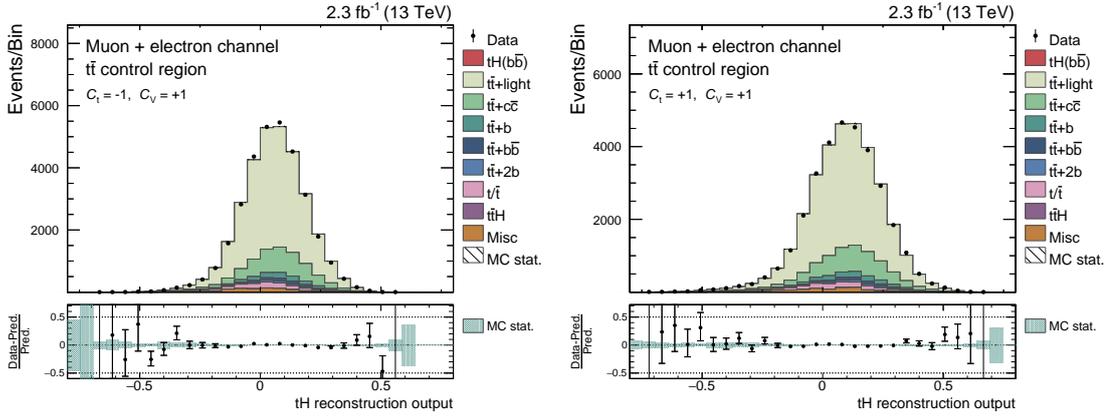


Figure 6.10.: Response of the tHq reconstruction comparing simulation to data. The highest output value (chosen jet assignment) per event is shown for the $C_t = -1$ point (left) and the SM prediction (right). Both diagrams are shown in the $t\bar{t}$ control region and a good agreement between simulation and data is observed. The simulation is scaled to match the event yields observed in data and all MC weights are applied. The corresponding distributions for the 3M and 4M region can be found in Appendix B.4.

forms clearly worse in the assignment for all central partons, outperforms the elaborate BDT reconstruction in the assignment of the light forward jet. The χ^2 reconstruction simply chooses the jet with the highest absolute pseudorapidity as light jet. This way the correct light jet is selected more often, but all correlations of the light quark to the other partons are ignored. By applying this procedure to background events the separation between signal events and background events actually becomes worse, as jets in the forward region that would have otherwise been ignored by the BDT reconstruction, are selected as light jet. A comparison of the two distributions, the absolute pseudorapidity of the jet chosen by the BDT as light jet and the highest absolute pseudorapidity of all jets in an event, can be seen for signal and background events in Figure 6.15. As the separation quality cannot be evaluated by a simple visual comparison, both variables have been separately used in the final classification, where the highest absolute pseudorapidity of an event showed a slightly worse behavior. As figure of merit the AUC was used and a small but noticeable reduction of 0.2% can be observed. Additionally, when both variables are used in the classification simultaneously, the BDT reconstructed pseudorapidity is chosen as the more important variable.

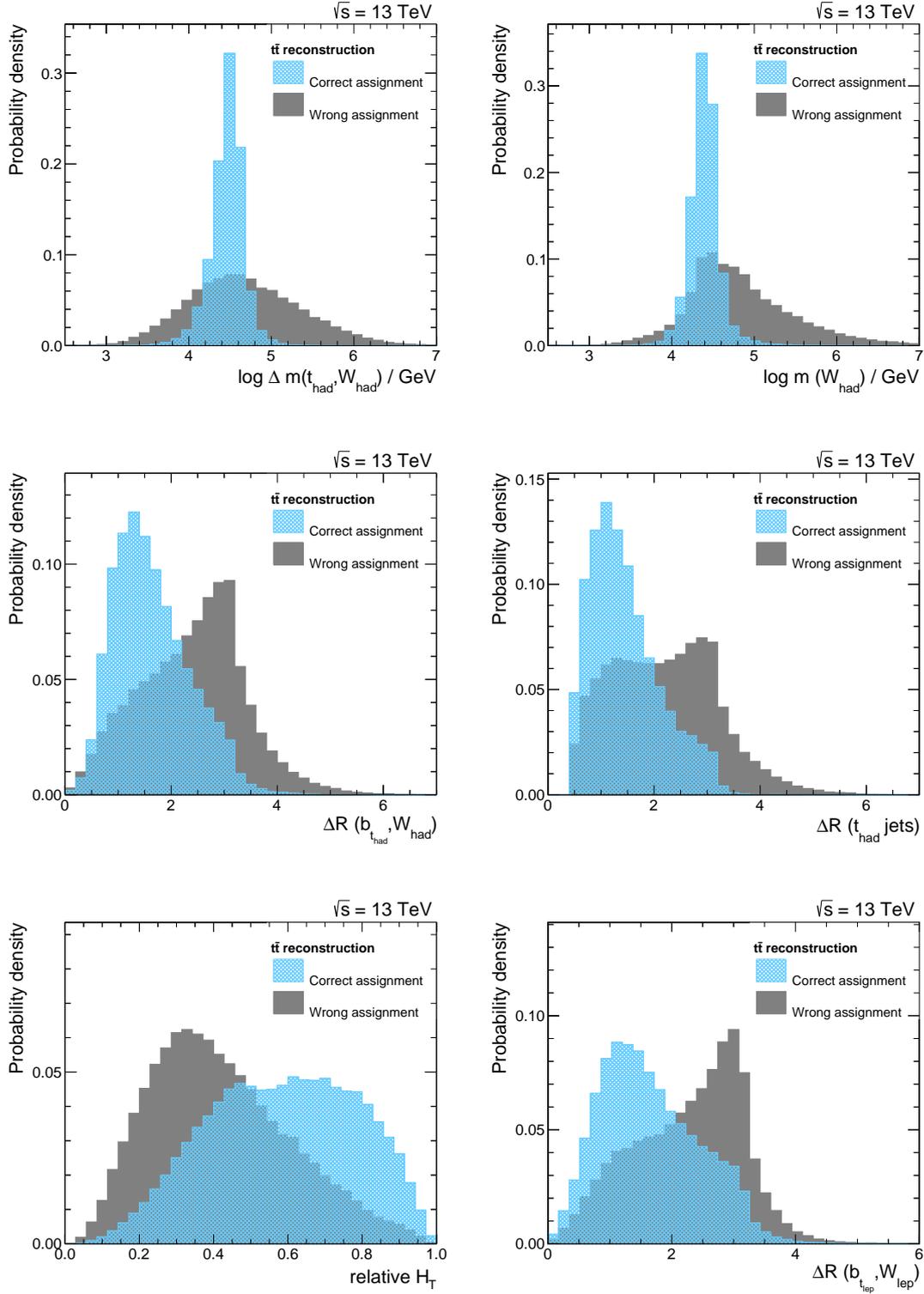


Figure 6.11.: The six most discriminating variables between correct and wrong jet assignments in the $t\bar{t}$ reconstruction at $\sqrt{s} = 13$ TeV are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.4. The remaining variables can be found in Appendix B.3.

6. Search for $t\bar{t}$ Production at $\sqrt{s} = 13$ TeV

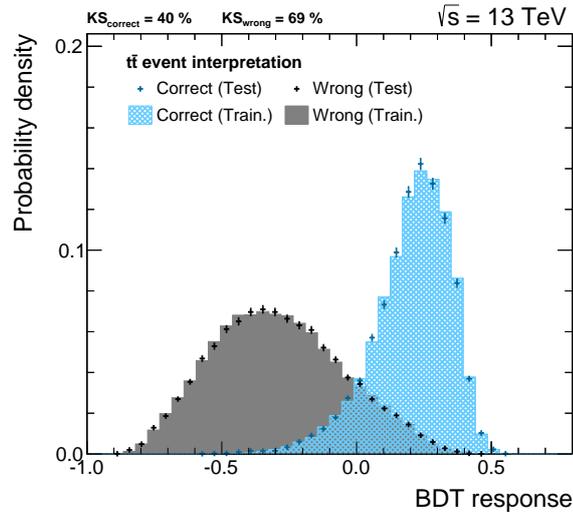


Figure 6.12.: Output values of the $t\bar{t}$ reconstruction BDTs for correct and wrong jet assignments. A clear separation is visible. The reconstruction training procedure is examined with an independent set of events that were not part of the training sample. Events used in the training are presented as histograms, whereas the independent testing sample is represented by markers. A good agreement between the testing and training sample is observed, verified by high KS-values of the two distributions. No sign of overtraining can be found.

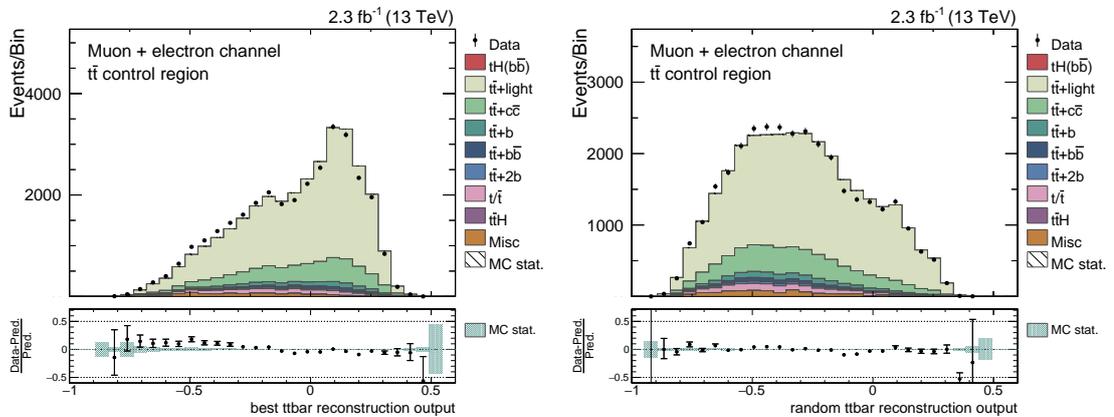


Figure 6.13.: Response of the $t\bar{t}$ reconstruction comparing simulation to data. On the left the highest output value (chosen jet assignment) per event is shown and on the right the BDT output value for a random assignment can be seen. Both diagrams are shown in the $t\bar{t}$ control region. In both distributions a reasonable agreement between simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied. The corresponding distributions in the 3M and 4M regions can be found in Appendix B.5.

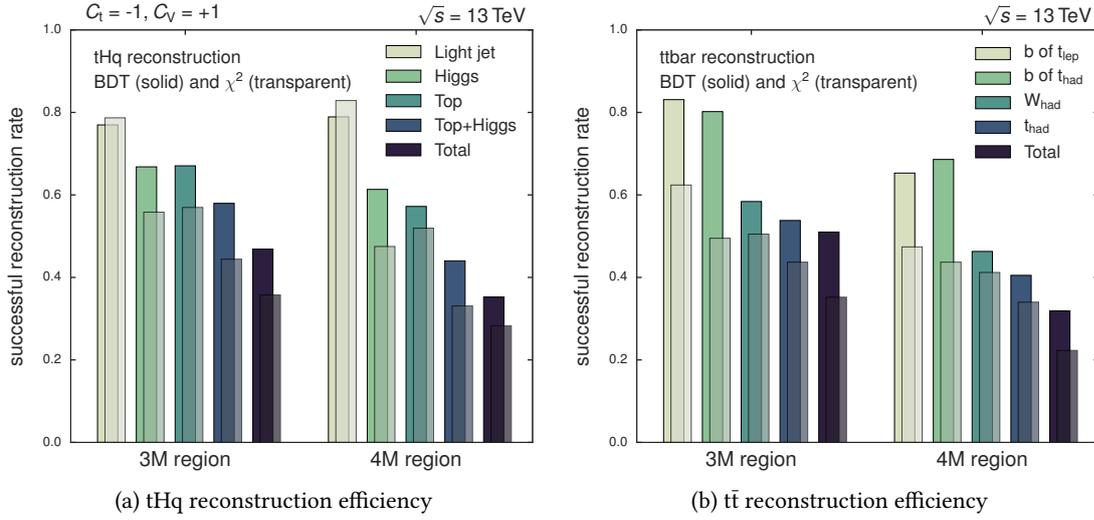


Figure 6.14.: Efficiencies of two reconstruction methods for the tHq hypothesis (left) and the $t\bar{t}$ hypothesis (right), one employing BDTs and the other using a simpler χ^2 approach, are shown in the 3M and 4M region. The bars represent the percentage that a certain object or the complete event is reconstructed correctly. The solid bars represent the reconstruction using BDTs for the assignment and the transparent bars represent the χ^2 reconstruction. The BDT reconstruction is clearly superior to the χ^2 method, as jets are assigned to their correct partons more frequently.

$t\bar{t}$ Reconstruction Evaluation

The performance of the $t\bar{t}$ reconstruction is evaluated as described earlier in Chapter 5.8.1. The respective rates to successfully assign jets to their original partons can be seen in Figure 6.14(b). The reconstruction based on MVA methods outperforms the χ^2 method for all individual partons. Whereas the assignment of the two light jets of W_{had} is slightly worse than in the $t\bar{t}$ reconstruction at $\sqrt{s} = 8$ TeV, the success rate of the jet assignment to bottom quarks of both top decays increased substantially. The probability to correctly assign jets to both b quarks has increased by 12.5% to 74.4% in the 3M region and by 15.3% to 52.6% in the 4M region. The culmination of these effects leads to a decrease of the successful total reconstruction rate in the 3M bin by 1.4% and an increase of 7.2% in the 4M bin. This will also be apparent in a gain in importance of the variables based on the $t\bar{t}$ reconstruction in the final classification, described in the next chapter.

6.10. Event Classification

After the two types of reconstructions every event is supplied with 51 jet assignments under the tHq hypothesis and one under the $t\bar{t}$ hypothesis. These reconstructed objects build the basis for most of the variables which are used in a multivariate classifier to separate signal events from background events. The classification has to be performed separately for each coupling point of the C_V - C_t plane, where for each point the identical variable set is used, but the observables

6. Search for tH Production at $\sqrt{s} = 13$ TeV

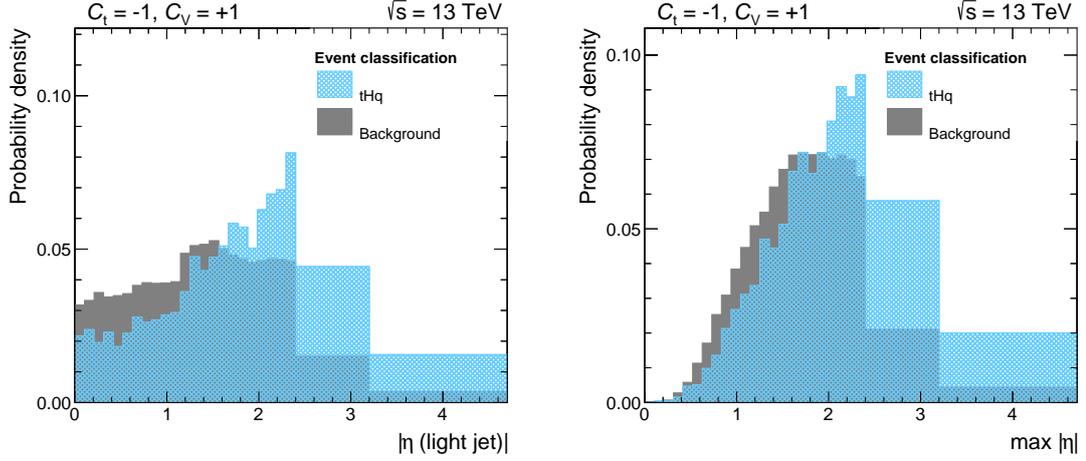


Figure 6.15.: On the left side the distribution of the absolute pseudorapidity of the jet assigned to the light quark by the tHq reconstruction is shown. On the right the distributions of the highest absolute pseudorapidity value of all jets in an event is shown. Both distributions are depicted for the signal process and the sum of background processes. The plots are shown for the $C_t = -1$ and $C_V = +1$ coupling pair and are normalized to unit area. Both variables show similar separation between signal and background. Although the choice of the most forward jet as light jet is more often correct than the jet chosen by the reconstruction, the variable obtained by the BDT reconstruction has a higher separation power and thus a larger impact in the final classification training.

based on reconstructed objects of the tHq process will differ for every studied coupling point. The classification BDTs are trained with the same parameters as the reconstruction BDTs, a setting which proved to be a good compromise between performance and robustness against overtraining. Similarly as in the tHq reconstruction, coupling pairs that correspond to the lowest cross sections and thus have fewer effective events in the training are prone to overtraining. Hence, these points are trained with a smaller number of trees (see Table 6.5) mitigating the overtraining effect.

The training is performed in the 3M region. Only the tHq process is used as signal input, as tHW is expected to share traits of the $t\bar{t}$ process and hence would decrease the separation power of the classifier. As background events a mixture of semi-leptonic $t\bar{t}$, full-leptonic $t\bar{t}$, $t\bar{t}H$ and single top events are used. The simulation samples of other processes do not contain enough events that they could be used in the training, as events used in the training are discarded and are not used further in the analysis. The background events are scaled to their predicted cross sections and the signal is scaled such that the integral coincides with the integral of the background events.

The set of classification variables has been optimized with regard to the analysis at $\sqrt{s} = 8$ TeV. A total of 15 different variables is used, four variables based on information gained from the $t\bar{t}$ reconstruction, seven variables based on information gained from the tHq reconstruction and four reconstruction-independent variables. A description of the variables and their averaged importance over all 51 trainings can be found in Table 6.6. The importance of the variables is quantified in the same way as in the tHq reconstruction described in Section 6.9.1.

Table 6.5.: Parameter settings used for the BDTs employed in the classification. A smaller number of trees is used for points in the C_V - C_t plane that have otherwise shown signs of overtraining. For $C_V = +1.5$ these include all points with $C_t \geq 0.5$, for $C_V = +1$ all points satisfying $C_t \geq 0.75$ and for points with $C_V = +0.5$ the tree number is reduced, if $0.25 \leq C_t \leq 2.0$. These points correspond to the points with the lowest cross sections and due to the reweighting lead to a smaller number of effective signal events in the training of the classification. After this reduction no sign of overtraining is found. The definitions of the configuration options can be found in Reference [117].

Parameter	Value
NTrees	400/100
MinNodeSize	1
MaxDepth	3
BoosteType	AdaBoost
nCuts	20
AdaBoostBeta	0.3
SeparationType	GiniIndex

The applied CSV reweighting procedure allows for the inclusion of the important CSV output distributions. Whereas in the predecessor analysis only variables could be used which provided information about the number of b-tagged jets assigned to a reconstructed object the complete shape of the CSVv2 discriminator for certain jets can be exploited in this iteration. The discriminator shapes for the two jets assigned to the two light quarks of the hadronically decaying W boson and for the two jets assigned to the Higgs boson decay products are exploited, out of which the CSV output for the light jets of the W_{had} decay are found to be very important in the training. When applying the $t\bar{t}$ reconstruction to a signal event in the signal region in most cases a b-tagged jet has to be assigned to the decay product of W_{had} due to missing alternatives, which leads to a good separation power of this variable. Two novel variables, which are independent of the reconstruction, also rank highly in the training: the aplanarity of the event, which contains information about the geometrical shape of the event in general (see Reference [186] for further details) and the variable m_3 , which corresponds to the invariant mass of the three hardest jets in the event, add valuable information to the classification. Another novel reconstruction-independent variable is the second Fox-Wolfram moment of the event.¹ The six most important variables are shown in Figure 6.18 for signal and background events, the remaining variables can be found in Appendix B.6. A good modeling of the employed variables by the simulation is validated by comparing simulation to data. The resulting distribution can be found in Figures 6.19, 6.20 and 6.21 for the $t\bar{t}$ control region, 3M and 4M region, respectively. The remaining comparisons can be found in the Appendix B.7-B.12.

Each of the 51 trainings is evaluated and checked for possible signs of overtraining. If overtraining occurred, the number of trees in the training is reduced until no overtraining is observed anymore. The response of the training sample and of a disjoint testing sample for the coupling point of $C_t = -1$ and $C_V = +1$ and for the couplings predicted by the standard model can be found in Figure 6.17. The area under the ROC curve for all 51 points is visualized in Figure 6.16.

¹Further information about Fox-Wolfram moments can be found in Reference [187].

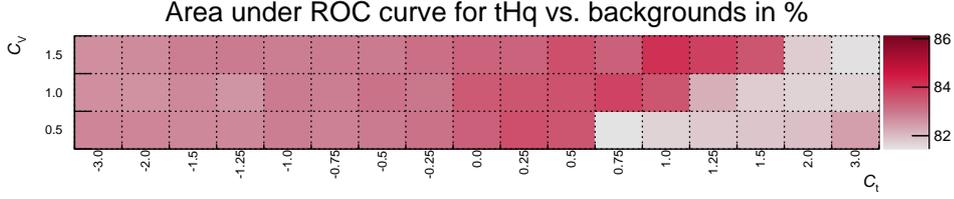


Figure 6.16.: The area under the ROC curve for all 51 classification trainings. With values above 80% good performance is observed for all trainings. A visible ridge on the right side of the plot corresponds to the points with the lowest cross sections. Effectively lower event counts lead to a deterioration of the separation power between signal and background events. This effect is enhanced as most of these points with low cross section had to be trained with smaller number of trees in the reconstruction and classification.

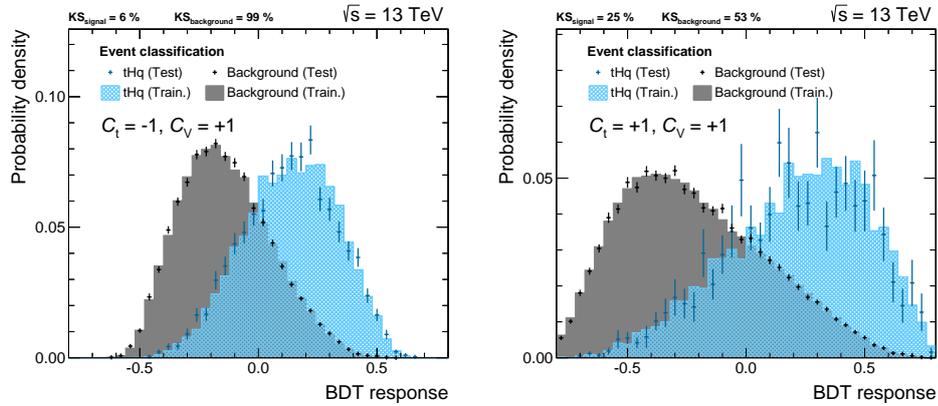


Figure 6.17.: Output values of the classification BDTs for the coupling case of $C_t = -1$ and $C_V = +1$ (left) and the SM prediction (right). A clear separation between signal and background is visible for both distributions. The classification is examined with an independent set of events. The events used in the training are presented as histograms, whereas the independent testing sample is represented by markers. A good agreement between the testing and training sample is seen, verified by high KS-values of the two distributions. No sign of overtraining is found.

The BDT responses show a clear separation between signal events and background events and the testing sample reproduces the shape of the output, hence giving confidence in the training of the classifier.

Subsequently, the classifier is applied to simulation samples and data samples alike. As the input variables already are well modeled in the simulation a good agreement between data and simulation in the classifier output is expected. The distributions in Figure 6.22, which show the classifier output for the $C_t = -1$ and $C_V = +1$ and the standard model scenario in the $t\bar{t}$ control region, confirm this expectation, as both classification outputs agree well for data and simulation. The different shape of the BDT output for the two depicted coupling points is a direct consequence of the reduced number of trees in the classification training for the point predicted by the standard model.

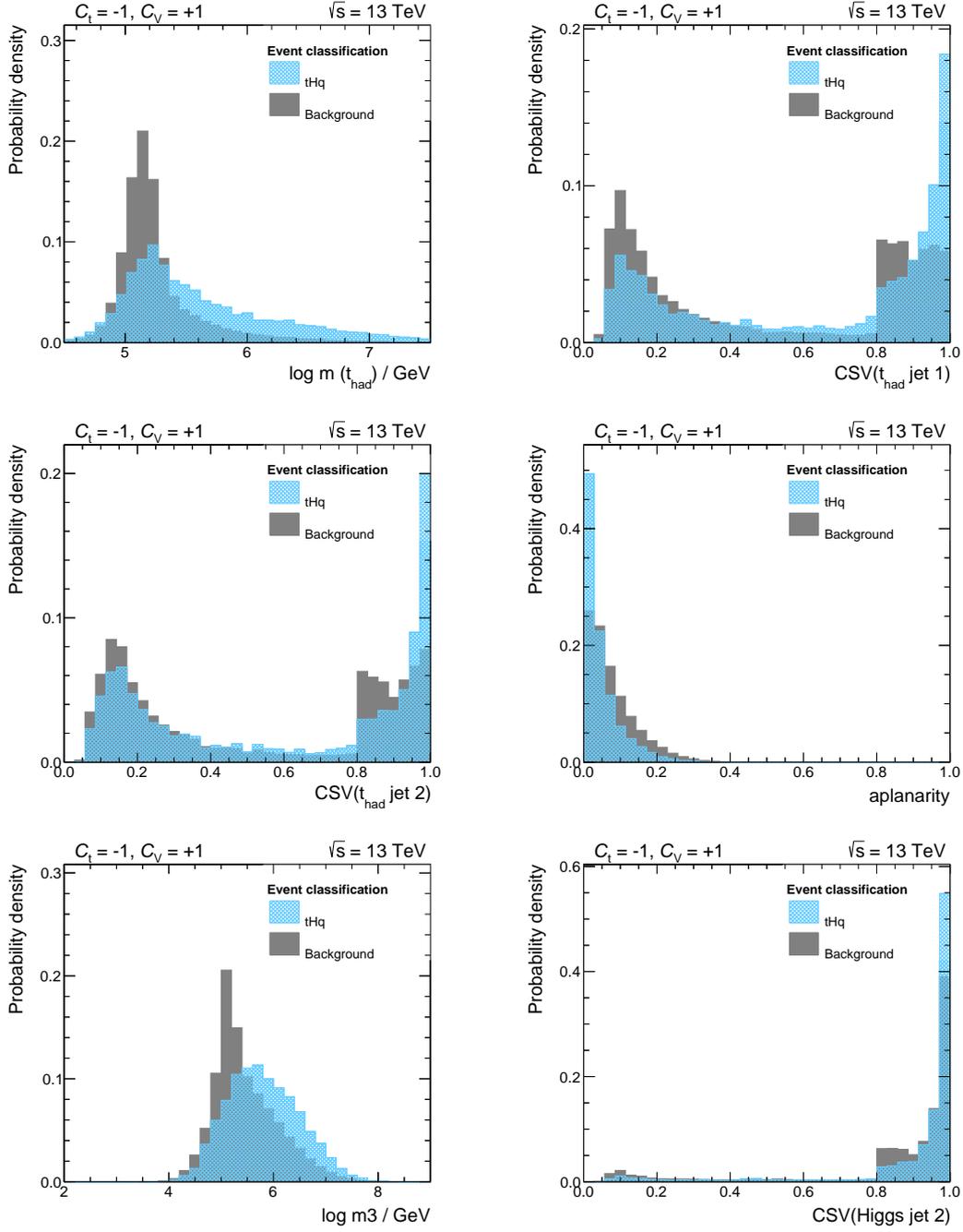


Figure 6.18.: Six of the most important variables used in the final classification sorted by their importance in the training at $\sqrt{s} = 13$ TeV. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.6. The remaining variables can be found in Appendix B.6. The third most important variable, the pseudorapidity of the light forward jet, is omitted here, as it is already shown in Figure 6.15.

6. Search for tH Production at $\sqrt{s} = 13$ TeV

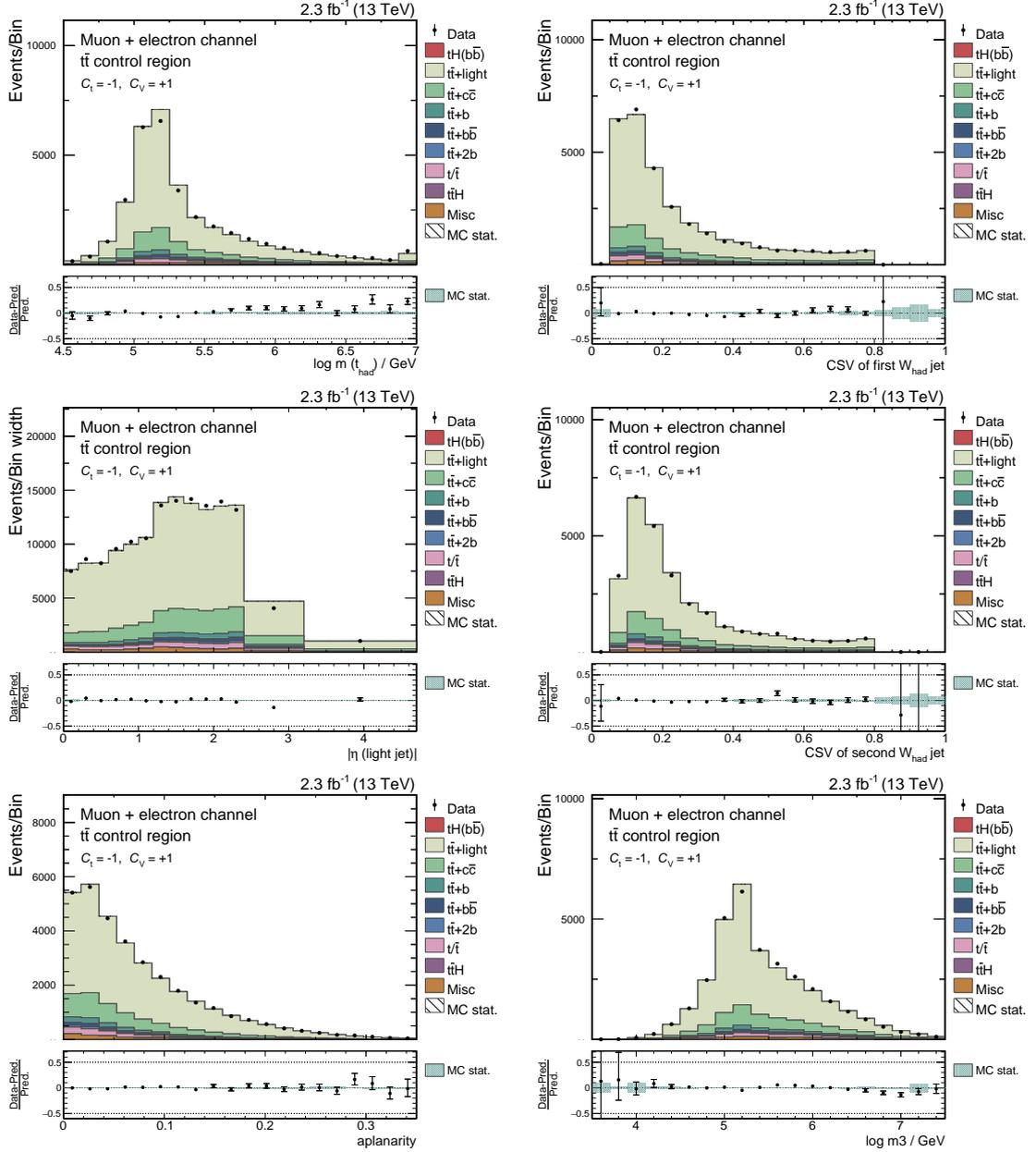


Figure 6.19.: Simulation to data comparisons for the six most important input variables of the classification at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the $t\bar{t}$ control region for the coupling case of $C_t = -1$ and $C_V = +1$. Besides a small normalization offset a good agreement of simulation and data is found. In all diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied. The distributions of the remaining variables of the classification can be found in Appendix B.7 and B.8.

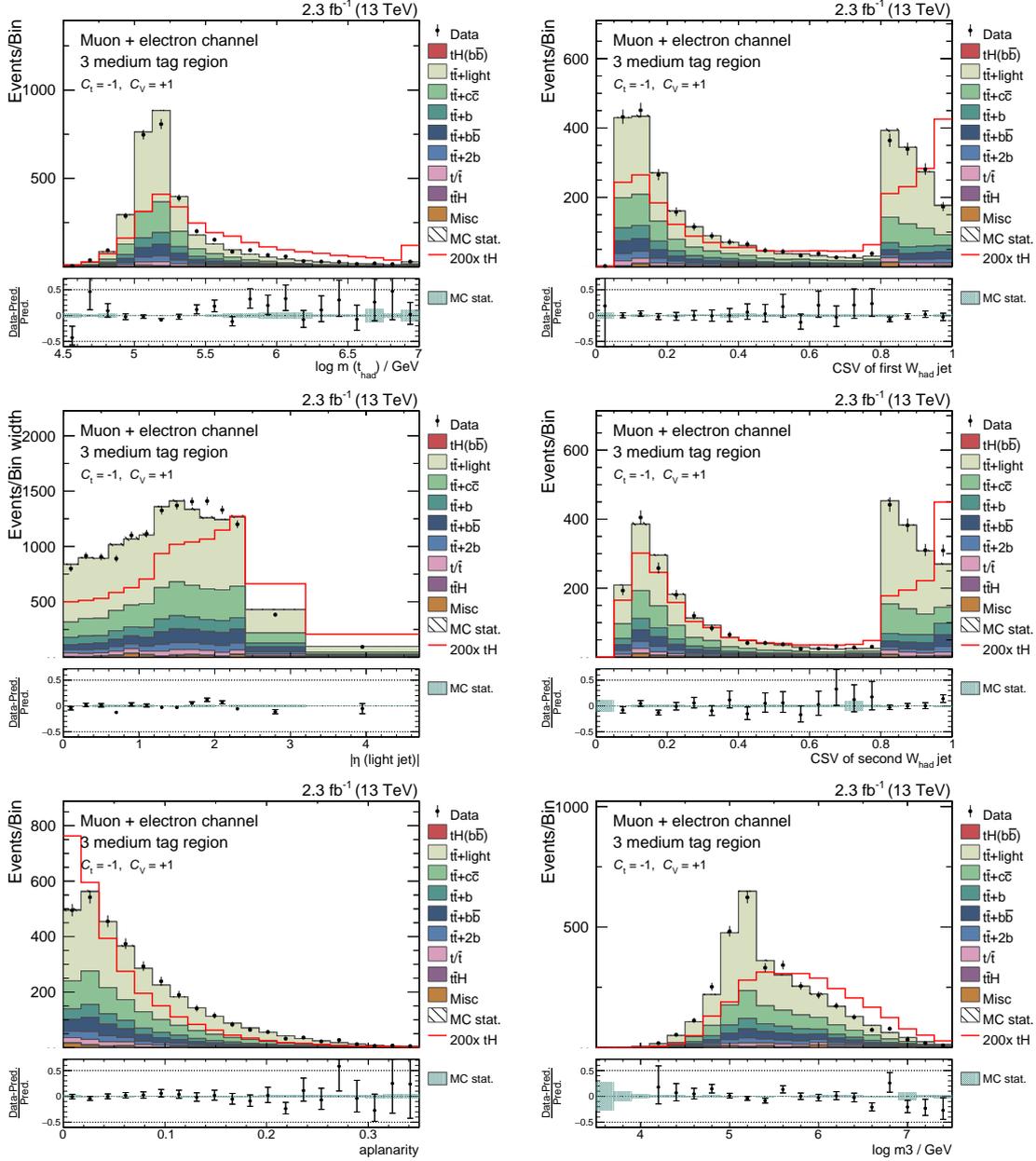


Figure 6.20.: Simulation to data comparisons for the six most important input variables of the classification at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the 3M region for the coupling case of $C_t = -1$ and $C_V = +1$. Besides a small normalization offset a good agreement of simulation and data is found. In all diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied. The distributions of the remaining variables of the classification can be found in Appendix B.9 and B.10.

6. Search for tH Production at $\sqrt{s} = 13$ TeV

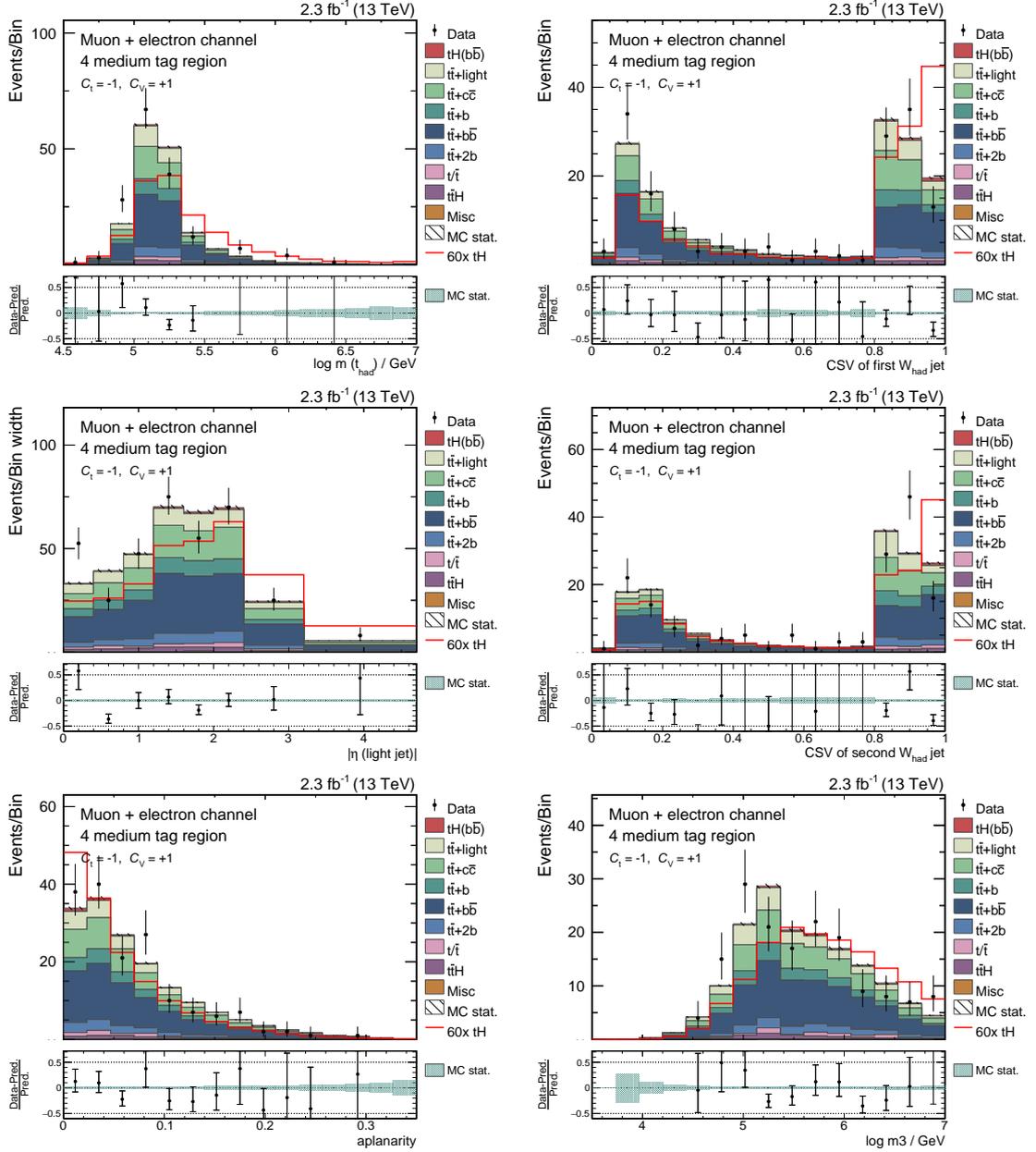


Figure 6.21.: Simulation to data comparisons for the six most important input variables of the classification at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the 4M region for the coupling case of $C_t = -1$ and $C_V = +1$. A good agreement of simulation and data is found. In all diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied. The distributions of the remaining variables of the classification can be found in Appendix B.11 and B.12.

Table 6.6.: Description of variables used in the classification and their importance ranking in the training. The variables are grouped into three categories: variables independent of any reconstruction (top), variables based on objects reconstructed under the $t\bar{t}$ hypothesis (center) and variables based on objects reconstructed under the tHq hypothesis (bottom). Instead of the transverse momenta the logarithm of these variables is used, as narrow distributions are better suited for the usage in MVA techniques than distributions with long tails.

Variable	Points	Description
aplanarity	9.62	aplanarity of the event
log m3	9.52	invariant mass of the three hardest jets in the event
Fox-Wolfram #1	2.70	first Fox-Wolfram moment of the event
q(ℓ)	2.29	electric charge of the lepton
log m(t_{had})	13.90	invariant mass of t_{had}
CSV(W_{had} jet 1)	12.03	CSV output of the first jet assigned to W_{had}
CSV(W_{had} jet 2)	9.68	CSV output of the second jet assigned to W_{had}
ΔR (W_{had} jets)	7.37	ΔR between the two light jets from the W_{had} decay
$\eta(\text{light jet})$	11.21	absolute pseudorapidity of the light forward jet
CSV(Higgs jet 2)	7.50	CSV output of the second jet assigned to the Higgs boson
CSV(Higgs jet 1)	5.66	CSV output of the first jet assigned to the Higgs boson
log p_T (light jet)	5.60	transverse momentum of the light forward jet
log p_T (Higgs)	5.33	transverse momentum of the Higgs boson
$\eta(t) - \eta(H)$	2.21	absolute difference of pseudorapidities of the reconstructed top quark and the reconstructed Higgs boson
cos $\theta(t, \ell)$	0.29	Cosine of the angle from the top quark vector to the sum vector of top quark and charged lepton in their common restframe

6. Search for tH Production at $\sqrt{s} = 13$ TeV

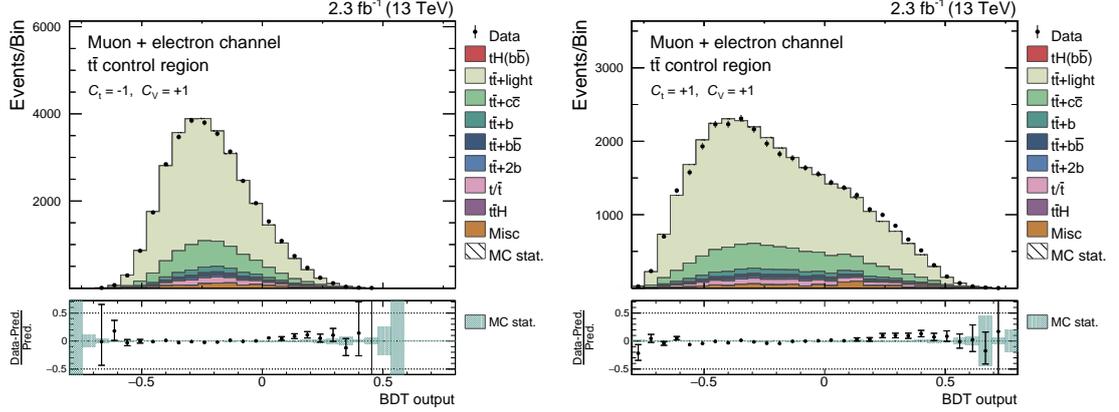


Figure 6.22.: Simulation to data comparisons for the output of the classification in the $t\bar{t}$ control region at $\sqrt{s} = 13$ TeV. The training output is shown for the $C_t = -1$ and $C_V = +1$ coupling point (left) and the coupling point predicted by the SM (right). In both distributions a good agreement between simulation and data is observed. The simulation is scaled to match the event yields observed in data and all MC weights are applied.

6.11. Systematic Uncertainties

A careful treatment of the systematic uncertainties applied to the analysis is necessary in order to lend substance to the final result. In comparison to the analysis at $\sqrt{s} = 8$ TeV, a similar set of uncertainties is employed.

Experimental Uncertainties

Luminosity (rate) The result of the most recent luminosity measurement at the time of writing is used in this analysis. By using the pixel cluster counting method and determining the absolute luminosity scale with Van der Meer scans it was possible to reduce the overall uncertainty on the luminosity measurement at $\sqrt{s} = 13$ TeV to 2.7% [188]. This uncertainty is applied to all processes in both signal regions.

Lepton Efficiencies (rate) As described in Section 6.7.2, lepton efficiencies have to be corrected for several different reasons, like trigger efficiencies or reconstruction efficiencies. To account for uncertainties in the estimation of these efficiencies a conservative overall uncertainty of 2% is applied to all processes.

Pileup (shape) The reweighting procedure implemented to reproduce the distribution of the number of reconstructed primary vertices is considered as source for a systematic uncertainty. The shape variation is evaluated by applying only 50% of the shift as down variation and 200% as an up variation to each event.

Jet Energy Resolution (shape and rate) The uncertainty covering the jet energy resolution is evaluated by increasing and decreasing the difference between reconstructed jet

energy and true jet energy on particle level. The scale factors and their uncertainties used in the smearing are provided in Reference [178]. The complete analysis chain is reiterated for samples with systematically changed jet energy resolution.

Jet Energy Scale (shape and rate) The applied jet energy corrections are varied within their provided uncertainties [155] and the complete analysis is repeated for the systematically shifted samples.

CSV Reweighting (shape and rate) The CSV reweighting procedure, which is described in Section 6.7.3, considers different uncertainty sources which are all treated separately. When changing the jet energy scale according to its uncertainties the change of the b-tagging scale factors is evaluated and taken as 100% correlated to the shift of the energy scale. Another uncertainty source is the purity of the sample from which the scale factors were derived. The third source, the impact of statistical uncertainties during the scale factor determination, is propagated to an alternative set of scale factors. The statistical impact is described by two different nuisance parameters, which both have a certain degree of control over distortions in the CSV distribution. All of the above described uncertainties are taken into account separately for heavy flavor and light flavor jets and are taken as fully uncorrelated. Additionally, two sets of weights are applied that change the contamination of charm jets in the samples used for the scale factor determination. More information on the uncertainties considered in the scale factor determination can be found in Reference [184]. The combined effect of these uncertainties on the final discriminator shape can be seen in Figure 6.23(a).

Theoretical Uncertainties

$t\bar{t}$ + Heavy Flavor Rates (rate) Similar as in the 8 TeV analysis an uncertainty of 50% is assigned to the $t\bar{t}+b$, $t\bar{t}+2b$, $t\bar{t}+b\bar{b}$ and $t\bar{t}+c\bar{c}$ samples.

PDFs/ Q^2 Scale (rate) Uncertainties applied to different simulation samples affecting the normalization of the different processes based on the choice of PDF set and the Q^2 scale can be found in Table 6.7. Uncertainties for processes with a common production mechanism are treated as fully correlated.

The treatment of the Q^2 scale uncertainty has been revised with respect to the analysis at 8 TeV as samples can be reweighted by employing the LHE reweighting to emulate a sample produced with a modified Q^2 scale.

These weights are stored in most of the officially produced samples, but are unfortunately missing for the used single top and diboson simulation samples. Therefore, the single top process is assigned with a 4.0% rate uncertainty and the diboson process is assigned with a 2.5% uncertainty to cover the effect of the variation of the Q^2 scale.

Q^2 scale (shape and rate) By utilizing the LHE reweighting procedure, events in a simulation sample can be reweighted such that they emulate a sample produced with a different Q^2 scale. For all simulation samples, but the single top and diboson samples, these weights are available, and for each process an uncorrelated Q^2 scale uncertainty is introduced by reweighting the events in the final classification output. The reweighted samples correspond

6. Search for tH Production at $\sqrt{s} = 13$ TeV

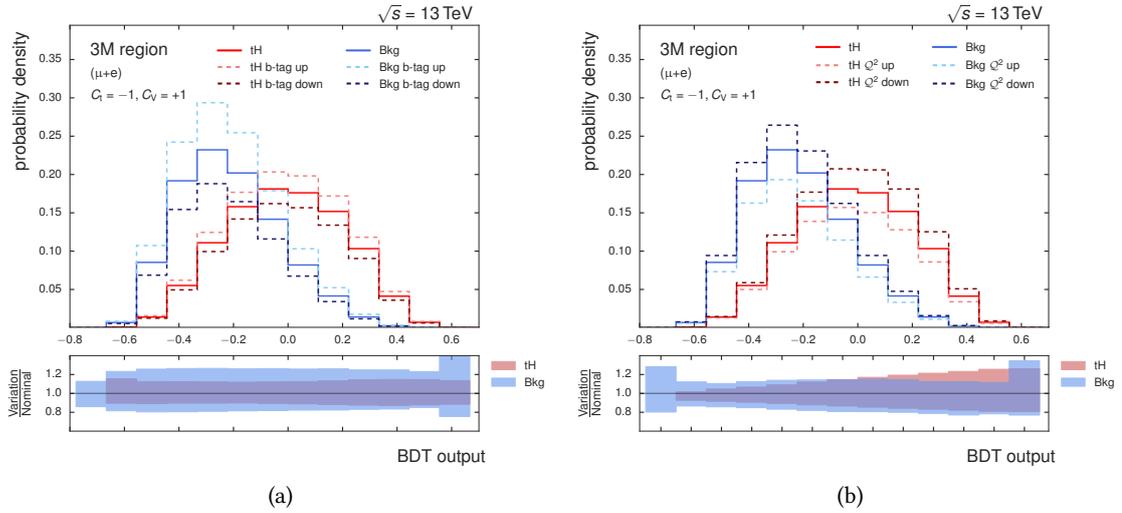


Figure 6.23.: Effect of the CSV reweighting uncertainty (left) and Q^2 scale uncertainties (right) on the final classifier for the background and signal shape in the 3M region are shown. The displayed uncertainties correspond squared sums of the contributing uncertainty sources. The signal and background histograms are normalized to unity, the systematically shifted templates are allowed to change shape as well as the normalization. Whereas the b-tagging uncertainty is mostly flat over the complete BDT output range, thus resulting in a change of normalization, Q^2 uncertainties grow for larger BDT output values.

to a Q^2 scale of fourfold and quarter of the initial value, respectively. The implementation of this uncertainty has a large effect on the shape of the classification output as well as on the normalization of the individual processes. The effect of this uncertainty on the final discriminator shape can be seen in Figure 6.23(b).

Statistical Uncertainties

Bin-by-bin uncertainties (shape) The bin-by-bin uncertainties are implemented as in the 8 TeV analysis, described in Section 5.10.

Table 6.7.: Cross section uncertainties based on the choice of PDF set and Q^2 scale applied for the different processes. The uncertainty values are obtained from References [37, 189–192].

Process	PDF (%)		Q^2 scale (%)	
	gg	q \bar{q}	qg	
tHq			3.7	
tHW			4.0	
t \bar{t} H	3.6			
t \bar{t}	3.0			
t \bar{t} W		2.0		
Single top			4.0	4.0
W+jets		4.0		
Diboson		2.0		2.5

6.12. Results

6.12.1. Fit of Final Discriminator

A simultaneous MLE fit is performed in the 3M and the 4M region with help of the `COMBINE` package for each of the studied points of the C_V - C_t parameter plane. The prefit and postfit distributions of both channels with full consideration of statistical and systematic uncertainties for the $C_t = -1$ and $C_V = +1$ coupling point can be seen in Figure 6.24.

6.12.2. Analysis of Nuisance Parameters

The values and uncertainties of all nuisance parameters are studied at the prefit and the postfit stage. A visualization of the behavior of all parameters sans the bin-by-bin uncertainties in an s+b fit can be found in Figure 6.25. It is visible that most of the nuisance parameters stay close to their initial values and their uncertainties are, if at all, only slightly decreased.

Already lightly constrained at the prefit stage and also pulled to its $+1\sigma$ boundary is the uncertainty covering the c quark treatment in the CSV reweighting. The 1σ upward fluctuation corresponds to a significant overall reduction of the background. This effect seems to be linked to an overestimation of the uncertainty to begin with as other analyses observe the same behavior [142].

The t \bar{t} backgrounds are affected by multiple different effects: The 0.5σ upwards fluctuation of the Q^2 systematic uncertainty assigned to the t \bar{t} +light background corresponds to a reduction of the normalization, whereas the t \bar{t} components containing b quarks are scaled upwards via their overall rate uncertainties. Overall the set of pulls and constraints do not give rise to concerns and increase the trust in the performance of the fit.

The effect of the single systematic uncertainty groups has been evaluated by checking how the removal or the exclusive usage of this systematic uncertainty affects the expected limit. The

6. Search for tH Production at $\sqrt{s} = 13$ TeV

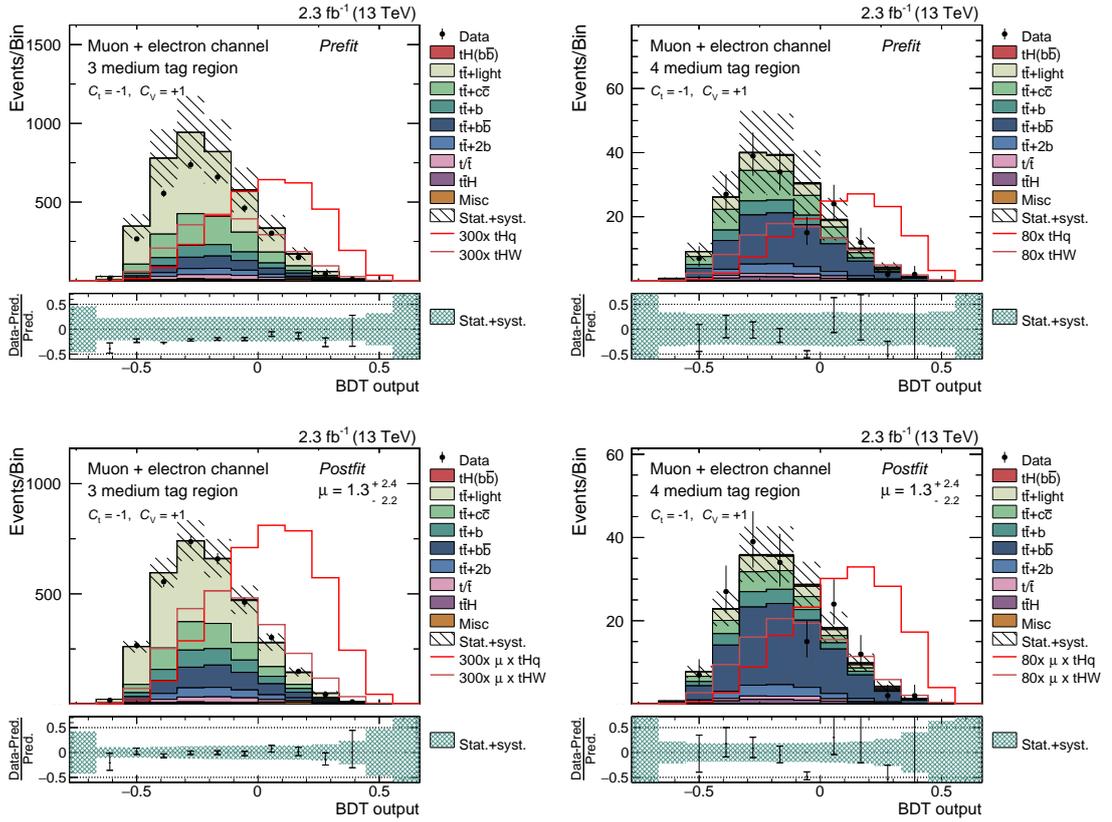


Figure 6.24.: Pre- and postfit distributions of the classifier output in the 3M and 4M region for the coupling pair of $C_t = -1$ and $C_V = +1$. These distributions are fitted simultaneously and a signal strength factor of $\mu = 1.3^{+2.4}_{-2.2}$ has been determined. A good agreement of simulation and data is observed after the fit. Pre- and postfit distributions for coupling parameters predicted by the SM can be found in Appendix B.13.

Table 6.8.: Pre- and postfit yields in the two simultaneously fitted signal regions. The quoted uncertainties include systematic and statistical components. Additionally, the sum of all expected background events and the observed number of events in data are quoted.

	3M		4M	
	Prefit	Postfit	Prefit	Postfit
$t\bar{t}$ +light	2185 \pm 809	1557 \pm 241	22.8 \pm 22.2	16.0 \pm 13.4
$t\bar{t}$ + $c\bar{c}$	827 \pm 477	530 \pm 223	39.6 \pm 34.3	19.8 \pm 9.6
$t\bar{t}$ +b	329 \pm 165	338 \pm 165	19.2 \pm 11.1	16.4 \pm 8.0
$t\bar{t}$ + $b\bar{b}$	337 \pm 157	437 \pm 114	72.3 \pm 38.2	89.0 \pm 23.4
$t\bar{t}$ +2b	179 \pm 118	195 \pm 111	13.5 \pm 10.7	12.1 \pm 7.9
Single Top	132 \pm 35	108 \pm 16	5.3 \pm 1.7	4.2 \pm 0.7
$t\bar{t}$ H	20.6 \pm 8.2	19.9 \pm 8.8	5.4 \pm 2.4	5.0 \pm 2.3
Diboson	1.9 \pm 1.0	1.2 \pm 0.5	0.1 \pm 0.0	0.1 \pm 0.1
W+jets	35.0 \pm 18.4	28.9 \pm 11.2	0.0 \pm 0.0	0.0 \pm 0.0
$t\bar{t}$ W	8.0 \pm 2.5	6.3 \pm 1.3	0.4 \pm 0.2	0.3 \pm 0.1
Σ Backgrounds	4054 \pm 975	3220 \pm 401	179 \pm 58	163 \pm 31
tHq	10.9 \pm 3.9	13.6 \pm 31.3	1.7 \pm 0.7	2.1 \pm 4.6
tHW	6.8 \pm 1.0	8.3 \pm 16.8	1.1 \pm 0.2	1.3 \pm 2.6
Observed	3199		162	

study is performed analogous to the study at 8 TeV, described in Section 5.11.2 and the outcome can be seen in Figure 6.26.

The Q^2 scale uncertainty and the b-tagging uncertainty have the largest impact on the limit. The inclusion of the CSV reweighting, which allows for the incorporation of the shape of the CSV output into the analysis, incurs large uncertainties, which are more than offset by the gain provided by them. The study shows that only a combined effort to reduce systematic uncertainties on all ends will improve the analysis. The reduction of only one uncertainty will help slightly, but other uncertainties will compensate the effect keeping the sensitivity limited.

The inclusion of the 50% rate uncertainties for the $t\bar{t}$ +heavy processes is expected to have a large impact on the limit, which is also shown by its large effect as sole uncertainty, but the impact of its removal from the complete set is almost negligible as other uncertainties cover the same uncertainty range.

The smaller rate uncertainties like the lepton efficiency or the luminosity uncertainty have almost no effect on the limit. The pileup systematic uncertainty also showed an impact close to zero, which is reasonable as the pileup scenario in MC was already very close to what is observed in data, and the reweighting corresponds only to a small correction. The fact that systematic uncertainties seem to show zero impact is owed to the application of the asymptotic method, when calculating the impact and the initial limit. The approximation used in the asymptotic formula causes the tiny impact of the removal of the lepton efficiency and the vanishing of luminosity and pileup uncertainty.

6. Search for tH Production at $\sqrt{s} = 13$ TeV

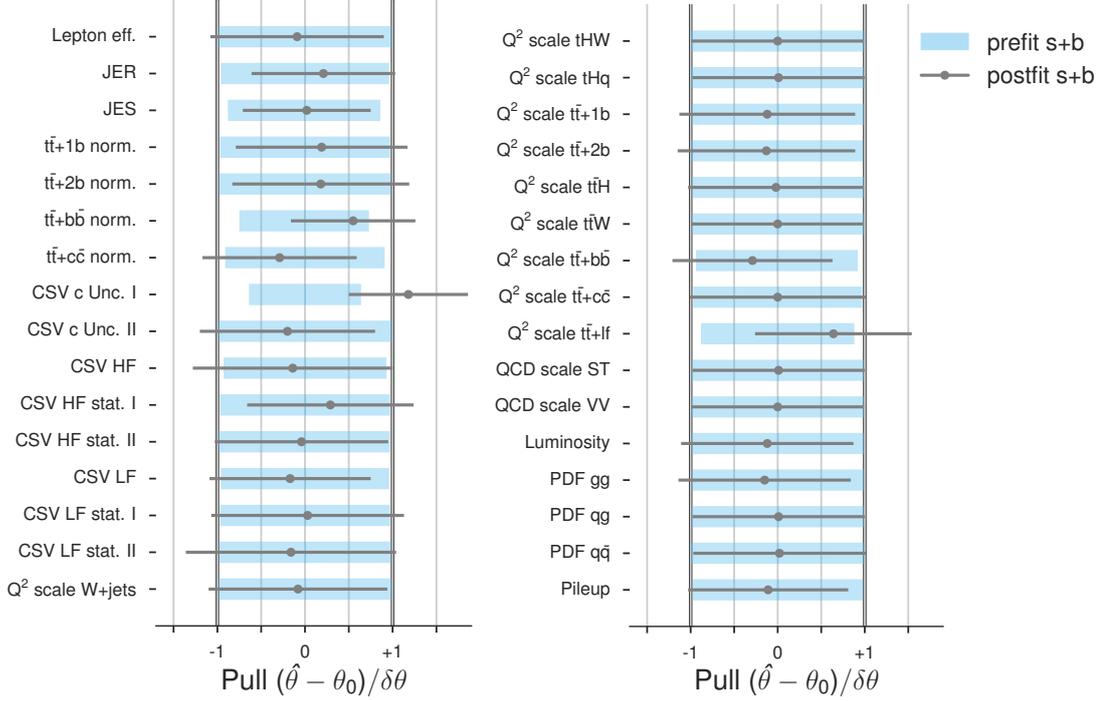


Figure 6.25.: Pre- and postfit pulls of all nuisance parameters of the analysis in an s+b fit (bin-by-bin uncertainties excluded).

6.12.3. CL_S Limits

Asymptotic limits at 95% confidence level are calculated for each of the 51 points in the C_V - C_t parameter plane. The expected and observed limits as a function of C_t for each of the three considered C_V values can be found in Figure 6.27. A smoothing of the limit bands is performed with help of a cubic spline fit interpolating between the actual calculated limit points. The upper limit values for the $C_t = -1$ and $C_V = +1$ coupling and for the standard model case can be found in Table 6.9. The remaining values can be found in the Appendix B.5. For the coupling point of $C_t = -1$ and $C_V = +1$ an upper limit of $\mu_{\text{obs}}^{\text{up}} = 7.4$ is observed, whereas $\mu_{\text{obs}}^{\text{up}} = 5.7$ is expected. This leads to an exclusion of a process with kinematics like the studied signal model for $C_t = -1$ with a cross section of $\mu_{\text{obs}}^{\text{up}}(C_t = -1, C_V = +1) \cdot \sigma_{\text{tHq+tHW}}(C_t = -1, C_V = +1) = 7.0$ pb. For the standard model case a production cross section of 9.3 pb can be excluded. In a first approximation the expected limits are anti-proportional to the cross sections of the sought signal process. On a closer look, small deviations of the maximal expected limit from the point of the lowest cross section can be noticed, which is caused by slightly differing cross section minima of the tHq and tHW production. The analysis is optimized towards the tHq process and therefore the limit is shifted closer towards the actual tHq cross section minimum.

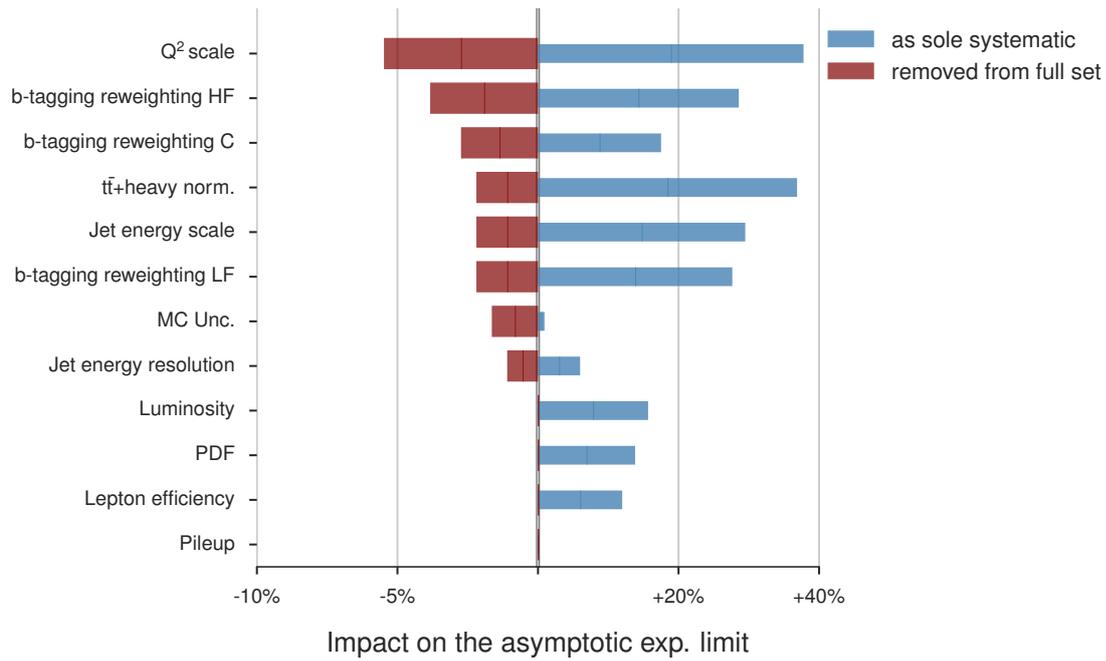


Figure 6.26.: Impact of groups of systematic uncertainties on the expected asymptotic limit. The groups of systematic uncertainties are either removed from the fit by fixing them to their postfit value, or used as single systematic uncertainty by fixing all other uncertainties to their postfit values. The changes displayed in this diagram are calculated relatively to the limit with all systematic uncertainties included (red bars) and to the limit, where all uncertainties are fixed to their best fit value (blue bars).

6. Search for tH Production at $\sqrt{s} = 13$ TeV

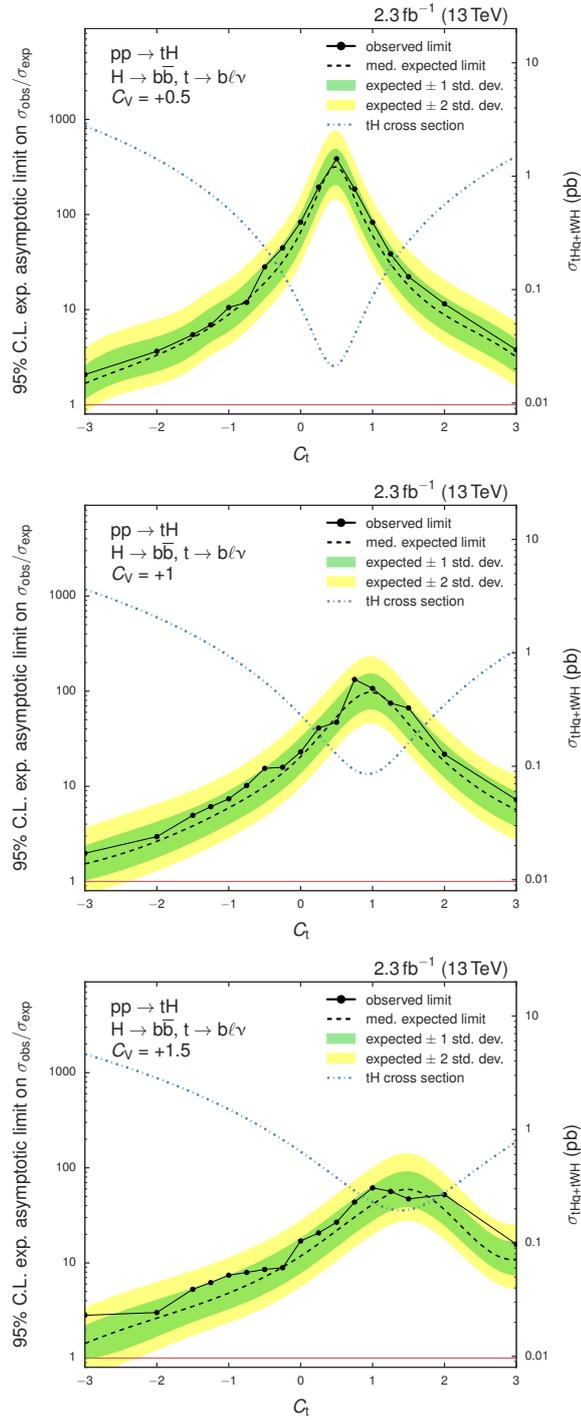


Figure 6.27.: Expected and observed asymptotic CL_S limits at 95% C.L. for the combination of 3M and 4M region as a function of C_t . The limits are shown for $C_V = +0.5$ (top), $C_V = +1.0$ (center) and $C_V = +1.5$ (bottom). The sum of the tHq and tHw cross sections as function of C_t is depicted as dotted blue line for each of the three C_V values. Neither an excess nor a strong downward fluctuation is observed. The corresponding limit values for all studied points can be found in the Appendix in Table B.5.

Table 6.9.: Expected and observed asymptotic CL_S limits at 95% C.L. in the 3M and 4M region and their combination for the coupling pair of $C_t = -1$ and $C_V = +1$ and the SM prediction. Also the 68% and 95% uncertainty band values are shown. The observed limit of $\mu_{\text{obs}}^{\text{up}} = 7.4$ for the studied point of $C_t = -1$ and $C_V = +1$ agrees well with the expectation. For the SM scenario a limit of $\mu_{\text{obs}}^{\text{up}} = 106.9$ is observed, also a value well within the one standard deviation uncertainty band of the expected limit. A graphical representation of the upper limits on the tH production can be found in Figure 6.27.

	Region	Observed Limit	Expected Limit		
			Median	$\pm 1\sigma$	$\pm 2\sigma$
$C_t = -1$	3M	8.7	7.0	[4.7 , 10.7]	[3.4 , 16.0]
	4M	12.0	9.7	[6.3 , 15.6]	[4.4 , 25.2]
	Combination	7.4	5.7	[3.9 , 8.8]	[2.8 , 13.4]
$C_t = +1$ (SM)	3M	121.6	111.3	[74.3 , 173.8]	[53.2 , 265.0]
	4M	198.9	169.6	[108.5 , 279.2]	[75.9 , 454.8]
	Combination	106.9	97.3	[65.0 , 152.7]	[46.5 , 234.3]

6.13. Search for $C\mathcal{P}$ -mixing in tHq

The tH production process is not only sensitive to the sign of the top-Yukawa coupling, but also to a possible $C\mathcal{P}$ -mixing in the Higgs boson sector, as explained in Chapter 1.3.1. Under this assumption, a search for a generic spin-0, $C\mathcal{P}$ -symmetry violating particle X_0 with SM-like coupling to the W boson is conducted. The LHE reweighting procedure simplifies the concurrent analysis of different signal hypotheses immensely. Whereas the previously described analysis studied 51 different points in the C_V - C_t plane, a similar study has been performed investigating a possible $C\mathcal{P}$ -mixture in the tX_0q production.

A privately produced tX_0q simulation sample contains 21 sets of event weights for different $C\mathcal{P}$ -mixing angles corresponding to points ranging from values of $\cos \alpha = 1$, the SM prediction, to $\cos \alpha = -1$, corresponding to the previously studied point of $C_t = -1$, in steps of 0.1. The inclusion of a tX_0W signal sample will improve the analysis sensitivity further, but is left for future improvements. The sample has been produced at LO with `MADGRAPH5_AMCATNLO` and a leptonic decay of the top quark and a decay of the studied boson into bottom quarks is enforced. The cross section is scaled to its NLO prediction. Different $C\mathcal{P}$ -mixing angles change normalization and kinematics of the tX_0q process, as has been described in Section 1.3.1. The effect of different $C\mathcal{P}$ -mixing angles on the shape of the p_T of the top quark and the difference of the pseudorapidity of top quark and X_0 can be found in Figure 1.5.

The complete analysis is repeated as described in the first part of this chapter, but for the training of the tX_0q reconstruction and subsequently the training of the classification. These trainings are repeated for each of the 21 new signal inputs. The set of variables in the training is kept constant, as they already were optimized with this study kept in mind. Especially, variables like the absolute difference of the pseudorapidity of top quark and X_0 are considered, as they are expected to change drastically for different assumed $C\mathcal{P}$ -mixing angles, thereby increasing their discrimination potential between signal inputs and background (see also Figure 1.5). The reconstruction BDTs are trained with the same parameters as listed in Table 6.3.

The ranking of the variables used in the training of the $C\mathcal{P}$ -specific classification is shown in Table 6.10. The BDT response of the tX_0q reconstruction for a $C\mathcal{P}$ -mixing angle $\alpha = 90^\circ$ for the training sample and the independent testing sample can be seen in Figure 6.28. The corresponding area under the ROC curve for each of the studied 21 $C\mathcal{P}$ -mixing angles can be found in Figure 6.29.

The BDTs for the classification are trained with only 100 trees as the privately produced tX_0q samples contains fewer events than the officially produced sample and hence shows symptoms of overtraining if 400 trees are used in the training. Other parameters are chosen according to Table 6.5. The importance of the variables has been averaged over the 21 trainings for reconstruction and classification each and the result can be found in Table 6.10. Asymptotic CL_s limits at 95% C.L. are calculated after performing a simultaneous fit in the 3M region and the 4M region for each of the studied points. A peculiarity is the treatment of the $t\bar{t}X_0$ background, as the cross section and the kinematics change for different $C\mathcal{P}$ -mixing angles. The cross section and the kinematics of the $t\bar{t}X_0$ process are degenerate for α and $\pi - \alpha$ as can be seen in the first chapter in Figure 1.5. The $t\bar{t}X_0$ background is scaled to its NLO cross section for each of the studied angles but shape variations are neglected, as the variations only become more

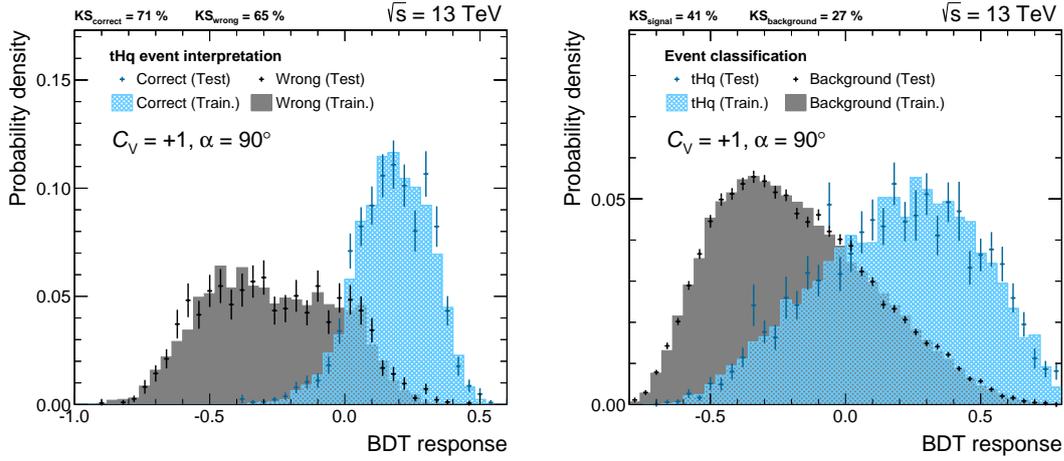


Figure 6.28.: The responses of the tX_0q reconstruction (left) and the classification (right) trainings for the training set (histograms) and the independent testing set (markers) under the assumption of a fully pseudoscalar CP boson. Training and testing sample agree well and no sign of overtraining is observed.

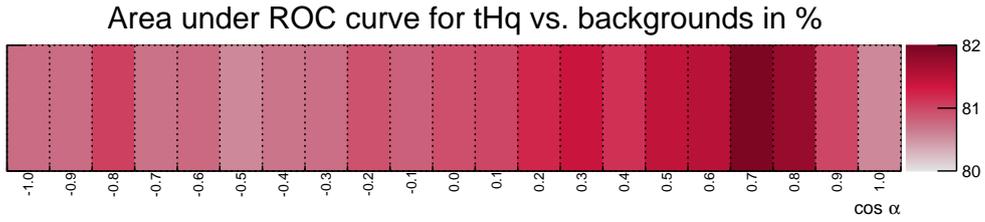


Figure 6.29.: The area under the ROC curve for all 21 tX_0q reconstruction trainings with different CP -mixing angles. With values above 80 good performance is observed for all trainings.

pronounced for the studied angles, where the impact of $t\bar{t}X_0$ in comparison to tX_0q becomes negligible, as they correspond to the lowest $t\bar{t}X_0$ cross sections. The cross sections for the tX_0q and $t\bar{t}X_0$ process for each of the studied mixing angles can be found in Appendix B.6.

The final classifier for the case of $\alpha = 90^\circ$ shows a very similar behavior as the classifier of the main analysis for simulation and data, as is expected. The output of the classifier for the 3M and 4M region before and after the simultaneous fit can be found in Figure 6.30. An illustration of the expected and observed limits for all investigated CP -mixing angles can be found in Figure 6.31. Again a cubic spline fit was used to interpolate between the estimated values. For the case of a fully pseudoscalar X_0 boson an upper limit of $\mu_{\text{obs}, \alpha=90^\circ}^{\text{up}} = 25.7$ has been observed, whereas a limit of $\mu_{\text{exp}, \alpha=90^\circ}^{\text{up}} = 22.7$ is expected. The complete set of limit values for a CP -mixing angle of $\alpha = 90^\circ$ can be found in Table 6.11.

The determined upper limits are anti-proportional to the cross section of the corresponding value of α . The lowest observed upper limit is therefore found for a value of $\alpha = 180^\circ$, which can exclude an enhancement of the cross section of $\mu_{\text{obs}, \alpha=180^\circ}^{\text{up}} = 11.7$. The limits for a CP -mixing angle of $\alpha = 0^\circ$ and $\alpha = 180^\circ$ are worse than the limits of their corresponding $C_t = +1$ and

6. Search for tH Production at $\sqrt{s} = 13$ TeV

Table 6.10.: Ranking of the input variables for the jet-assignment BDTs under the tX_0q hypothesis and classification averaged over the 21 studied CP -mixing angles.

tX_0q reco. Variable	Points	Class. Variable	Points
log m(X_0)	14.00	log m3	10.61
log m(t)	12.61	aplanarity	5.38
$\Delta R(X_0 \text{ jets})$	12.33	q(ℓ)	5.23
$\Delta R(b_t, W)$	10.33	Fox-Wolfram M. #1	2.90
$\Delta E(\text{light jet}, b_t)$	8.61	log m(t_{had})	13.85
relative H_T	7.23	CSV(W_{had} jet 1)	11.95
CSV(b_t)	7.19	CSV(W_{had} jet 2)	8.71
$\eta(\text{light jet}) - \eta(b_t)$	6.38	$\Delta R(W_{\text{had}} \text{ jets})$	3.80
cos $\theta(t, \ell)$	6.09	$\eta(\text{light jet})$	11.38
$\eta(t) - \eta(X_0)$	5.33	CSV(X_0 jet 2)	8.61
$\eta(\text{light jet})$	4.61	CSV(X_0 jet 1)	7.76
CSV(X_0 jet 2)	4.57	log p_T (light jet)	3.04
CSV(X_0 jet 1)	3.28	log p_T (X_0)	8.47
log min(p_T (X_0 jets))	1.38	$\eta(t) - \eta(X_0)$	3.04
$\eta(b_t)$	1.00	cos $\theta(t, \ell)$	0.19

$C_t = -1$ points of the main analysis, caused by the omission of the tHW signal sample and the training with fewer trees for reconstruction and classification. The analysis of a dataset corresponding to 2.3 fb^{-1} is not yet able to exclude any of the studied CP -mixing angles and no large deviations from the expectations have been found. For the pure CP -odd case the analysis is able to exclude an enhanced production cross for tX_0q section of $\mu_{\text{obs}}^{\text{up}} \cdot \sigma_{\alpha=90^\circ} = 25.7 \cdot 0.275 \text{ pb} = 7.1 \text{ pb}$.

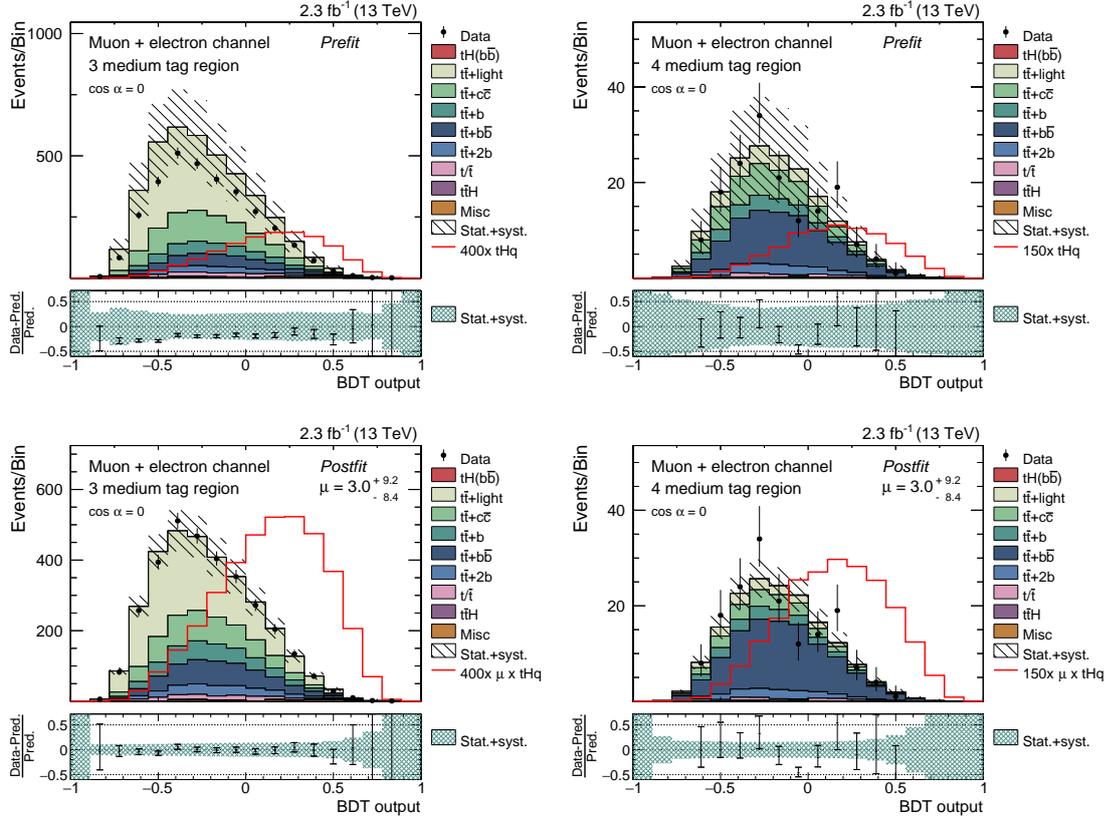


Figure 6.30.: Pre- and postfit distributions of the classifier output in the 3M and 4M region for a CP -mixing angle of 90° . These distributions are fitted simultaneously and a signal strength factor of $\mu = 3.0^{+9.2}_{-8.4}$ has been determined. A good agreement of simulation and data is observed after the fit.

Table 6.11.: Expected and observed asymptotic CL_S limits at 95% C.L. in the 3M and 4M region and their combination for the case of a fully pseudoscalar boson. Also the 68% and 95% uncertainty band values are shown. A limit of $\mu_{\text{obs}}^{\text{up}} = 25.7$ for the combination is observed, which agrees well with the expectation. A graphical representation of the upper limits on the tX_0q production can be found in Figure 6.31.

	Region	Observed Limit	Expected Limit		
			Median	$\pm 1\sigma$	$\pm 2\sigma$
$\alpha = 90^\circ$	3M	35.7	26.9	[18.1, 42.0]	[13.0, 64.7]
	4M	36.4	41.4	[26.5, 68.1]	[18.3, 111.8]
	Combination	25.7	22.7	[15.2, 35.6]	[10.9, 55.7]

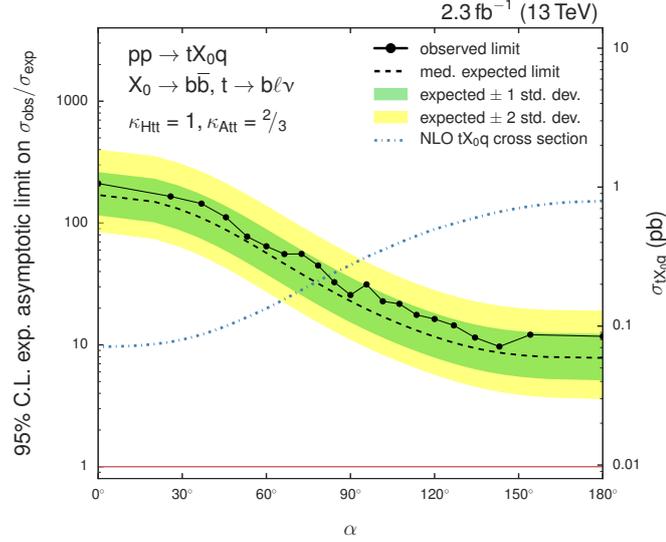


Figure 6.31.: Expected and observed asymptotic CL_s limits at 95% C.L. for the combination of the 3M and 4M region as function of the $C\mathcal{P}$ -mixing angle α . Corresponding numbers for the case of a fully pseudoscalar boson can be found in Table 6.11. The tX_0q cross section as function of α is depicted as dotted blue line. No $C\mathcal{P}$ -mixing angle can be excluded yet, hence only upper limits are set. Neither an excess nor a strong downward fluctuation is observed for any studied angle. The limit values for all mixing angles can be found in the Appendix in Table B.6.

6.14. Summary

In this chapter the search for tH with the Higgs boson decaying into $b\bar{b}$ at $\sqrt{s} = 13$ TeV has been presented. The analyzed data corresponds to 2.3 fb^{-1} recorded in 2015 by the CMS detector. The analysis is performed in parallel under 51 different top-Yukawa coupling hypotheses in order to be able to exclude points in the C_V - C_t parameter plane via this direct search.

Additionally, a search for $C\mathcal{P}$ -violation in the Higgs sector has been conducted exploiting the sensitivity of the associated production of a single top quark with a $C\mathcal{P}$ -symmetry violating boson X_0 . A search for tX_0q for 21 different $C\mathcal{P}$ -mixing angles has been conducted.

With the currently analyzed data neither a point in the coupling plane nor of the studied $C\mathcal{P}$ -mixing angles can be excluded. The updated analysis, which has been greatly improved with respect to the analysis at $\sqrt{s} = 8$ TeV, is already able to reach an expected limit only slightly worse than what was achieved at 8 TeV despite the substantially smaller available data sample. The observed limit is even better than what was achieved at 8 TeV, as no upward fluctuation of the limit as in its predecessor analysis is observed. This is made possible by the inclusion of tHW as additional signal process, the migration to an improved b -tagging algorithm and the optimization of the multivariate reconstructions and classification.

For the coupling point of most interest, the case of $C_t = -1$ and $C_V = +1$, an upper limit of $\mu_{\text{obs}}^{\text{up}} = 7.4$ at 95% C.L. can be set, whereas an upper limit of $\mu_{\text{exp}}^{\text{up}} = 5.7$ is expected. This value lies well within the uncertainty corresponding to one standard deviation.

In this analysis the first search for a $C\mathcal{P}$ -mixture state of the Higgs boson has been conducted

for 21 different mixing angles and as none of the studied angles can be excluded upper limits at 95% C.L. are set. For the case of a purely pseudoscalar boson ($\alpha = 90^\circ$) an upper limit of $\mu_{\text{obs}}^{\text{up}} = 25.7$ is observed, whereas a limit of $\mu_{\text{exp}}^{\text{up}} = 22.7$ is expected. Also here no strong fluctuation of the observed limit can be found, as the observed limit agrees well with the expectation within one standard deviation.

7. Conclusion and Outlook

The discovery of the Higgs boson in 2012 has led to a shift of the analyzer's focus. Only precise measurements of the Higgs boson properties can reveal small deviations from the SM expectations about said boson. One place, where such deviations could surface is the value of the top-Yukawa coupling y_t , the coupling strength of the Higgs boson to the top quark. Whereas measurements seem to indicate an absolute value of y_t close to one, a degeneracy regarding the sign of the coupling is apparent. The search for associated single top quark production with Higgs bosons could help lifting this degeneracy.

In this thesis a thorough search for exactly this process, where a Higgs boson is produced in association with a single top quark in the $H \rightarrow b\bar{b}$ channel at $\sqrt{s} = 8$ TeV and $\sqrt{s} = 13$ TeV is conducted. The analyzed data corresponds to 19.7 fb^{-1} and 2.3 fb^{-1} at 8 TeV and 13 TeV, respectively.

The analysis at $\sqrt{s} = 8$ TeV has been made public by the CMS collaboration in a Physics Analysis Summary [137] and is part of a combination of multiple tH analyses within CMS [138]. Multivariate analysis tools, which are employed during the object reconstruction of tHq events and semi-leptonic $t\bar{t}$ events, and during the classification of signal and background events, made it possible to separate signal events from an overwhelming background mainly consisting of semi-leptonic $t\bar{t}$ events. An upper limit of $\mu_{\text{obs}}^{\text{up}} = 7.5$ at 95% C.L. on a tHq process under the assumption of a flipped sign of the top-Yukawa coupling has been set, whereas a limit of $\mu_{\text{exp}}^{\text{up}} = 5.0$ was expected. In the combination with other tH analyses which exploit different Higgs boson decay channels an observed upper limit of $\mu_{\text{obs}}^{\text{up}} = 2.8$ with an expected upper limit of $\mu_{\text{exp}}^{\text{up}} = 2.0$ has been obtained for the case of an anomalous coupling.

The analysis at $\sqrt{s} = 13$ TeV constitutes a clear improvement over its predecessor analysis. It is already possible to reach a sensitivity comparable to that of the analysis in Run I with only a fraction of the previously recorded amount of data. Optimizations of the employed MVA methods and selection criteria, as well as the inclusion of an additional signal process made it possible to set an upper limit of $\mu_{\text{obs}}^{\text{up}} = 7.4$ at 95% C.L. for the case of a flipped sign of the top-Yukawa coupling, whereas an upper limit of $\mu_{\text{exp}}^{\text{up}} = 5.7$ was expected. Developments in the production of Monte Carlo simulation samples made it possible to extend the analysis to several different signal models. As such, the analysis was able to probe a whole plane of different possible Higgs boson coupling configurations and study the possibility of a $C\mathcal{P}$ -mixed state of the studied boson. This analysis constitutes the overall first search for $C\mathcal{P}$ -mixing in the Higgs boson sector in the tH production channel and was able to set the first upper limits. For a pseudoscalar boson X_0 an upper limit of $\mu_{\text{obs}}^{\text{up}} = 25.7$ has been set, whereas $\mu_{\text{exp}}^{\text{up}} = 22.7$ was expected.

The amount of data collected in 2015 is not yet sufficient to exclude any of the analyzed points, but this analysis will serve as a great starting basis for the ongoing Run II of the LHC and will

make it possible to exclude the first coupling scenarios in the months and years to come.

According to the optimistic long-term schedule of the LHC [193] a total of up to 150 fb^{-1} of data can be collected until the LS II in 2018 and with an upgrade to the High-Luminosity-LHC (HL-LHC) up to 3 ab^{-1} of data could potentially be recorded. A projection of the expected limit for these integrated luminosities can be found in Figure 7.1. The projection is performed under the assumption of a constant beam energy of $\sqrt{s} = 13 \text{ TeV}$, although an increase of the center-of-mass energy to $\sqrt{s} = 14 \text{ TeV}$ is imminent in the coming years, which would increase the analysis sensitivity even further. The nominal expected limit visualizes, how much the analysis gains from simply recording more data, without any improvements on the analysis side. The illustration also shows the effect of a 50% reduction of the overall systematic uncertainties, as e.g. foreseeable developments on the generator side would shrink the Q^2 scale uncertainty and more collected data will help to improve the CSV reweighting, therefore decreasing its associated uncertainties. With an interplay of reduced uncertainties and improved methods this analysis alone will be able to lift the degeneracy regarding the sign of the top-Yukawa coupling during Run II of the LHC. Other points in the C_V - C_t parameter plane will already be excluded much earlier in Run II. The projection is also performed for the coupling predicted by the SM. It is obvious that due to its very low cross section a direct search for tH as predicted by the SM is not viable, as huge improvements would be necessary to put this process into reach, even in the HL-LHC era.

The projection is only performed for the $tH \rightarrow b\bar{b}\ell\nu$ channel. An extension with a dilepton channel could be worthwhile, optimized for a tHW process, where both W bosons decay leptonically.

The more data is available the more can be gained from determining backgrounds directly from data, thus being independent from MC simulation samples. This could lead to a significant reduction of uncertainties and subsequently to an increased sensitivity of the analysis. A data-driven background estimation of the $t\bar{t}$ background has been studied during the course of the analysis at $\sqrt{s} = 8 \text{ TeV}$ [150], but was dismissed, as it would have introduced larger systematic uncertainties compared to the approach relying on MC simulation samples with the available integrated luminosity. However, the uncertainties associated with a data-driven background estimation will keep shrinking with more recorded data and will ultimately lead to a situation, where no MC simulation sample can keep up and a migration to a data-driven background estimation will be unavoidable to decrease uncertainties further.

For the realisation of the full potential of the tH production process a combination with other tH analyses within CMS is necessary and will increase its capability to exclude points in the C_V - C_t plane. As the analysis of this thesis is at the moment of writing the by far most progressed tH analysis in the still young Run II of the LHC, it is unclear how or when a combination will happen.

A further goal is the full inclusion of the tH analysis into the combined coupling fits of the CMS collaboration. At the moment tH is incorporated into the fit as possible signal contribution to the $t\bar{t}H$ analyses, but no dedicated tH analysis is included. This is at the moment hindered by an event overlap with the $t\bar{t}H$ analysis at CMS. A method will have to be developed to remove

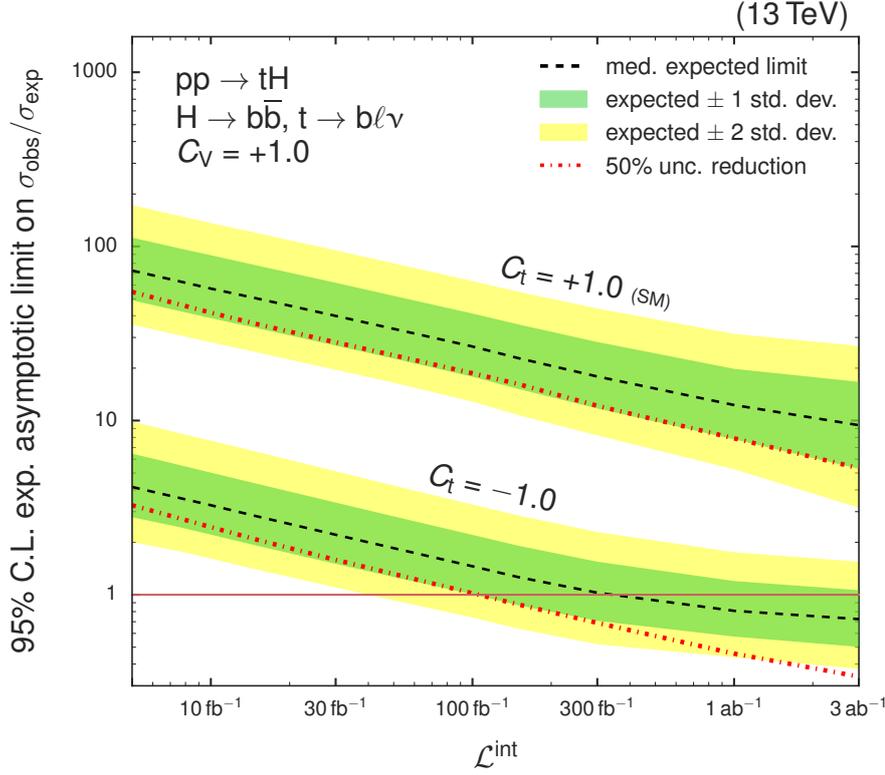


Figure 7.1.: Projection of the expected CL_s limits at 95% C.L. up to an integrated luminosity of 3 ab^{-1} for the nominal coupling pair of $C_t = -1$ and $C_V = +1$ and the SM prediction. The effect of a 50% reduction of the systematic uncertainties is depicted in red. Without improvements of the analysis nor reduction of the uncertainties the $C_t = -1$ coupling could be excluded at $\sim 300 \text{ fb}^{-1}$. A reduction of the systematic uncertainties could decrease the necessary amount of data to $\sim 100 \text{ fb}^{-1}$. The standard model tH production is out of reach even for $\mathcal{L}_{\text{int}} = 3 \text{ ab}^{-1}$.

this overlap without decreasing the sensitivity reach of either analysis too drastically.

At the time of writing of this thesis, no comparable tH ($H \rightarrow b\bar{b}$) analysis is published by the ATLAS collaboration, making this one of the few analyses which are exclusively performed and made public by the CMS collaboration.

This analysis of the associated production of a single top quark with a Higgs boson at $\sqrt{s} = 13 \text{ TeV}$ shows that the tH process is still a viable tool for the exclusion of possible top-Yukawa coupling parameters, thereby measuring the properties of the Higgs boson with increasing precision. By the end of Run II analyses studying the tH process will have contributed to our knowledge about the Higgs boson and aid in discovering nature's blueprint.

A. Appendix - Search for tHq Production at $\sqrt{s} = 8 \text{ TeV}$

Table A.1.: Experimental data exploited in the analysis. The second column shows the recorded pixel-based integrated luminosity calculated with the “golden” JSON file.

Dataset name	Int. luminosity
/SingleMu/Run2012A-22Jan2013-v1/AOD	876 pb ⁻¹
/SingleMu/Run2012B-22Jan2013-v1/AOD	4 412 pb ⁻¹
/SingleMu/Run2012C-22Jan2013-v1/AOD	7 055 pb ⁻¹
/SingleMu/Run2012D-22Jan2013-v1/AOD	7 369 pb ⁻¹
/SingleElectron/Run2012A-22Jan2013-v1/AOD	876 pb ⁻¹
/SingleElectron/Run2012B-22Jan2013-v1/AOD	4 412 pb ⁻¹
/SingleElectron/Run2012C-22Jan2013-v1/AOD	7 055 pb ⁻¹
/SingleElectron/Run2012D-22Jan2013-v1/AOD	7 369 pb ⁻¹

A. Appendix - Search for tHq Production at $\sqrt{s} = 8 \text{ TeV}$

Table A.2.: Utilized nominal and systematically varied simulation samples. All samples have been processed starting from the “/AODSIM” data format, therefore the common label is omitted. The common notation for the production era “Summer12_DR53X-PU_S10_START53_V7A” is also omitted everywhere. If only specific top quark decays are needed, the cross section is multiplied with the top quark branching ratio of $\mathcal{BR}(t) = 0.1080 \pm 0.0009$. The cross sections are obtained from the provided references, or, if no reference is provided, from Reference [197] or the generator itself.

Dataset name	Cross section, pb
/THTo3BLNu_t-channel-AnomPhase_8TeV-madgraph/. . . -v1	$36.8 \cdot 10^{-3}$ (NLO) [38]
/TTH_Inclusive_M-125_8TeV_pythia6/. . . -v1	$130.2 \cdot 10^{-3}$ (NLO) [194]
/TTJets_SemiLeptMGDecays_8TeV-madgraph/. . . -v1	107.7 (NNLO) [195]
/TTJets_SemiLeptMGDecays_8TeV-madgraph/. . . _ext1-v1	107.7 (NNLO) [195]
/TTJets_SemiLeptMGDecays_8TeV-madgraph/. . . _ext2-v1	107.7 (NNLO) [195]
/TTJets_FullLeptMGDecays_8TeV-madgraph/. . . -v2	25.8 (NNLO) [195]
/TToLeptons_t-channel_8TeV-powheg-tauola/. . . -v1	18.27 (approx. NNLO) [196]
/TBarToLeptons_t-channel_8TeV-powheg-tauola/. . . -v1	9.95 (approx. NNLO) [196]
/T_tW-channel-DR_TuneZ2star_8TeV-powheg-tauola/. . . -v1	11.1 (approx. NNLO) [196]
/Tbar_tW-channel-DR_TuneZ2star_8TeV-powheg-tauola/. . . -v1	11.1 (approx. NNLO) [196]
/TToLeptons_s-channel_8TeV-powheg-tauola/. . . -v1	1.23 (approx. NNLO) [196]
/TToLeptons_s-channel_8TeV-powheg-tauola/. . . -v1	0.57 (approx. NNLO) [196]
/WJetsToLNu_TuneZ2Star_8TeV-madgraph-tarball/. . . -v2	35 509 (NNLO)
/W2JetsToLNu_TuneZ2Star_8TeV-madgraph/. . . -v1	2116 (NNLO)
/W3JetsToLNu_TuneZ2Star_8TeV-madgraph/. . . -v1	637 (NNLO)
/W4JetsToLNu_TuneZ2Star_8TeV-madgraph/. . . -v1	262 (NNLO)
/WW_TuneZ2star_8TeV_pythia6_tauola/. . . -v1	54.8 (NLO)
/WZ_TuneZ2star_8TeV_pythia6_tauola/. . . -v1	12.6 (LO)
/ZZ_TuneZ2star_8TeV_pythia6_tauola/. . . -v1	5.2 (LO)
/DYJetsToLL_M-50_TuneZ2Star_8TeV-madgraph-tarball/. . . -v1	3504 (NNLO)
/THTo3BLNu_t-channel-AnomPhase_scaleup_8TeV-madgraph/. . . -v1	$36.8 \cdot 10^{-3}$ (NLO)
/THTo3BLNu_t-channel-AnomPhase_scaledown_8TeV-madgraph/. . . -v1	$36.8 \cdot 10^{-3}$ (NLO)
/TTJets_scaleup_TuneZ2star_8TeV-madgraph-tauola/. . . -v1	245.8 (NNLO)
/TTJets_scaledown_TuneZ2star_8TeV-madgraph-tauola/. . . -v1	245.8 (NNLO)
/TTJets_matchingup_TuneZ2star_8TeV-madgraph-tauola/. . . -v1	245.8 (NNLO)
/TTJets_matchingdown_TuneZ2star_8TeV-madgraph-tauola/. . . -v1	245.8 (NNLO)

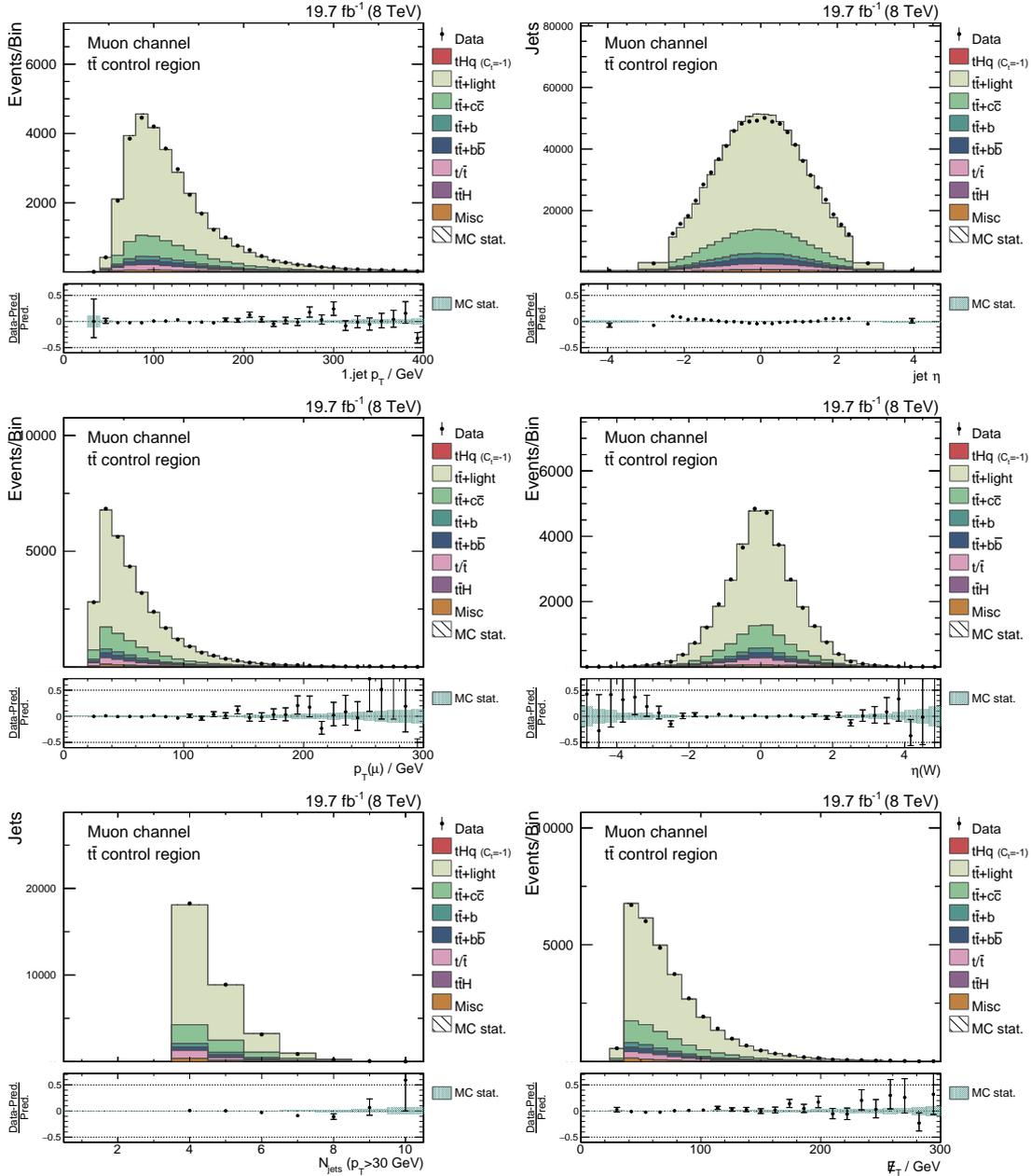


Figure A.1: A set of control plot in the $t\bar{t}$ control (2T) region in the muon channel is shown. The diagrams show the transverse momentum of the hardest jet of the event, the pseudorapidity of all jets, the transverse momentum of the muon, the pseudorapidity of the reconstructed W boson, the number of jets with a transverse momentum greater than 30 GeV and the missing transverse energy of the event. For all distributions the simulation is scaled to match the observed event yields in data. A good agreement between simulation and data is observed.

A. Appendix - Search for tHq Production at $\sqrt{s} = 8 \text{ TeV}$

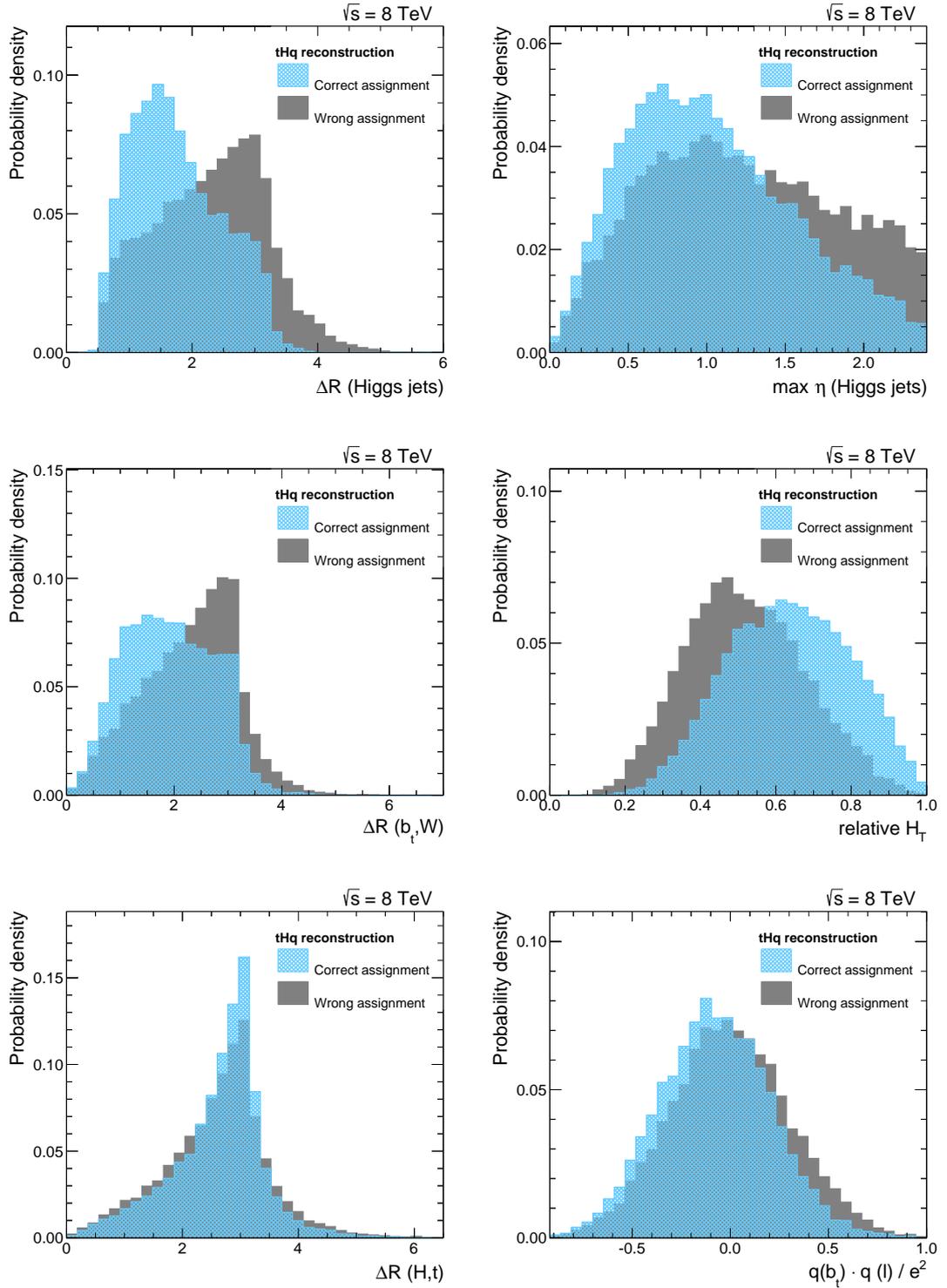


Figure A.2.: The remaining six variables between correct and wrong jet assignments in the tHq reconstruction are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 5.3.

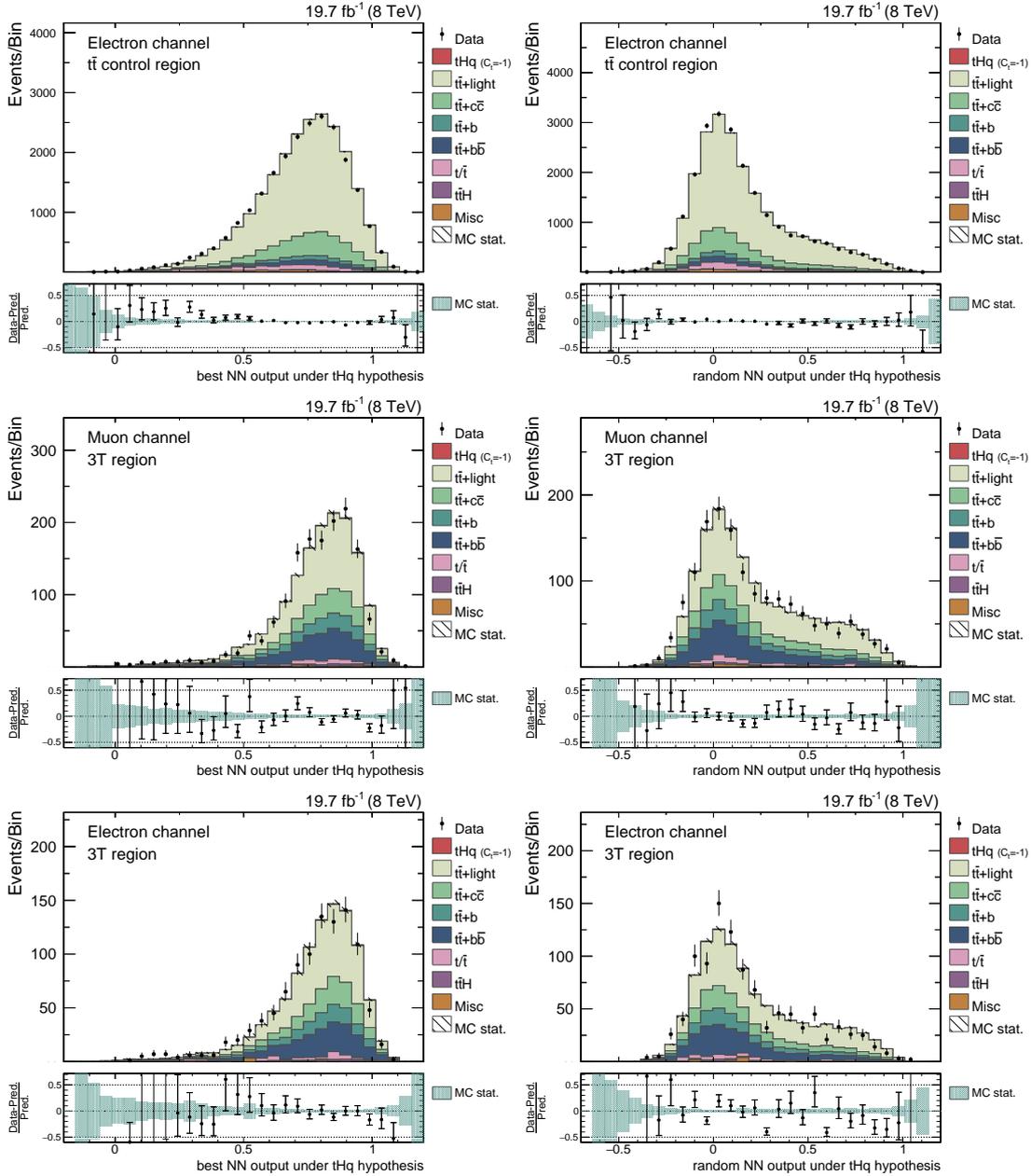


Figure A.3.: The distributions show the assignment for each event, which is assigned with the highest output value of the neural net (left column) and a random assignment (right column). The diagrams show the situation in the electron channel in the 2T region (top row), the muon channel in the 3T region (middle row) and the electron channel in the 3T region (bottom row). In all distributions a good agreement between simulated samples and recorded data is found. The simulation is scaled to match the event yields observed in data and all MC weights are applied.

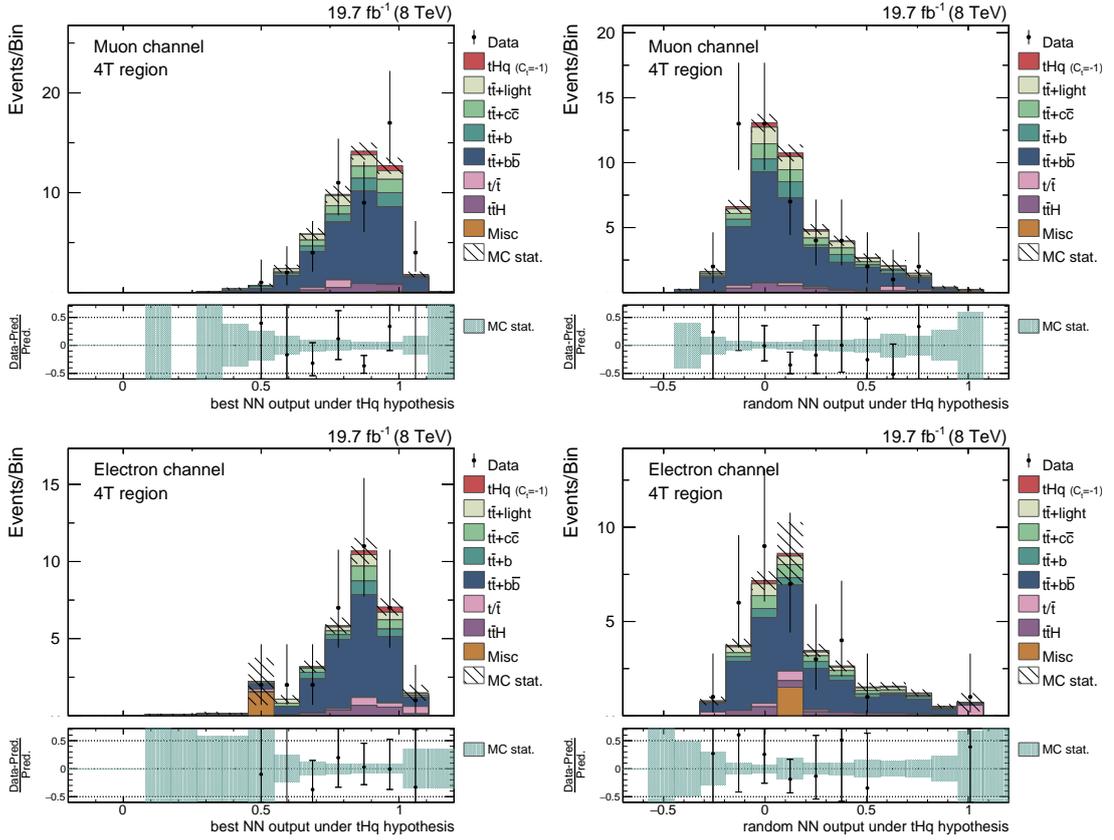


Figure A.4.: Response of the $t\bar{t}H$ reconstruction comparing simulation to data. In the left column the highest output value (chosen jet assignment) per event is shown and on the right the NN output value for a random assignment is shown. The diagrams show the situation in the 4T region in the muon channel (top row) and the electron channel (bottom row). In all distributions a good agreement between simulation and data is found. The simulation is scaled to match the event yields observed in data and all MC weights are applied.

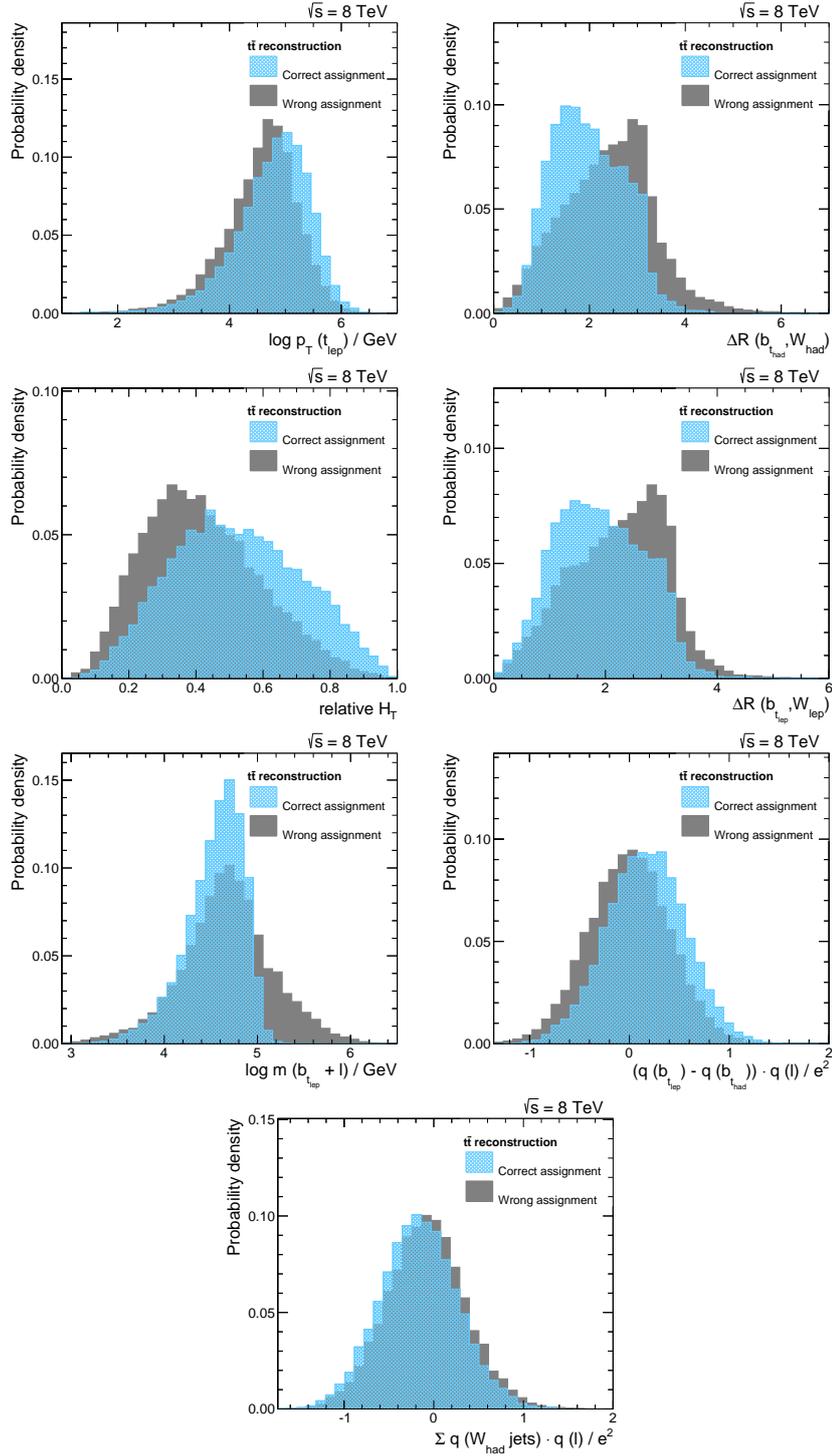


Figure A.5.: The remaining seven variables between correct and wrong jet assignments in the $t\bar{t}$ reconstruction are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 5.4.

A. Appendix - Search for $t\bar{t}H$ Production at $\sqrt{s} = 8 \text{ TeV}$

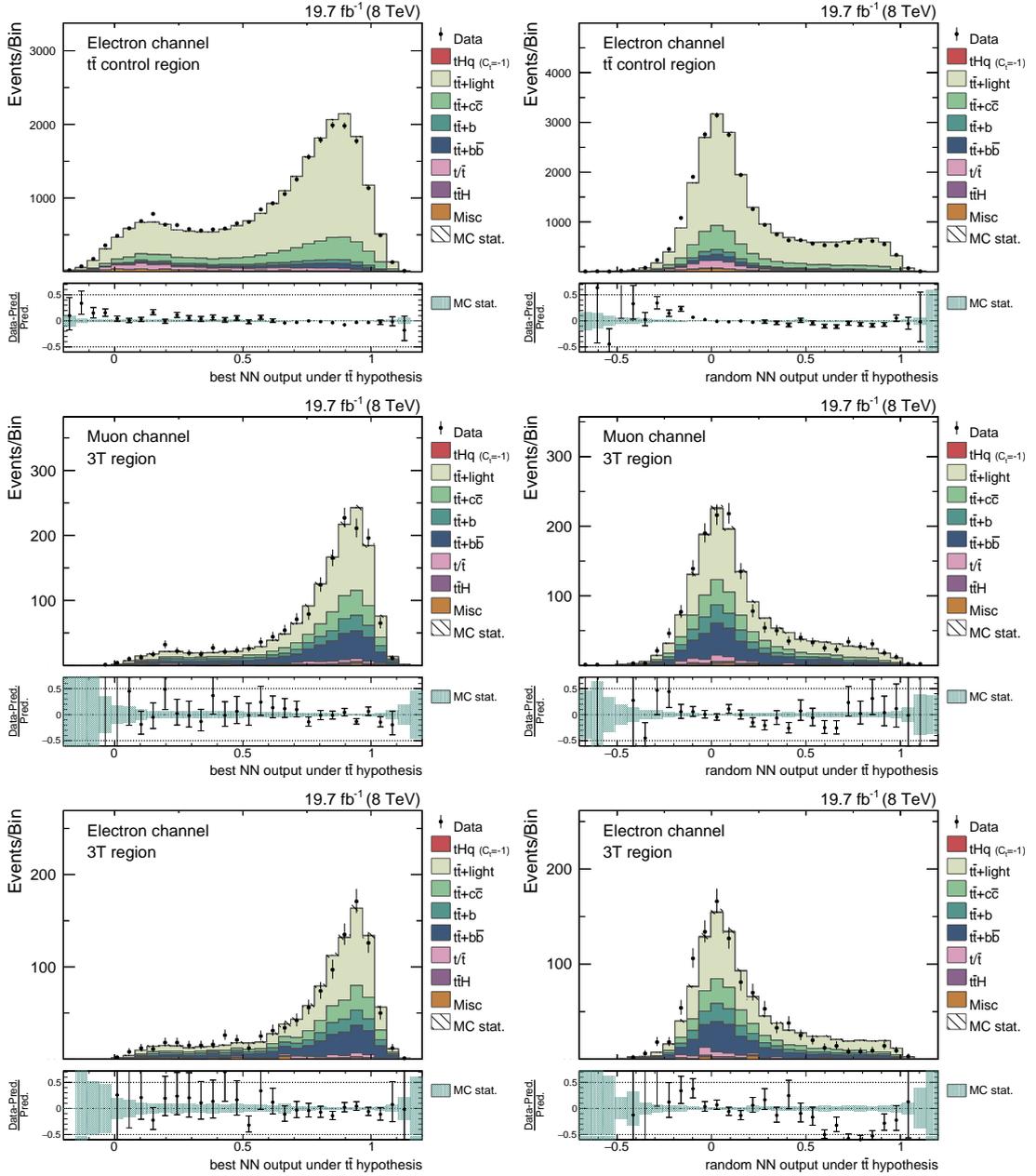


Figure A.6.: Response of the $t\bar{t}$ reconstruction comparing simulation to data. In the left column the highest output value (chosen jet assignment) per event is shown and on the right the NN output value for a random assignment is shown. The diagrams show the situation in the electron channel in the 2T region (top row), in the muon channel in the 3T region (middle row) and in the electron channel in the 3T region (bottom row). In all distributions a good agreement between simulation and data is found. The simulation is scaled to match the event yields observed in data and all MC weights are applied.

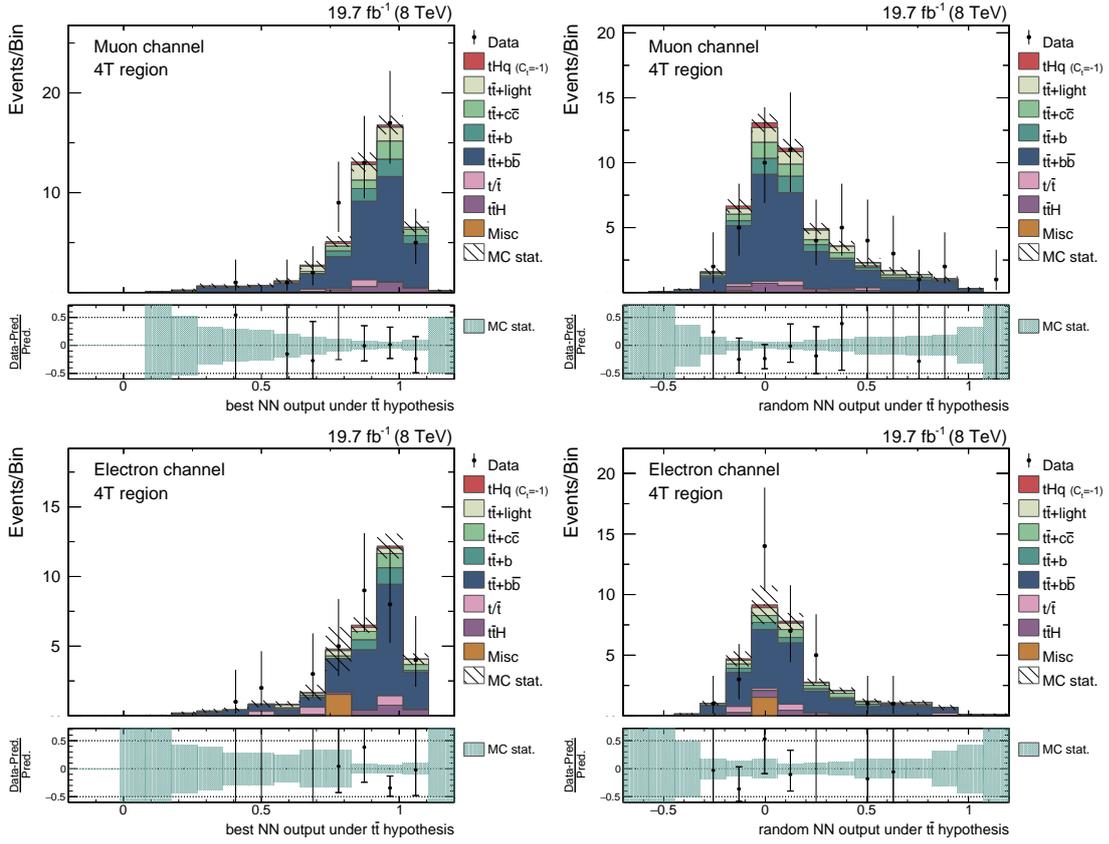


Figure A.7.: Response of the $t\bar{t}$ reconstruction comparing simulation to data. In the left column the highest output value (chosen jet assignment) per event is shown and on the right the NN output value for a random assignment is shown. The diagrams show the situation in the 4T region for the muon channel (top row) and the electron channel (bottom row). In all distributions a good agreement between simulation and data is found. The simulation is scaled to match the event yields observed in data and all MC weights are applied.

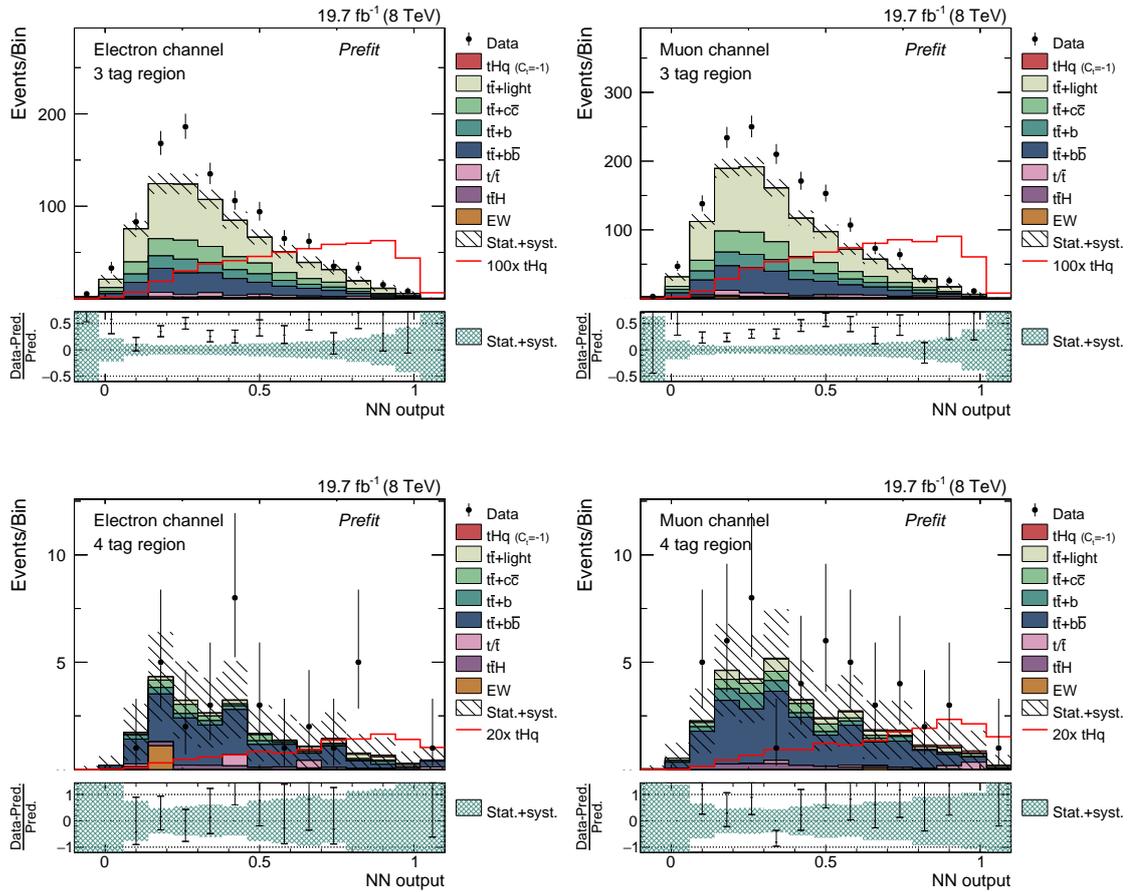


Figure A.8.: Prefit distributions of the classifier output in both lepton channels and both signal regions. Uncertainty bands include statistical and systematic uncertainties. Besides a clear offset a fair agreement in all four channels is observed, deviations from the prediction are covered by the uncertainties. The simulation samples are scaled to their expected event yields.

B. Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV

Table B.1.: The CMS datasets utilized in this thesis, which are selected according to the GoldenJSON.

Dataset name	Int. luminosity
/SingleMuon/Run2015D-16Dec2015-v1/MINIAOD	2300.5 pb ⁻¹
/SingleElectron/Run2015D-16Dec2015-v1/MINIAOD	2300.5 pb ⁻¹

Table B.2.: Utilized simulation samples are listed with their corresponding cross sections. The common labels “RunIIFall15MiniAODv2-PU25nsData2015v1_76X_mcRun2_asymptotic_v12” and “/MINIAODSIM” are omitted for all samples. When only specific decays of a top quark are considered in a dataset, the inclusive cross section is scaled using the corresponding $t\bar{t}$ branching ratio.

Dataset name	Cross section (pb)
/THQ_Hinc1_13TeV-madgraph-pythia8_TuneCUETP8M1/. . . -v1	see Table B.3
/THW_Hinc1_13TeV-madgraph-pythia8_TuneCUETP8M1/. . . -v1	see Table B.4
/ST_t-channel_antitop_4f_leptonDecays_13TeV-powheg-pythia8_TuneCUETP8M1/. . . -v1	27 (NLO [190])
/ST_t-channel_top_4f_leptonDecays_13TeV-powheg-pythia8_TuneCUETP8M1/. . . -v1	45.3 (NLO [190])
/ST_tW_antitop_5f_inclusiveDecays_13TeV-powheg-pythia8_TuneCUETP8M1/. . . -v1	35.9 (NNLO [190])
/ST_tW_top_5f_inclusiveDecays_13TeV-powheg-pythia8_TuneCUETP8M1/. . . -v1	35.9 (NNLO [190])
/TTWJetsToLNu_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8/. . . -v1	0.21 (NLO [191])
/TTWJetsToQQ_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8/. . . -v1	0.435 (NLO [191])
/TTToSemiLeptonic_13TeV-powheg/. . . _ext1-v1	831.76 (NNLO [189])
/TTTo2L2Nu_13TeV-powheg/. . . -v1	831.76 (NNLO [189])
/TTTo2L2Nu_13TeV-powheg/. . . _ext1-v1	831.76 (NNLO [189])
/WJetsToLNu_TuneCUETP8M1_13TeV-madgraphMLM-pythia8/. . . -v1	61,526.7 (NNLO [198])
/WW_TuneCUETP8M1_13TeV-pythia8/. . . -v1	118.7 (NNLO [192])
/WZ_TuneCUETP8M1_13TeV-pythia8/. . . -v1	47.13 (NLO [192])
/ZZ_TuneCUETP8M1_13TeV-pythia8/. . . -v1	16.523 (NLO [198])
/ttHTobb_M125_13TeV-powheg-pythia8/. . . -v1	0.2934 (NLO [32])

B. Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV

Table B.3.: Production cross sections for tHq at $\sqrt{s} = 13$ TeV, depending on C_t and C_V . Obtained with MADGRAPH5_AMC@NLO at NLO in the 4F scheme. The quoted uncertainties on the cross section correspond to scale variations in %. Values are taken from Reference [36].

C_t	C_V	σ (pb)	C_t	C_V	σ (pb)	C_t	C_V	σ (pb)
-3.0	0.5	$2.260^{+1.9}_{-2.7}$	-3.0	1.0	$2.991^{+2.1}_{-3.1}$	-3.0	1.5	$3.845^{+2.6}_{-3.2}$
-2.0	0.5	$1.160^{+2.0}_{-2.9}$	-2.0	1.0	$1.706^{+2.6}_{-3.2}$	-2.0	1.5	$2.371^{+2.5}_{-3.6}$
-1.5	0.5	$0.748^{+2.1}_{-3.1}$	-1.5	1.0	$1.205^{+2.5}_{-3.6}$	-1.5	1.5	$1.784^{+2.7}_{-3.9}$
-1.25	0.5	$0.573^{+2.1}_{-3.0}$	-1.25	1.0	$0.987^{+2.6}_{-3.4}$	-1.25	1.5	$1.518^{+2.8}_{-3.9}$
-1.0	0.5	$0.472^{+2.3}_{-3.3}$	-1.0	1.0	$0.793^{+2.7}_{-3.9}$	-1.0	1.5	$1.287^{+3.0}_{-4.3}$
-0.75	0.5	$0.300^{+2.5}_{-3.5}$	-0.75	1.0	$0.621^{+2.9}_{-4.1}$	-0.75	1.5	$1.067^{+3.1}_{-4.4}$
-0.5	0.5	$0.198^{+2.8}_{-3.9}$	-0.5	1.0	$0.472^{+3.2}_{-4.4}$	-0.5	1.5	$0.874^{+3.4}_{-4.7}$
-0.25	0.5	$0.119^{+3.1}_{-4.6}$	-0.25	1.0	$0.351^{+3.5}_{-5.0}$	-0.25	1.5	$0.703^{+3.6}_{-5.0}$
0.0	0.5	$0.062^{+3.8}_{-5.6}$	0.0	1.0	$0.248^{+3.9}_{-5.5}$	0.0	1.5	$0.558^{+3.8}_{-5.4}$
0.25	0.5	$0.028^{+5.0}_{-7.1}$	0.25	1.0	$0.169^{+4.4}_{-6.2}$	0.25	1.5	$0.437^{+4.2}_{-6.1}$
0.5	0.5	$0.018^{+4.2}_{-6.7}$	0.5	1.0	$0.113^{+5.0}_{-7.1}$	0.5	1.5	$0.334^{+4.6}_{-6.5}$
0.75	0.5	$0.030^{+1.4}_{-2.9}$	0.75	1.0	$0.081^{+5.7}_{-7.6}$	0.75	1.5	$0.256^{+5.2}_{-7.2}$
1.0	0.5	$0.066^{+1.0}_{-3.6}$	1.0	1.0	$0.071^{+4.1}_{-6.7}$	1.0	1.5	$0.200^{+5.7}_{-7.6}$
1.25	0.5	$0.124^{+0.9}_{-3.7}$	1.25	1.0	$0.084^{+2.3}_{-4.6}$	1.25	1.5	$0.167^{+5.5}_{-7.5}$
1.5	0.5	$0.205^{+0.8}_{-3.7}$	1.5	1.0	$0.120^{+1.2}_{-2.9}$	1.5	1.5	$0.159^{+4.1}_{-6.7}$
2.0	0.5	$0.436^{+1.0}_{-3.6}$	2.0	1.0	$0.260^{+1.0}_{-3.6}$	2.0	1.5	$0.211^{+2.0}_{-3.9}$
3.0	0.5	$1.177^{+1.2}_{-3.2}$	3.0	1.0	$0.821^{+0.8}_{-3.7}$	3.0	1.5	$0.589^{+0.9}_{-3.7}$

Table B.4.: Production cross sections for tHW at $\sqrt{s} = 13$ TeV, depending on C_t and C_V . Obtained with MADGRAPH5_AMC@NLO at NLO in the 5F scheme. The quoted uncertainties on the cross section correspond to scale variations in %. Values are taken from Reference [36].

C_t	C_V	σ (pb)	C_t	C_V	σ (pb)	C_t	C_V	σ (pb)
-3.0	0.5	$0.514^{+2.3}_{-3.0}$	-3.0	1.0	$0.641^{+2.3}_{-2.7}$	-3.0	1.5	$0.783^{+2.1}_{-2.1}$
-2.0	0.5	$0.255^{+2.3}_{-2.8}$	-2.0	1.0	$0.346^{+2.2}_{-2.5}$	-2.0	1.5	$0.457^{+2.1}_{-2.1}$
-1.5	0.5	$0.159^{+2.3}_{-2.8}$	-1.5	1.0	$0.253^{+2.1}_{-2.2}$	-1.5	1.5	$0.329^{+1.9}_{-1.8}$
-1.25	0.5	$0.120^{+2.2}_{-2.5}$	-1.25	1.0	$0.188^{+2.0}_{-2.0}$	-1.25	1.5	$0.275^{+1.9}_{-1.6}$
-1.0	0.5	$0.087^{+2.1}_{-2.3}$	-1.0	1.0	$0.147^{+2.0}_{-1.8}$	-1.0	1.5	$0.224^{+1.9}_{-1.5}$
-0.75	0.5	$0.059^{+2.0}_{-2.1}$	-0.75	1.0	$0.110^{+2.0}_{-1.7}$	-0.75	1.5	$0.180^{+1.8}_{-1.3}$
-0.5	0.5	$0.037^{+1.9}_{-1.8}$	-0.5	1.0	$0.080^{+1.7}_{-1.4}$	-0.5	1.5	$0.141^{+1.6}_{-1.2}$
-0.25	0.5	$0.020^{+1.8}_{-1.3}$	-0.25	1.0	$0.055^{+1.6}_{-1.1}$	-0.25	1.5	$0.108^{+1.6}_{-1.2}$
0.0	0.5	$0.009^{+1.6}_{-1.3}$	0.0	1.0	$0.036^{+1.5}_{-1.2}$	0.0	1.5	$0.081^{+1.5}_{-1.2}$
0.25	0.5	$0.004^{+2.1}_{-2.0}$	0.25	1.0	$0.022^{+1.6}_{-1.5}$	0.25	1.5	$0.059^{+1.5}_{-1.4}$
0.5	0.5	$0.004^{+4.6}_{-6.1}$	0.5	1.0	$0.014^{+2.1}_{-2.0}$	0.5	1.5	$0.043^{+1.8}_{-1.7}$
0.75	0.5	$0.010^{+4.7}_{-6.3}$	0.75	1.0	$0.012^{+3.2}_{-3.9}$	0.75	1.5	$0.033^{+2.1}_{-2.0}$
1.0	0.5	$0.021^{+4.0}_{-5.5}$	1.0	1.0	$0.016^{+4.6}_{-6.1}$	1.0	1.5	$0.028^{+2.8}_{-3.0}$
1.25	0.5	$0.038^{+3.7}_{-5.2}$	1.25	1.0	$0.025^{+4.8}_{-5.4}$	1.25	1.5	$0.029^{+3.6}_{-4.7}$
1.5	0.5	$0.061^{+3.5}_{-4.9}$	1.5	1.0	$0.039^{+4.6}_{-6.3}$	1.5	1.5	$0.035^{+4.6}_{-6.0}$
2.0	0.5	$0.125^{+3.0}_{-4.3}$	2.0	1.0	$0.086^{+4.0}_{-5.5}$	2.0	1.5	$0.065^{+4.8}_{-6.5}$
3.0	0.5	$0.317^{+2.8}_{-4.0}$	3.0	1.0	$0.247^{+3.3}_{-4.6}$	3.0	1.5	$0.193^{+4.0}_{-5.6}$

B. Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV

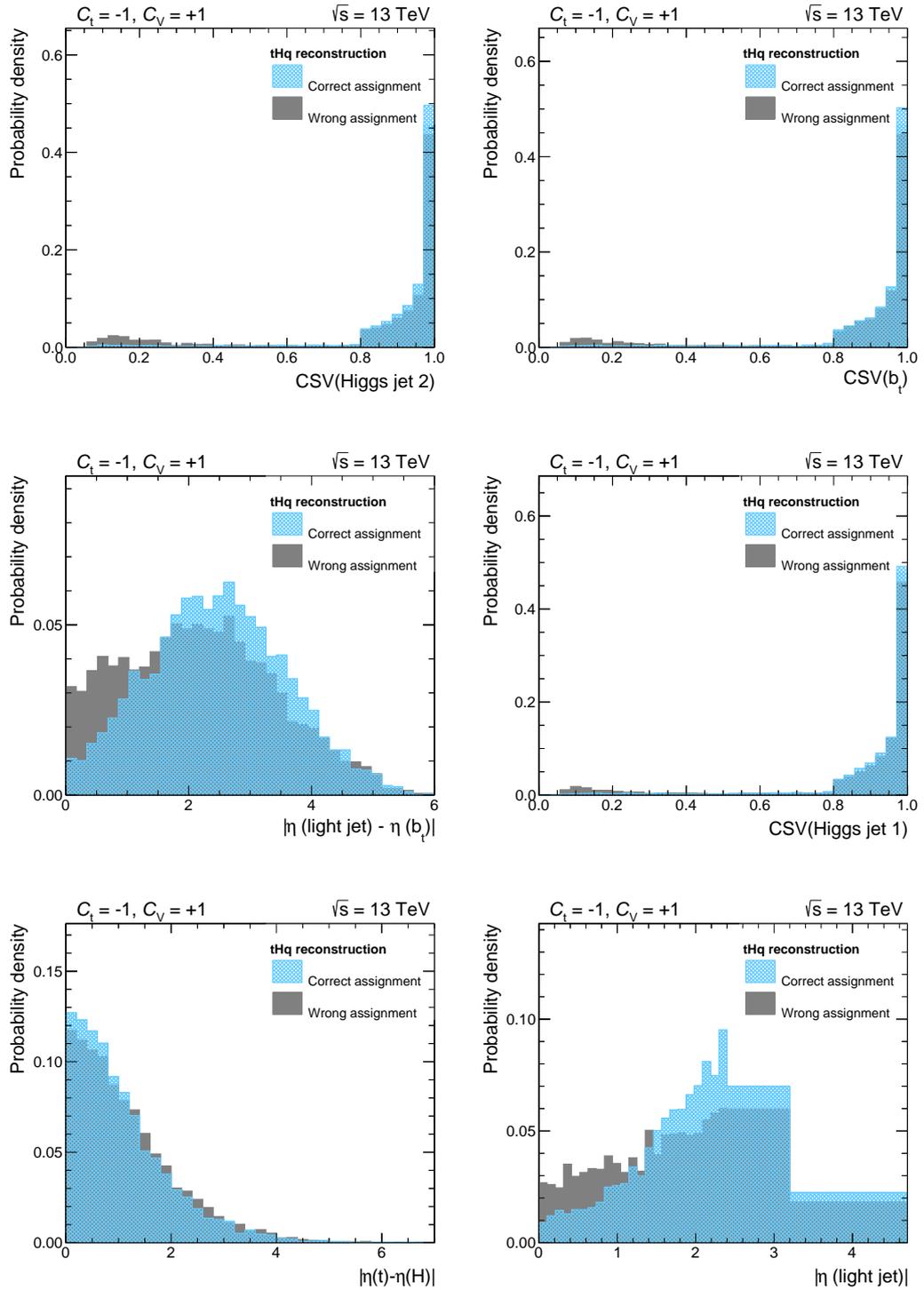


Figure B.1.: Distributions of variables ranked 7th to 12th place are shown for correct and wrong jet assignments in the tHq reconstruction. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.2. The three remaining variables can be found in Appendix B.2.

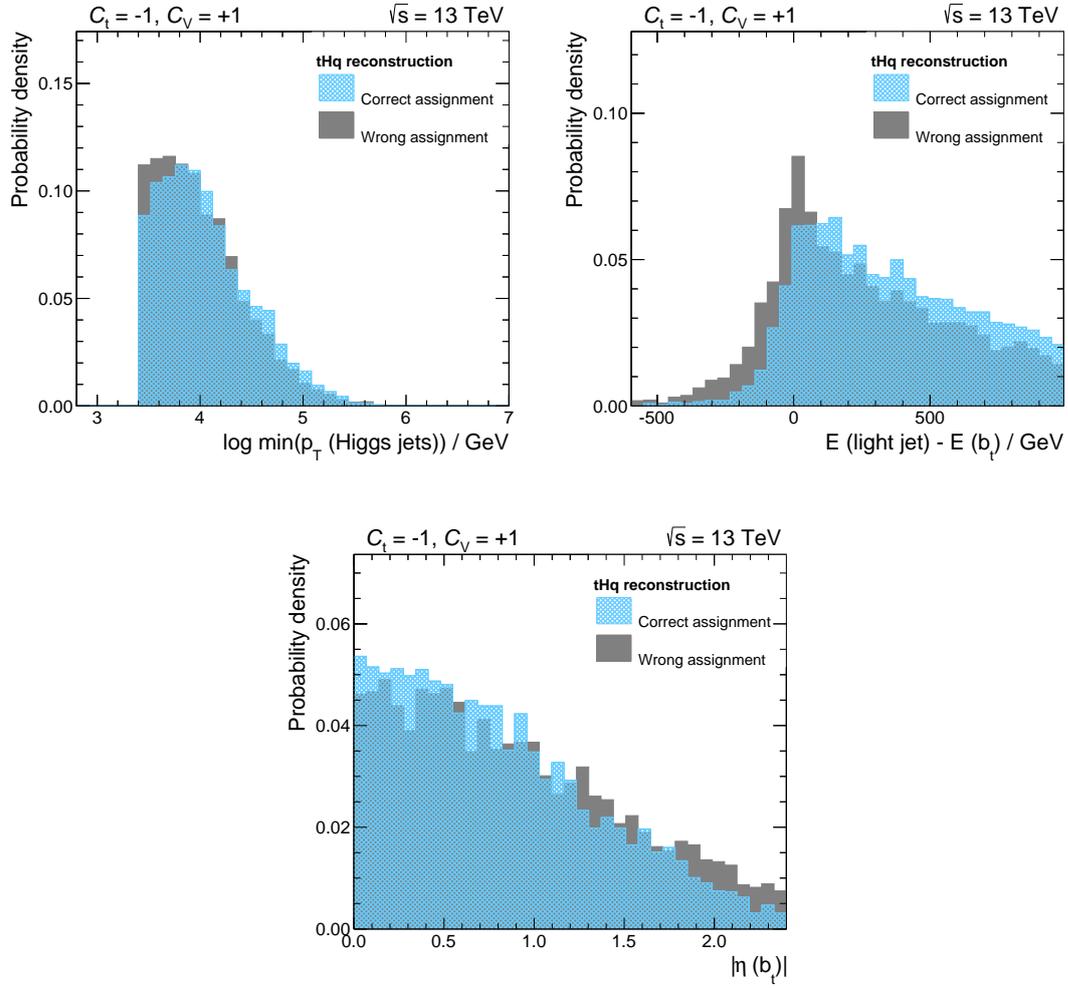


Figure B.2.: Distributions of variables ranked 13th to 15th place are shown for correct and wrong jet assignments in the tHq reconstruction. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.2.

B. Appendix - Search for $t\bar{t}H$ Production at $\sqrt{s} = 13$ TeV

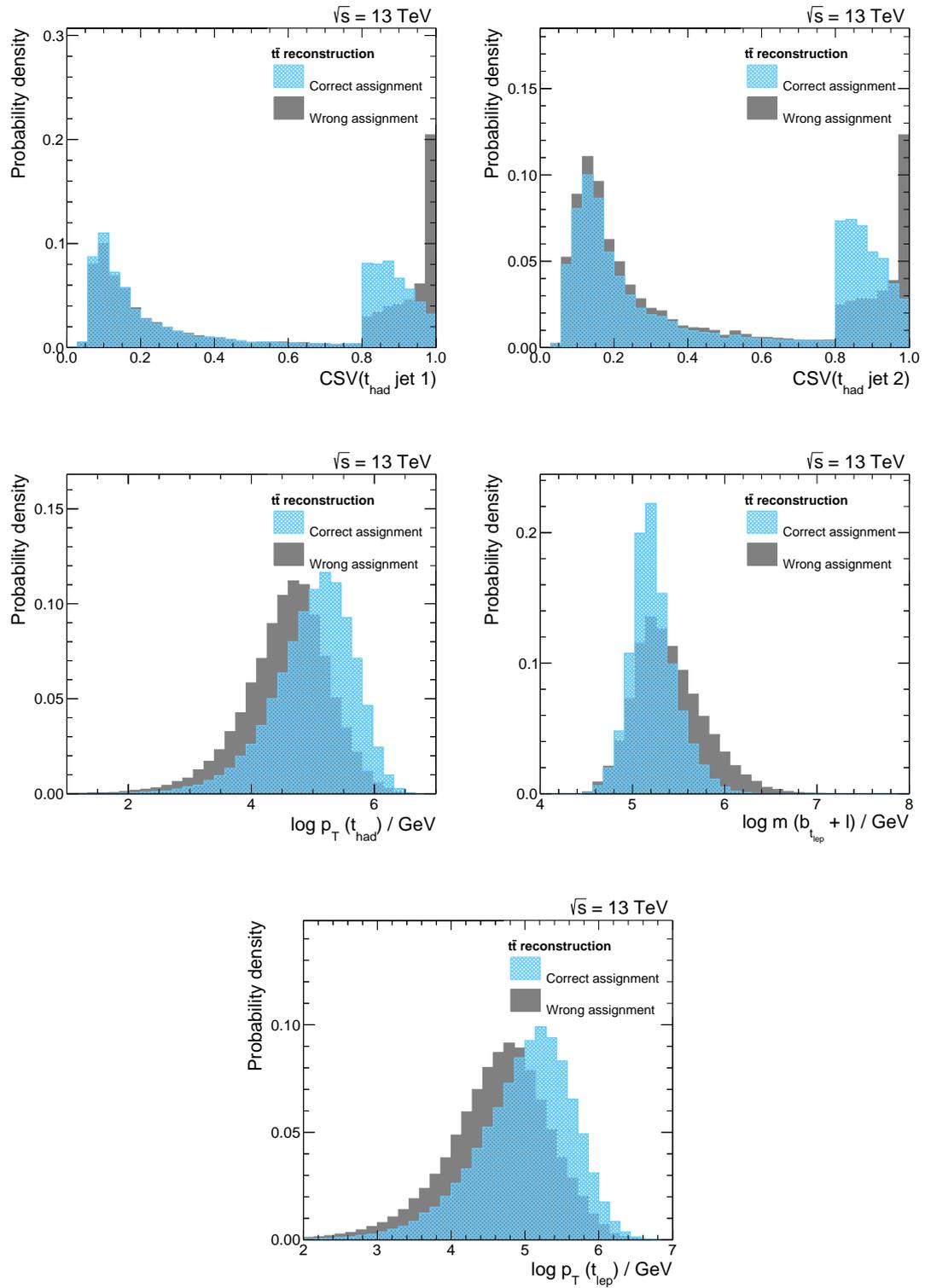


Figure B.3.: The remaining five variables between correct and wrong jet assignments in the $t\bar{t}$ reconstruction are shown sorted by their importance in the training. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.4.

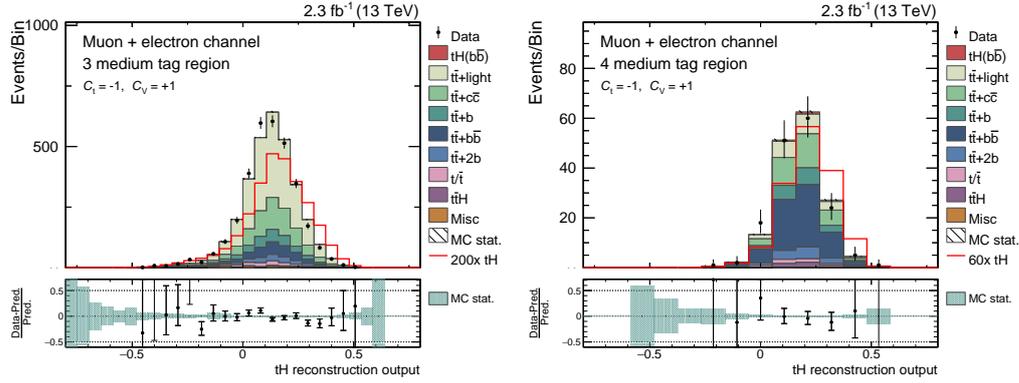


Figure B.4.: Response of the tHq reconstruction comparing simulation to data. The distribution of the highest output value (chosen jet assignment) per event is shown for the $C_t = -1$ and $C_V = +1$ point in the 3M region (left) and in the 4M region (right). In both distributions a good agreement between simulation and data is found. The simulation is scaled to match the event yields observed in data and all MC weights are applied.

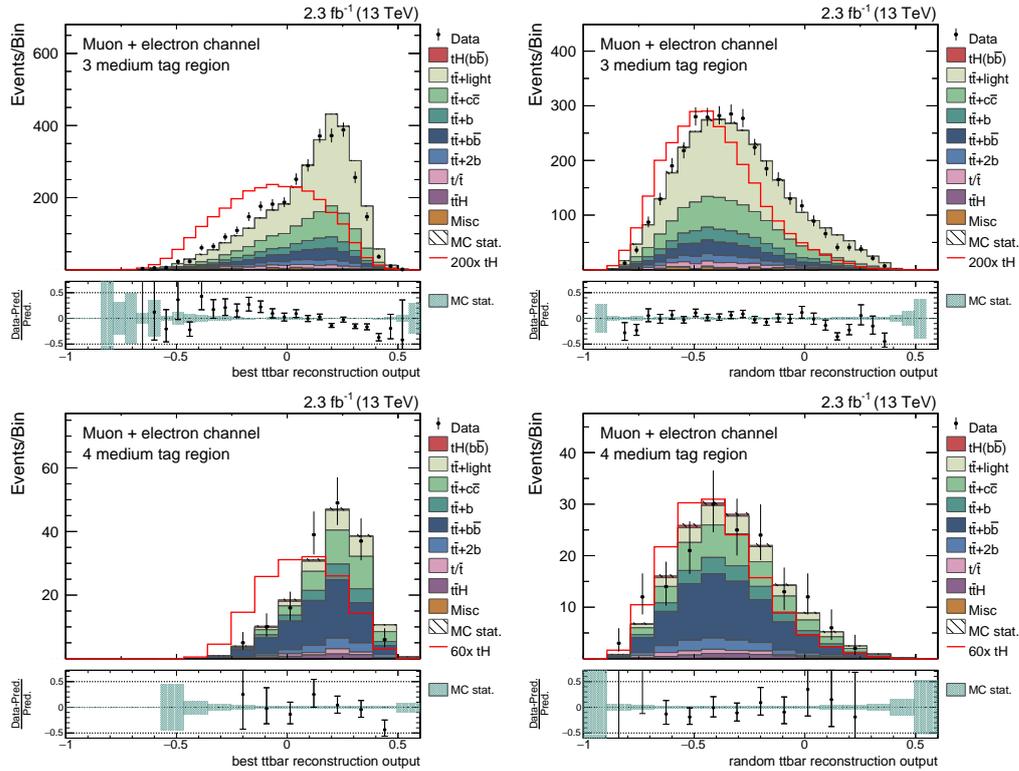


Figure B.5.: Response of the $t\bar{t}$ reconstruction comparing simulation to data in the 3M and 4M region. On the left the highest output value (chosen jet assignment) per event is shown and on the right the BDT output value for a random assignment is shown. In all distributions a good agreement between simulation and data is found. All diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

B. Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV

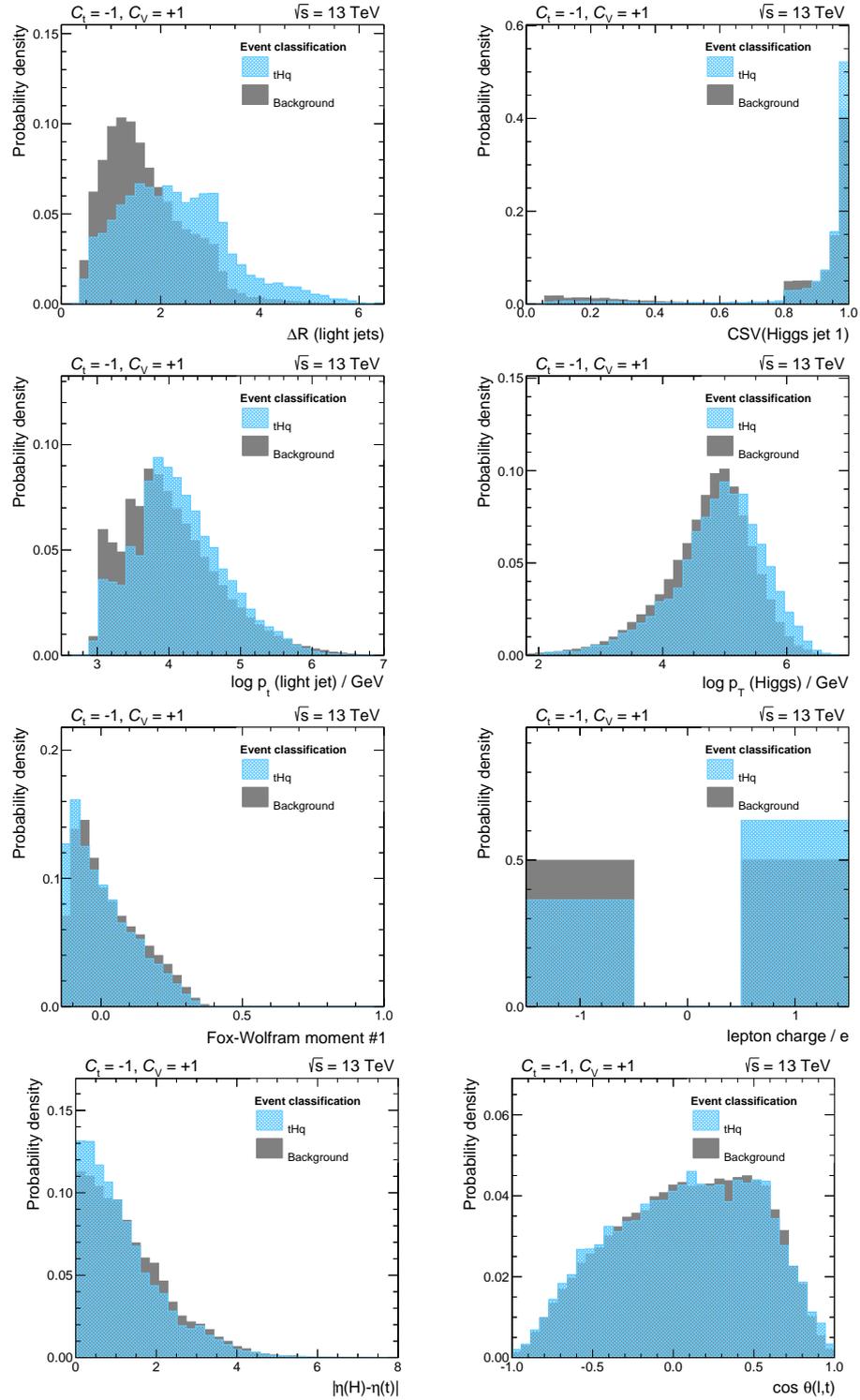


Figure B.6.: Variables ranked 8th to 15th used in the final classification of events at $\sqrt{s} = 13$ TeV. All distributions are normalized to unit area. The corresponding variable descriptions can be found in Table 6.6.

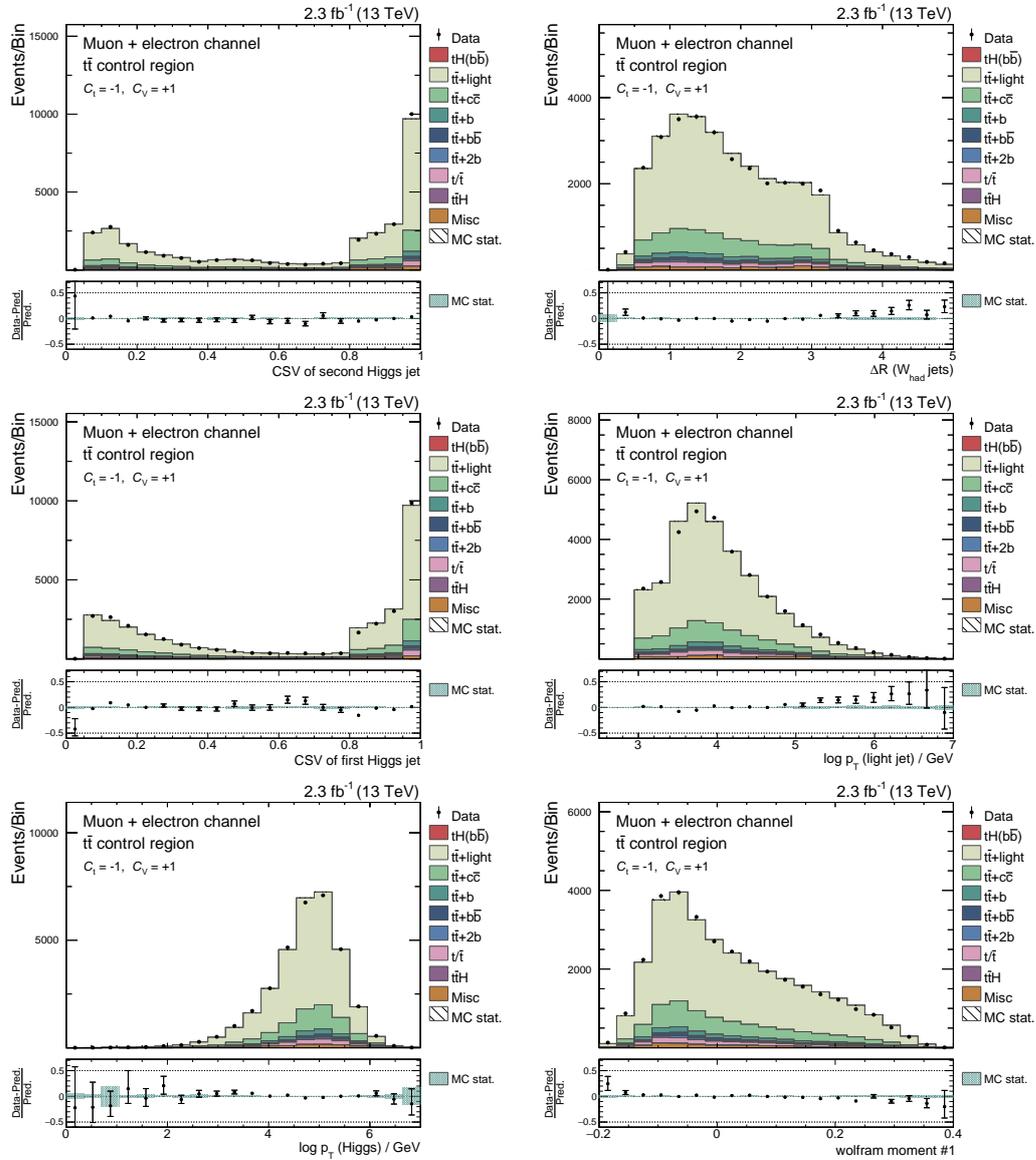


Figure B.7.: Simulation to data comparisons for input variables of the classification ranked 7th to 12th at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the $t\bar{t}$ control region for the coupling case of $C_t = -1$ and $C_V = 1$. A good agreement of simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

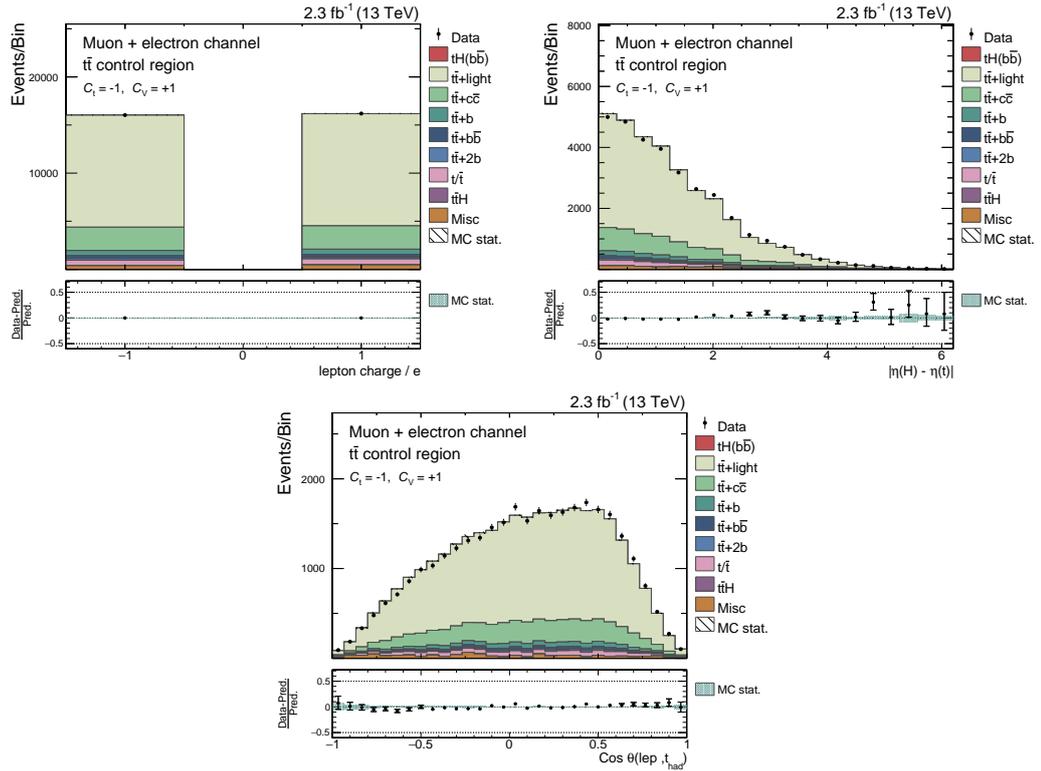


Figure B.8.: Simulation to data comparisons for input variables of the classification ranked 13th to 15th at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the $t\bar{t}$ control region for the coupling case of $C_t = -1$ and $C_V = 1$. A good agreement of simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

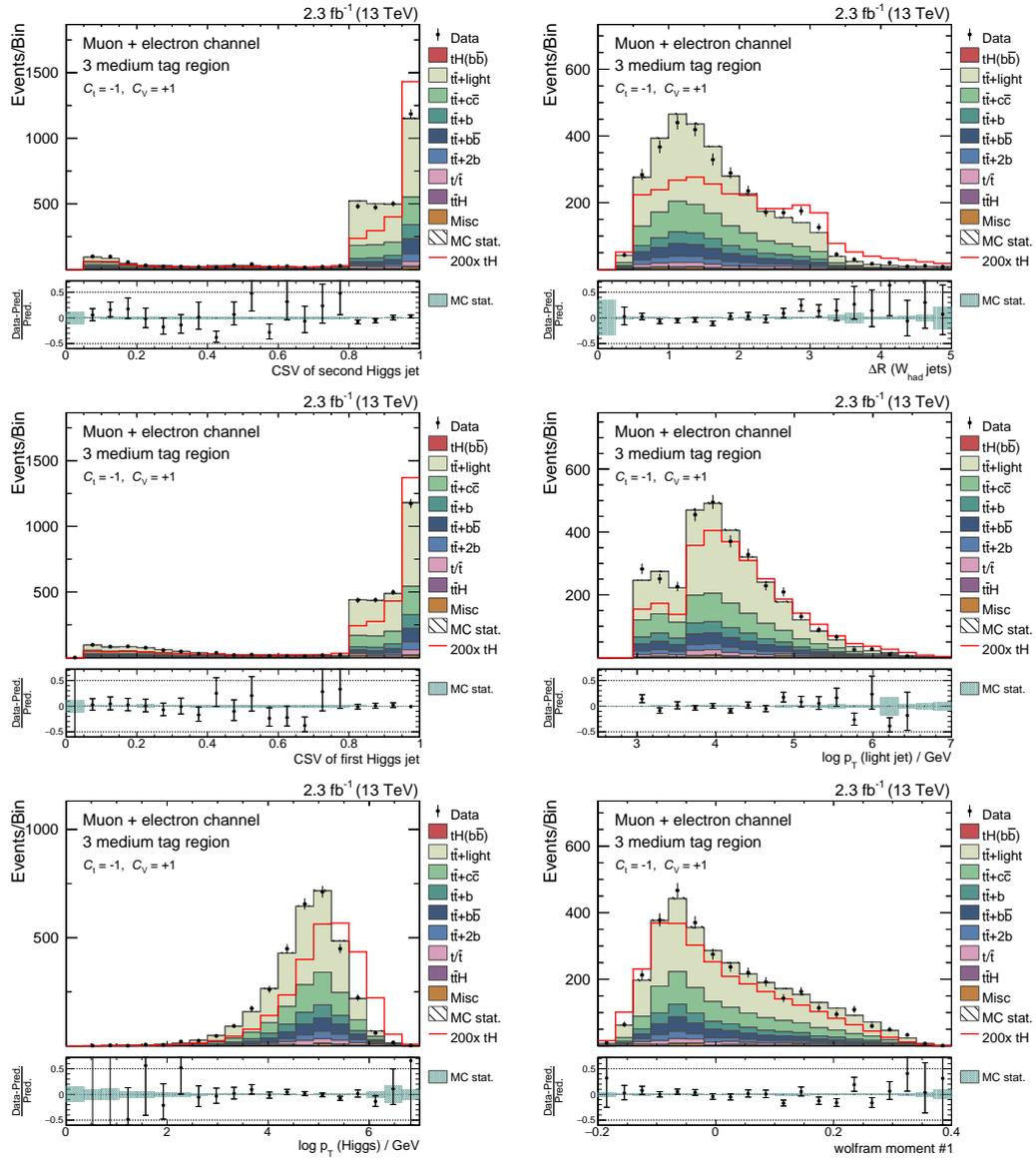


Figure B.9.: Simulation to data comparisons for input variables of the classification ranked 7th to 12th at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the 3M region for the coupling case of $C_t = -1$ and $C_V = 1$. A good agreement of simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

B. Appendix - Search for $t\bar{t}H$ Production at $\sqrt{s} = 13$ TeV

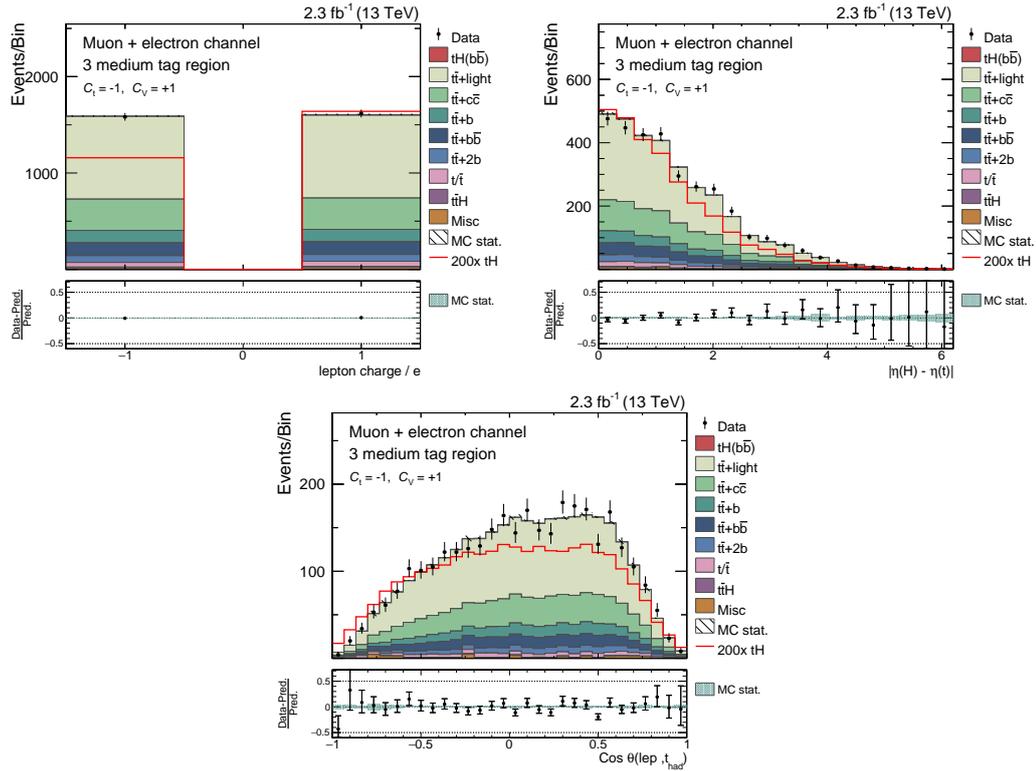


Figure B.10.: Simulation to data comparisons for input variables of the classification ranked 13th to 15th at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the 3M region for the coupling case of $C_t = -1$ and $C_V = 1$. A good agreement of simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

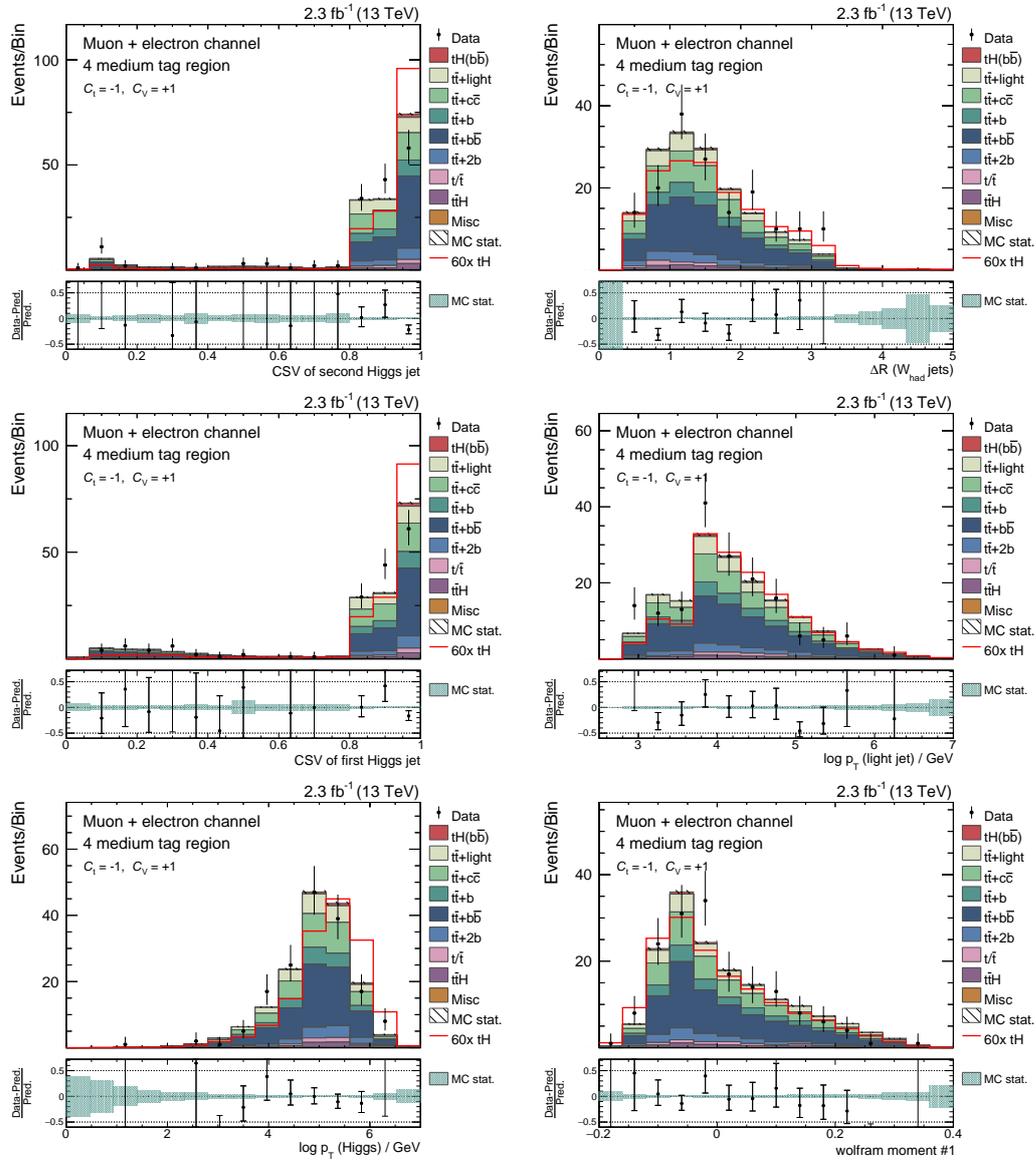


Figure B.11.: Simulation to data comparisons for input variables of the classification ranked 7th to 12th at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the 4M region for the coupling case of $C_t = -1$ and $C_V = 1$. A good agreement of simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

B. Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV

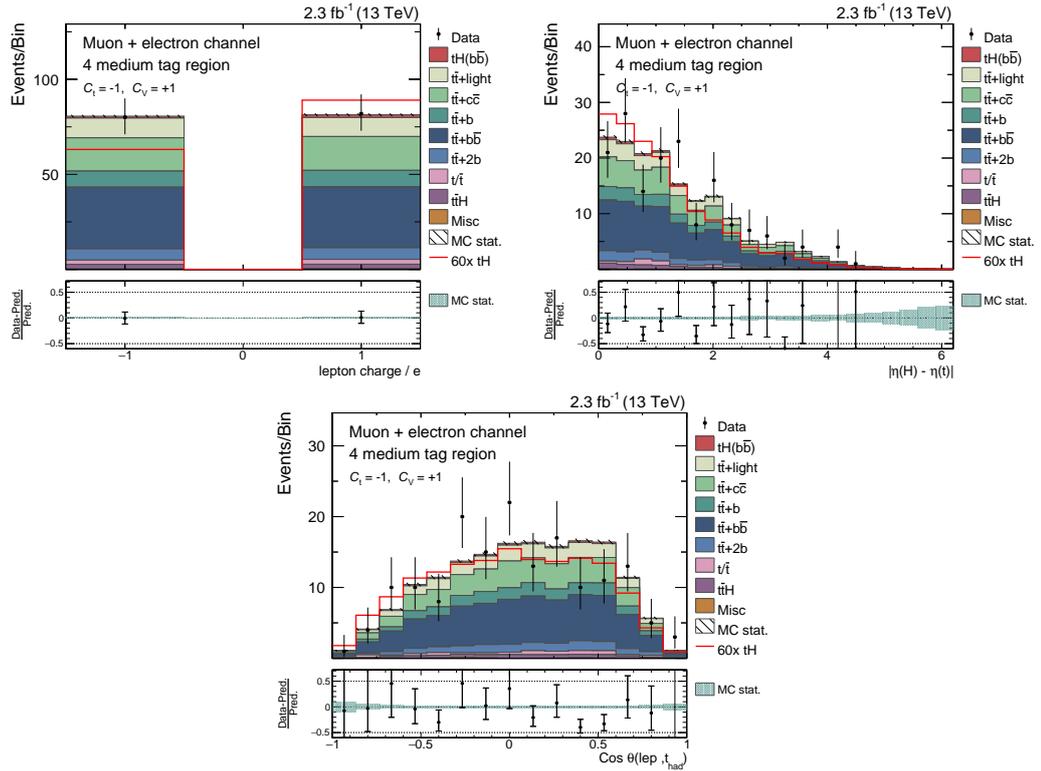


Figure B.12.: Simulation to data comparisons for input variables of the classification ranked 13th to 15th at $\sqrt{s} = 13$ TeV. Distributions are sorted by their importance in the training and are shown in the 4M region for the coupling case of $C_t = -1$ and $C_V = 1$. A good agreement of simulation and data is found. In both diagrams the simulation is scaled to match the event yields observed in data and all MC weights are applied.

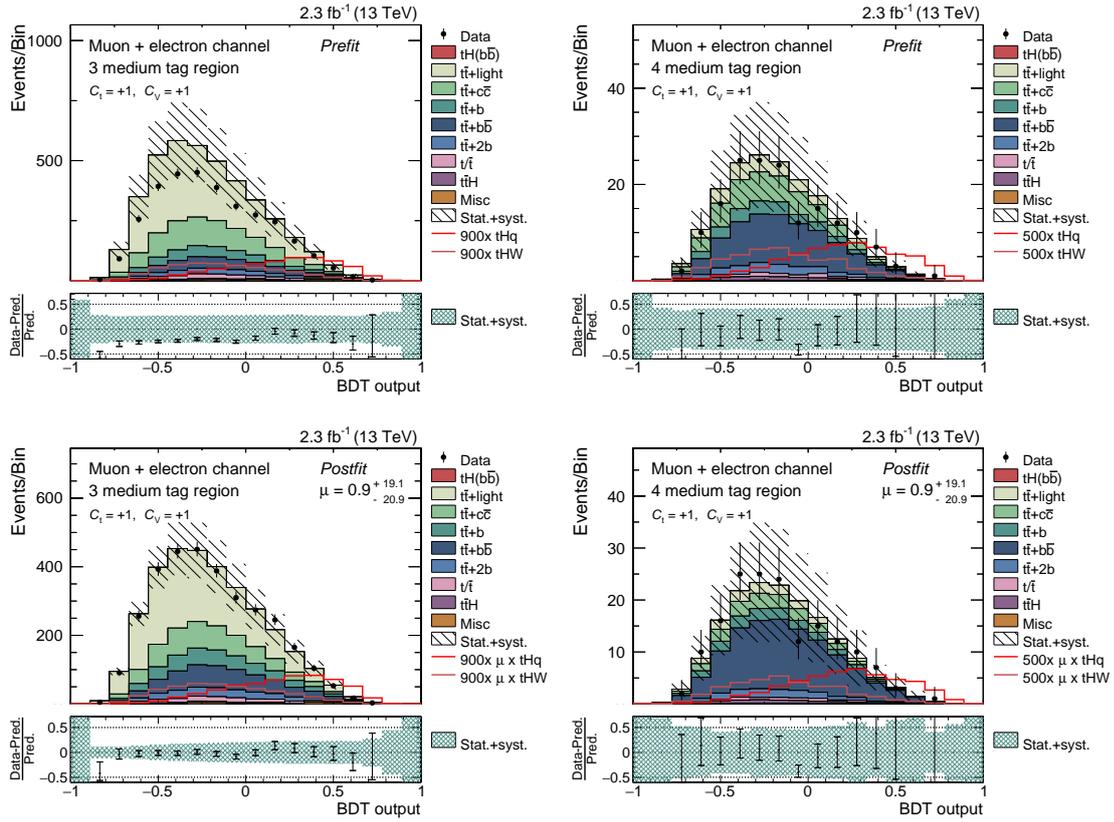


Figure B.13.: Pre- and postfit distributions of the classifier output in the 3M and 4M region for the coupling pair predicted by the SM of $C_t = +1$ and $C_V = +1$. These distributions are fitted simultaneously and a signal strength factor or $\mu = 0.9^{+19.1}_{-20.9}$ has been determined. A good agreement of simulation and data is observed after the fit.

B. Appendix - Search for tH Production at $\sqrt{s} = 13$ TeV

Table B.5.: List of all expected and observed asymptotic limits at 95% C.L. for all studied points in the C_V - C_t plane. The super- and subscribed values for the expected limit correspond to the $\pm 1\sigma$ uncertainty band values for the studied points.

C_t	$C_V = 0.5$		$C_V = 1.0$		$C_V = 1.5$	
	obs.	exp.	obs.	exp.	obs.	exp.
-3.00	2.1	$1.7^{+0.9}_{-0.5}$	2.0	$1.5^{+0.9}_{-0.5}$	2.8	$1.4^{+0.8}_{-0.5}$
-2.00	3.6	$3.3^{+1.9}_{-1.1}$	3.0	$2.6^{+1.5}_{-0.9}$	3.0	$2.5^{+1.4}_{-0.8}$
-1.50	5.5	$5.0^{+2.7}_{-1.6}$	4.9	$3.9^{+2.1}_{-1.3}$	5.3	$3.7^{+2.0}_{-1.2}$
-1.25	6.9	$6.8^{+3.8}_{-2.2}$	6.1	$4.8^{+2.6}_{-1.6}$	6.2	$4.1^{+2.3}_{-1.4}$
-1.00	10.5	$8.7^{+4.7}_{-2.9}$	7.4	$5.7^{+3.1}_{-1.8}$	7.4	$4.8^{+2.7}_{-1.6}$
-0.75	11.9	$12.7^{+7.1}_{-4.2}$	10.2	$7.5^{+4.3}_{-2.5}$	7.9	$5.9^{+3.2}_{-1.9}$
-0.50	28.1	$19.2^{+10.8}_{-6.4}$	15.4	$10.3^{+5.7}_{-3.4}$	8.5	$7.0^{+4.0}_{-2.3}$
-0.25	44.6	$32.4^{+17.9}_{-10.8}$	15.9	$13.9^{+7.9}_{-4.6}$	8.9	$9.2^{+5.0}_{-3.0}$
0.00	83.3	$64.9^{+37.0}_{-21.7}$	23.0	$19.9^{+11.6}_{-6.6}$	17.0	$12.1^{+6.8}_{-4.0}$
0.25	194.5	$171.2^{+99.0}_{-57.3}$	40.9	$32.8^{+18.4}_{-11.1}$	20.7	$15.2^{+8.5}_{-5.0}$
0.50	385.0	$324.2^{+190.0}_{-110.6}$	47.2	$55.2^{+30.6}_{-18.5}$	26.8	$21.3^{+12.0}_{-7.2}$
0.75	185.5	$144.5^{+80.1}_{-48.0}$	132.8	$80.2^{+47.0}_{-27.2}$	43.6	$31.0^{+17.7}_{-10.4}$
1.00	82.9	$62.2^{+33.5}_{-20.4}$	106.9	$97.2^{+55.4}_{-32.3}$	61.5	$41.6^{+24.1}_{-14.1}$
1.25	38.6	$29.9^{+16.1}_{-9.8}$	74.7	$77.9^{+43.8}_{-25.9}$	56.5	$54.2^{+31.4}_{-18.1}$
1.50	22.2	$18.0^{+9.8}_{-5.9}$	66.6	$43.6^{+23.8}_{-14.3}$	47.1	$59.0^{+34.1}_{-20.0}$
2.00	11.5	$8.8^{+4.8}_{-2.9}$	21.7	$18.4^{+9.9}_{-6.0}$	52.2	$36.1^{+19.4}_{-11.9}$
3.00	3.8	$3.2^{+1.7}_{-1.0}$	7.2	$5.6^{+3.1}_{-1.9}$	15.7	$11.0^{+5.9}_{-3.6}$

Table B.6.: Production cross sections for tHq, tHW and ttH at $\sqrt{s} = 13$ TeV, depending on $\cos(\alpha_{CP})$. The tHq cross sections are obtained with MADGRAPH5_AMC@NLO at NLO in the 4F scheme, whereas tHW cross sections are obtained at NLO in the 5F scheme. The quoted uncertainties on the cross section correspond to scale variations in %. The used ttH NLO cross sections are obtained from the authors of Reference [43] and are interpolated to the angles, for which the LHE weights in the signal MC samples are available. The table also lists all expected and observed asymptotic limits at 95% C.L. for all studied CP -mixing angles. The super- and subscribed values correspond to the $\pm 1\sigma$ uncertainties on the expected limit for the studied points.

$\cos(\alpha_{CP})$	Cross sections			95% C.L. limits	
	$\sigma_{\text{NLO},4F}^{\text{tHq}}$ (pb)	$\sigma_{\text{NLO},5F}^{\text{tHW}}$ (pb)	$\sigma_{\text{EXTRP. NLO}}^{\text{ttH}}$ (pb)	obs.	exp.
-1.00	$0.794^{+2.8}_{-4.0}$	$0.146^{+0.2}_{-0.2}$	0.293	11.7	$8.0^{+4.6}_{-2.8}$
-0.90	$0.728^{+2.7}_{-4.1}$	$0.135^{+0.2}_{-0.2}$	0.248	12.1	$8.0^{+4.6}_{-2.7}$
-0.80	$0.664^{+2.7}_{-4.2}$	$0.123^{+0.2}_{-0.2}$	0.207	9.7	$8.4^{+4.8}_{-2.8}$
-0.70	$0.601^{+2.8}_{-4.0}$	$0.112^{+0.2}_{-0.2}$	0.172	11.5	$9.2^{+5.1}_{-3.0}$
-0.60	$0.546^{+2.9}_{-4.3}$	$0.102^{+0.2}_{-0.2}$	0.141	14.5	$10.4^{+5.8}_{-3.4}$
-0.50	$0.497^{+3.1}_{-4.2}$	$0.092^{+0.2}_{-0.2}$	0.115	16.3	$11.8^{+6.6}_{-3.9}$
-0.40	$0.446^{+3.1}_{-4.5}$	$0.083^{+0.2}_{-0.2}$	0.093	17.6	$13.6^{+7.5}_{-4.4}$
-0.30	$0.398^{+3.2}_{-4.6}$	$0.074^{+0.2}_{-0.2}$	0.077	21.7	$15.2^{+8.5}_{-4.9}$
-0.20	$0.353^{+3.5}_{-4.8}$	$0.066^{+0.2}_{-0.2}$	0.065	22.8	$17.2^{+9.8}_{-5.6}$
-0.10	$0.314^{+3.7}_{-4.9}$	$0.059^{+0.2}_{-0.2}$	0.058	31.4	$19.8^{+11.1}_{-6.5}$
0.00	$0.275^{+3.6}_{-5.2}$	$0.052^{+0.2}_{-0.2}$	0.055	25.7	$22.7^{+12.9}_{-7.5}$
0.10	$0.242^{+4.0}_{-5.5}$	$0.045^{+0.2}_{-0.2}$	0.058	32.7	$27.1^{+15.2}_{-9.1}$
0.20	$0.211^{+4.1}_{-5.8}$	$0.040^{+0.2}_{-0.2}$	0.065	44.8	$30.9^{+17.6}_{-10.3}$
0.30	$0.182^{+4.1}_{-6.1}$	$0.035^{+0.2}_{-0.2}$	0.077	55.8	$38.1^{+21.7}_{-12.7}$
0.40	$0.156^{+4.4}_{-6.5}$	$0.030^{+0.2}_{-0.2}$	0.093	55.6	$44.9^{+25.6}_{-14.9}$
0.50	$0.134^{+4.5}_{-6.6}$	$0.026^{+0.2}_{-0.2}$	0.115	64.4	$55.6^{+31.7}_{-18.5}$
0.60	$0.116^{+4.7}_{-6.9}$	$0.023^{+0.2}_{-0.2}$	0.141	77.4	$74.2^{+42.3}_{-24.5}$
0.70	$0.100^{+5.0}_{-7.1}$	$0.020^{+0.2}_{-0.2}$	0.172	111.5	$87.2^{+49.0}_{-28.8}$
0.80	$0.087^{+4.8}_{-7.1}$	$0.018^{+0.2}_{-0.2}$	0.208	144.3	$112.2^{+63.1}_{-36.7}$
0.90	$0.077^{+4.7}_{-7.0}$	$0.017^{+0.2}_{-0.2}$	0.248	165.3	$139.0^{+77.0}_{-45.5}$
1.00	$0.071^{+4.2}_{-6.7}$	$0.016^{+0.2}_{-0.2}$	0.293	211.3	$166.2^{+89.5}_{-52.9}$

List of Figures

1.1	Higgs potential	6
1.2	Higgs boson production cross sections	8
1.3	tH production mechanisms	10
1.4	tX ₀ q and tX ₀ production cross section as function of α	12
1.5	tHq kinematics for different CP-mixing angles	13
1.6	Combined coupling fit of ATLAS and CMS	13
2.1	CERN's accelerator complex	16
2.2	CERN's dipole magnet	17
2.3	Luminosity overview	18
2.4	CMS tracker overview	21
2.5	ECAL	23
2.6	Muon system	25
2.7	Computing model	27
3.1	Factorization overview	30
3.2	CTEQ 61 PDF sets	31
3.3	Collinear and infrared safety	38
3.4	anti- k_t algorithm	39
3.5	Factorization overview	41
3.6	Displaced vertices	42
3.7	b tagging and JEC	43
4.1	Neural network illustration	46
4.2	Decision tree illustration	48
4.3	Illustration of overtraining	49
4.4	Test statistic illustration	52
5.1	Analysis workflow	54
5.2	tHq production Feynman diagram	56
5.3	Background production Feynman diagrams	57
5.4	Pileup reweighting	63
5.5	Top quark p_T reweighting	64
5.6	Jet pseudorapidity mismodeling	66
5.7	Control distributions in 2T region - electron channel	68

5.8	Most important variables in tHq reconstruction	72
5.9	Correlations of variables in tHq reconstruction	73
5.10	tHq reconstruction performance on training and test sample	74
5.11	Output of the tHq reconstruction NN for random and best hypothesis	75
5.12	Most important variables in $t\bar{t}$ reconstruction	77
5.13	Correlations of variables in $t\bar{t}$ reconstruction	78
5.14	$t\bar{t}$ reconstruction performance on training and test sample	79
5.15	Output of the $t\bar{t}$ reconstruction NN for random and best hypothesis	80
5.16	Efficiencies of the $t\bar{t}$ reconstruction	82
5.17	Signal and background shapes of classification variables	83
5.18	Classification variables in 2T region	84
5.19	Classification variables in 3T region	85
5.20	Classification variables in 4T region	86
5.21	Correlations of variables in classification	87
5.22	Classification response for training and testing set	89
5.23	Output of the classification NN	89
5.24	Shape variations caused by systematic shape uncertainties	94
5.25	Postfit distributions of the classifier output	95
5.26	Prefit and postfit pulls of nuisances	96
5.27	Impact of single nuisances on the expected limit	98
5.28	Expected and observed limit	99
5.29	Combination limit plot	102
6.1	Kinematics for different C_t values	107
6.2	tH cross sections	109
6.3	tHq kinematics	110
6.4	Pileup reweighting	113
6.5	CSV reweighting	114
6.6	Control plots in $t\bar{t}$ control region	117
6.7	Area under ROC curve for all 51 parameter points	118
6.8	Response of the tHq reconstruction classifier at 13 TeV	121
6.9	Variables in tHq reconstruction at 13 TeV - Pt. I	122
6.10	Output of the tHq reconstruction BDT for random and best hypothesis in $t\bar{t}$ control region	124
6.11	Variables in $t\bar{t}$ reconstruction at 13 TeV - Pt. I	125
6.12	Response of the $t\bar{t}$ reconstruction classifier at 13 TeV	126
6.13	Output of the $t\bar{t}$ reconstruction classifier for random and best hypothesis	126
6.14	Efficiencies of the tHq and $t\bar{t}$ reconstructions at 13 TeV	127
6.15	Comparison of different light jet assignments	128
6.16	Area under ROC curve for classification at 13 TeV	130
6.17	Response of the classification at 13 TeV	130
6.18	Signal and background shapes of classification variables at 13 TeV - Pt. I	131
6.19	Classification variables in $t\bar{t}$ control region at 13 TeV - Pt. I	132

6.20	Classification variables in 3M region at 13 TeV - Pt. I	133
6.21	Classification variables in 4M region at 13 TeV - Pt. I	134
6.22	Classification output in $t\bar{t}$ control region at 13 TeV	136
6.23	Shape variations caused by systematic shape uncertainties	138
6.24	Pre- and postfit distributions of the classifier output in the 3M and 4M region at 13 TeV	140
6.25	Pre- and postfit pulls of nuisances at 13 TeV	142
6.26	Impact of single nuisances on the expected limit at 13 TeV	143
6.27	Expected and observed limits as a function of C_t	144
6.28	Response of tX_{0q} reconstruction and classification for $\alpha = 90^\circ$	147
6.29	Area under ROC curve for 21 $C\mathcal{P}$ -mixing angles	147
6.30	Pre- and postfit distributions of the classifier output in the 3M and 4M region at 13 TeV	149
6.31	Expected and observed limits as function of the $C\mathcal{P}$ -mixing angle	150
7.1	Projection study for $tH \rightarrow b\bar{b}$	155
A.1	Control distributions in 2T region - muon channel	159
A.2	Remaining variables in tHq reconstruction	160
A.3	Output of the tHq reconstruction NN for random and best hypothesis - 2T & 3T	161
A.4	Output of the tHq reconstruction NN for random and best hypothesis - 4T	162
A.5	Remaining variables in $t\bar{t}$ reconstruction	163
A.6	Output of the $t\bar{t}$ reconstruction NN for random and best hypothesis - 2T & 3T	164
A.7	Output of the $t\bar{t}$ reconstruction NN for random and best hypothesis - 4T	165
A.8	Pre- and postfit distributions of the classifier output	166
B.1	Variables in tHq reconstruction at 13 TeV- Pt. II	170
B.2	Variables in tHq reconstruction at 13 TeV- Pt. III	171
B.3	Variables in $t\bar{t}$ reconstruction at 13 TeV- Pt. II	172
B.4	Output of the tHq reconstruction BDT for best hypothesis in 3M and 4M region	173
B.5	Output of the $t\bar{t}$ reconstruction classifier for random and best hypothesis	173
B.6	Signal and background shapes of classification variables at 13 TeV- Pt. II	174
B.7	Classification variables in 2M region at 13 TeV- Pt. II	175
B.8	Classification variables in 2M region at 13 TeV- Pt. III	176
B.9	Classification variables in 3M region at 13 TeV- Pt. II	177
B.10	Classification variables in 3M region at 13 TeV- Pt. III	178
B.11	Classification variables in 4M region at 13 TeV- Pt. II	179
B.12	Classification variables in 4M region at 13 TeV- Pt. III	180
B.13	Pre- and postfit classifier distributions for SM couplings in the 3M and 4M region	181

List of Tables

1.1	Gauge bosons of the standard model	3
1.2	Leptons in the standard model	4
5.1	Event selection at $\sqrt{s} = 8$ TeV	69
5.2	Neural network settings	70
5.3	Description of input variables of the tHq reconstruction	71
5.4	Description of input variables of the $t\bar{t}$ reconstruction	76
5.5	Description of input variables of the classification	88
5.6	PDF and QCD scale uncertainties	92
5.7	Postfit yields in the four signal channels	93
5.8	Expected and observed limits at $\sqrt{s} = 8$ TeV	98
5.9	Expected and observed limits of tH combination	101
6.1	Event selection at $\sqrt{s} = 13$ TeV	115
6.2	Description of input variables of the tHq reconstruction	119
6.3	BDT settings in reco	120
6.4	Description of input variables of the $t\bar{t}$ reconstruction	123
6.5	BDT settings in reco	129
6.6	Description of input variables of the classification	135
6.7	PDF and QCD scale uncertainties	139
6.8	Pre- and postfit yields in the 3M and 4M region at 13 TeV	141
6.9	Expected and observed limits at $\sqrt{s} = 13$ TeV	145
6.10	Description of input variables of the tX ₀ q reconstruction and classification for different $C\mathcal{P}$ -mixing angles	148
6.11	Expected and observed limits at $\sqrt{s} = 13$ TeV for fully pseudoscalar boson	149
A.1	Utilized datasets at $\sqrt{s} = 8$ TeV	157
A.2	Utilized MC simulation samples at $\sqrt{s} = 8$ TeV	158
B.1	Utilized datasets at $\sqrt{s} = 13$ TeV	167
B.2	Utilized simulation samples at $\sqrt{s} = 13$ TeV	167
B.3	Production cross sections for tHq at $\sqrt{s} = 13$ TeV	168
B.4	Production cross sections for tHW at $\sqrt{s} = 13$ TeV	169
B.5	Expected and observed limits for all 51 coupling points	182
B.6	Production cross sections for tHq for different $C\mathcal{P}$ angles	183

Bibliography

- [1] CMS COLLABORATION, “Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC,” *Physics Letters B*, vol. 716, no. 1, pp. 30–61, (2012).
- [2] ATLAS COLLABORATION, “Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC,” *Physics Letters B*, vol. 716, pp. 1–29, (2012).
- [3] F. ENGLERT AND R. BROUT, “Broken symmetry and the mass of gauge vector mesons,” *Physical Review Letters*, vol. 13, no. 9, p. 321, (1964).
- [4] P. W. HIGGS, “Broken symmetries and the masses of gauge bosons,” *Physical Review Letters*, vol. 13, no. 16, p. 508, (1964).
- [5] G. S. GURALNIK, C. R. HAGEN, AND T. W. KIBBLE, “Global conservation laws and massless particles,” *Physical Review Letters*, vol. 13, no. 20, p. 585, (1964).
- [6] T. W. KIBBLE, “Symmetry breaking in non-Abelian gauge theories,” *Physical Review*, vol. 155, no. 5, p. 1554, (1967).
- [7] B. POVH, K. RITH, C. SCHOLZ, F. ZETSCHKE, AND W. RODEJOHANN, *Teilchen und Kerne: eine Einführung in die physikalischen Konzepte*. Springer-Verlag, (2013).
- [8] M. PESKIN AND D. SCHROEDER, *An introduction to quantum field theory*. Perseus Books Publishing L.L.C., (1995).
- [9] L. SHELDON AND E. GLASHOW, “Partial-symmetries of weak interactions,” *Selected Papers on Gauge Theory of Weak and Electromagnetic Interactions*, p. 171, (1981).
- [10] S. WEINBERG, “A model of leptons,” *Physical Review Letters*, vol. 19, no. 21, p. 1264, (1967).
- [11] A. SALAM AND J. C. WARD, “Weak and electromagnetic interactions,” *Il Nuovo Cimento (1955-1965)*, vol. 11, no. 4, pp. 568–577, (1959).
- [12] E. NOETHER, “Invariante Variationsprobleme,” *Nachrichten von der Gesellschaft der Wissenschaften zu Göttingen, mathematisch-physikalische Klasse*, vol. 1918, pp. 235–257, (1918).
- [13] K. A. OLIVE ET AL., “Review of Particle Physics,” *Chinese Physics*, vol. C38, p. 090001, (2014).

- [14] SNO COLLABORATION, “Measurement of the Rate of $\nu_e + d \rightarrow p + p + e^-$ Interactions Produced by ^8B Solar Neutrinos at the Sudbury Neutrino Observatory,” *Physical Review Letters*, vol. 87, no. 7, p. 071301, (2001).
- [15] SUPER-KAMIOKANDA COLLABORATION, “Evidence for oscillation of atmospheric neutrinos,” *Physical Review Letters*, vol. 81, no. 8, p. 1562, (1998).
- [16] V. GRIBOV AND B. PONTECORVO, “Neutrino astronomy and lepton charge,” *Physics Letters B*, vol. 28, no. 7, pp. 493–496, (1969).
- [17] N. CABIBBO, “Unitary symmetry and leptonic decays,” *Physical Review Letters*, vol. 10, no. 12, p. 531, (1963).
- [18] M. KOBAYASHI AND T. MASKAWA, “CP-violation in the renormalizable theory of weak interaction,” *Progress of Theoretical Physics*, vol. 49, no. 2, pp. 652–657, (1973).
- [19] CDF COLLABORATION, “Observation of top quark production in p p collisions with the collider detector at fermilab,” *Physical Review Letters*, vol. 74, no. 14, p. 2626, (1995).
- [20] DØ COLLABORATION, “Observation of the top quark,” *Physical Review Letters*, vol. 74, no. 14, p. 2632, (1995).
- [21] ATLAS AND CMS COLLABORATIONS, “Combination of single top-quark cross-sections measurements in the t -channel at $\sqrt{s} = 8$ TeV with the ATLAS and CMS experiments,” (2013). CMS-PAS-TOP-12-002, ATLAS-CONF-2013-098.
- [22] CMS COLLABORATION, “Observation of the Associated Production of a Single Top Quark and a W Boson in pp Collisions at $\sqrt{s} = 8$ TeV,” *Physical Review Letters*, vol. 112, p. 231802, (2014).
- [23] CMS COLLABORATION, “Search for s channel single top quark production in pp collisions at $\sqrt{s} = 7$ and 8 TeV,” (2016). CMS-PAS-TOP-13-009-003.
- [24] J. GOLDSTONE, A. SALAM, AND S. WEINBERG, “Broken symmetries,” *Physical Review*, vol. 127, no. 3, p. 965, (1962).
- [25] F. BEZRUKOV AND M. SHAPOSHNIKOV, “Why should we care about the top quark Yukawa coupling?,” *Journal of Experimental and Theoretical Physics*, vol. 120, pp. 335–343, (2015).
- [26] G. DEGRASSI *ET AL.*, “Higgs mass and vacuum stability in the Standard Model at NNLO,” *Journal of High Energy Physics*, vol. 08, p. 098, (2012).
- [27] ATLAS AND CMS COLLABORATIONS, “Combined Measurement of the Higgs Boson Mass in pp Collisions at $\sqrt{s} = 7$ and 8 TeV with the ATLAS and CMS Experiments,” *Physical Review Letters*, vol. 114, p. 191803, (2015).

-
- [28] CMS COLLABORATION, “Search for Higgs boson off-shell production in proton-proton collisions at 7 and 8 TeV and derivation of constraints on its total decay width,” (2016). CMS-PAS-HIG-14-032.
- [29] L. LANDAU, “The moment of a 2-photon system,” *Proceedings of the USSR Academy of Sciences*, vol. 60, no. 207, pp. 12–3, (1948).
- [30] C.-N. YANG, “Selection rules for the dematerialization of a particle into two photons,” *Physical Review*, vol. 77, no. 2, p. 242, (1950).
- [31] CMS COLLABORATION, “Measurement of the properties of a Higgs boson in the four-lepton final state,” *Physical Review D*, vol. 89, no. 9, p. 092007, (2014).
- [32] LHC HIGGS CROSS SECTION WORKING GROUP, “Higgs cross sections and decay branching ratios.” <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/LHCHXSWG>, last accessed on 26.05.2016.
- [33] CMS COLLABORATION, “SM Higgs Branching Ratios and Total Decay Widths (update in CERN Report4 2016).” <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/CERNYellowReportPageBR>, last accessed on 05.05.2016.
- [34] F. MALTONI, K. PAUL, T. STELZER, AND S. WILLENBROCK, “Associated production of Higgs and single top at hadron colliders,” *Physical Review D*, vol. 64, p. 094023, (2001).
- [35] S. BISWAS, E. GABRIELLI, F. MARGAROLI, AND B. MELE, “Direct constraints on the top-Higgs coupling from the 8 TeV LHC data,” *Journal of High Energy Physics*, vol. 07, p. 073, (2013).
- [36] CMS COLLABORATION, “Modelling of the single top-quark production in association with the Higgs boson at 13 TeV.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/SingleTopHiggsGeneration13TeV>, last accessed on 07.05.2016.
- [37] CMS COLLABORATION, “SM Higgs production cross sections at $\sqrt{s} = 13$ TeV (update in CERN Report4 2016).” <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/CERNYellowReportPageAt13TeV>, last accessed on 05.05.2016.
- [38] M. FARINA, C. GROJEAN, F. MALTONI, E. SALVIONI, AND A. THAMM, “Lifting degeneracies in Higgs couplings using single top production in association with a Higgs boson,” *Journal of High Energy Physics*, vol. 05, p. 022, (2013).
- [39] S. BISWAS, E. GABRIELLI, AND B. MELE, “Single top and Higgs associated production as a probe of the Htt coupling sign at the LHC,” *Journal of High Energy Physics*, vol. 01, p. 088, (2013).
- [40] J. ELLIS, D. S. HWANG, K. SAKURAI, AND M. TAKEUCHI, “Disentangling Higgs-top couplings in associated production,” *Journal of High Energy Physics*, vol. 2014, no. 4, pp. 1–19, (2014).

- [41] J. YUE, “Enhanced h signal at the LHC with $h \rightarrow \gamma\gamma$ decay and CP-violating top–Higgs coupling,” *Physics Letters B*, vol. 744, pp. 131–136, (2015).
- [42] A. KOBAKHIDZE, L. WU, AND J. YUE, “Anomalous Top-Higgs Couplings and Top Polarisation in Single Top and Higgs Associated Production at the LHC,” *Journal of High Energy Physics*, vol. 10, p. 100, (2014).
- [43] F. DEMARTIN, F. MALTONI, K. MAWATARI, AND M. ZARO, “Higgs production in association with a single top quark at the LHC,” *European Physical Journal C*, vol. 75, p. 267, (2015).
- [44] J. ELLIS AND T. YOU, “Updated Global Analysis of Higgs Couplings,” *Journal of High Energy Physics*, vol. 06, p. 103, (2013).
- [45] ATLAS AND CMS COLLABORATIONS, “Measurements of the Higgs boson production and decay rates and constraints on its couplings from a combined ATLAS and CMS analysis of the LHC pp collision data at $\sqrt{s} = 7$ and 8 TeV,” (2016). CERN-EP-2016-100, ATLAS-HIGG-2015-07, CMS-HIG-15-002.
- [46] ATLAS COLLABORATION, “Search for $H \rightarrow \gamma\gamma$ produced in association with top quarks and constraints on the Yukawa coupling between the top quark and the Higgs boson using data taken at 7 TeV and 8 TeV with the ATLAS detector,” *Physics Letters B*, vol. 740, pp. 222–242, (2015).
- [47] ATLAS COLLABORATION, “The ATLAS Experiment at the CERN Large Hadron Collider,” *Journal of Instrumentation*, vol. 3, no. 08, p. S08003, (2008).
- [48] CMS COLLABORATION, “The CMS experiment at the CERN LHC,” *Journal of Instrumentation*, vol. 3, no. 08, p. S08004, (2008).
- [49] LHCb COLLABORATION, “The LHCb Detector at the LHC,” *Journal of Instrumentation*, vol. 3, no. 08, p. S08005, (2008).
- [50] ALICE COLLABORATION, “The ALICE experiment at the CERN LHC,” *Journal of Instrumentation*, vol. 3, no. 08, p. S08002, (2008).
- [51] C. LEFÈVRE, “LHC: the guide (English version). Guide du LHC (version anglaise),” (2009). CERN-Brochure-2009-003-Eng.
- [52] C. LEFÈVRE, “The CERN accelerator complex,” (2008). CERN-DI-0812015.
- [53] AC TEAM, “Diagram of an LHC dipole magnet. Schéma d’un aimant dipôle du LHC,” (1999). CERN-DI-9906025.
- [54] K. FORAZ ET AL., “LS1 “First Long Shutdown of LHC and its Injector Chains”,” *Proceedings, 5th International Particle Accelerator Conference (IPAC 2014)*, p. TUPRO007, (2014).
- [55] CMS COLLABORATION, “Long Shutdown 1 | CERN timelines.” <http://timeline.web.cern.ch/timelines/Long-Shutdown-1>, last accessed on 14.04.2016.

-
- [56] CMS COLLABORATION, “Public CMS Luminosity Information.” <https://twiki.cern.ch/twiki/bin/view/CMSPublic/LumiPublicResults>, last accessed on 14.04.2016.
- [57] CMS COLLABORATION, “CMS Physics : Technical Design Report Volume 1: Detector Performance and Software,” (2006). CMS-TDR-8-1.
- [58] A. DOMINGUEZ *ET AL.*, “CMS Technical Design Report for the Pixel Detector Upgrade,” (2012). CERN-LHCC-2012-016. CMS-TDR-11.
- [59] CMS COLLABORATION, “A cool Tracker! Cold operation established for subdetector | CMS Experiment.” <http://cms.web.cern.ch/news/cool-tracker-cold-operation-established-subdetector>, last accessed on 14.04.2016.
- [60] CMS COLLABORATION, “Description and performance of track and primary-vertex reconstruction with the CMS tracker,” *Journal of Instrumentation*, vol. 9, no. 10, p. P10009, (2014).
- [61] CMS COLLABORATION, “ECAL preshower removed from detector for repairs | CMS Experiment.” <http://cms.web.cern.ch/news/ecal-preshower-removed-detector-repairs>, last accessed on 14.04.2016.
- [62] W. ADAM *ET AL.*, “The CMS high level trigger,” *European Physical Journal C*, vol. 46, pp. 605–667, (2006).
- [63] WLCG PROJECT OFFICE, “Documents & Reference - Tiers,” (2014). <http://wlcg.web.cern.ch/documents-reference>, last accessed on 14.04.2016.
- [64] M. A. DOBBS *ET AL.*, “Les Houches guidebook to Monte Carlo generators for hadron collider physics,” pp. 411–459, (2004).
- [65] P. NASON AND B. WEBBER, “Next-to-Leading-Order Event Generators,” *Annual Review of Nuclear and Particle Science*, vol. 62, pp. 187–213, (2012).
- [66] Y. L. DOKSHITZER, “Calculation of structure functions of deep-inelastic scattering and e^+e^- annihilation by perturbation theory in quantum chromodynamics,” *Soviet Journal of Nuclear Physics*, vol. 46, p. 641, (1977).
- [67] V. N. GRIBOV AND L. N. LIPATOV, “Deep inelastic ep-scattering in a perturbation theory,” *Soviet Journal of Nuclear Physics*, vol. 15, pp. 438–450, (1972).
- [68] G. ALTARELLI AND G. PARISI, “Asymptotic Freedom in Parton Language,” *Nuclear Physics B*, vol. 126, p. 298, (1977).
- [69] S. DULAT *ET AL.*, “New parton distribution functions from a global analysis of quantum chromodynamics,” *Physical Review D*, vol. 93, no. 3, p. 033006, (2016).
- [70] A. D. MARTIN, W. J. STIRLING, R. S. THORNE, AND G. WATT, “Parton distributions for the LHC,” *European Physical Journal C*, vol. 63, pp. 189–285, (2009).

- [71] S. FORTE, L. GARRIDO, J. I. LATORRE, AND A. PICCIONE, “Neural network parametrization of deep inelastic structure functions,” *Journal of High Energy Physics*, vol. 05, p. 062, (2002).
- [72] A. BUCKLEY *ET AL.*, “LHAPDF6: parton density access in the LHC precision era,” *European Physical Journal C*, vol. 75, p. 132, (2015).
- [73] DURHAM HEPDATA PROJECT, “The Durham HepData Project - PDF Plotter.” <http://hepdata.cedar.ac.uk/pdf/pdf3.html>, last accessed on 26.05.2016.
- [74] V. V. SUDAKOV, “Vertex parts at very high energies in quantum electrodynamics,” *Journal of Experimental and Theoretical Physics*, vol. 3, pp. 65–71, (1956).
- [75] S. CATANI, F. KRAUSS, R. KUHN, AND B. R. WEBBER, “QCD matrix elements + parton showers,” *Journal of High Energy Physics*, vol. 11, p. 063, (2001).
- [76] M. L. MANGANO, M. MORETTI, F. PICCININI, AND M. TRECCANI, “Matching matrix elements and shower evolution for top-quark production in hadronic collisions,” *Journal of High Energy Physics*, vol. 01, p. 013, (2007).
- [77] R. FREDERIX AND S. FRIXIONE, “Merging meets matching in MC@NLO,” *Journal of High Energy Physics*, vol. 12, p. 061, (2012).
- [78] F. MALTONI, G. RIDOLFI, AND M. UBIALI, “b-initiated processes at the LHC: a reappraisal,” *Journal of High Energy Physics*, vol. 07, p. 022, (2012).
- [79] B. ANDERSSON, G. GUSTAFSON, G. INGELMAN, AND T. SJÖSTRAND, “Parton Fragmentation and String Dynamics,” *Physics Reports*, vol. 97, pp. 31–145, (1983).
- [80] T. STELZER AND W. F. LONG, “Automatic generation of tree level helicity amplitudes,” *Computer Physics Communications*, vol. 81, pp. 357–371, (1994).
- [81] J. ALWALL *ET AL.*, “MadGraph 5 : Going Beyond,” *Journal of High Energy Physics*, vol. 06, p. 128, (2011).
- [82] P. ARTOISENET, R. FREDERIX, O. MATTELAER, AND R. RIETKERK, “Automatic spin-entangled decays of heavy resonances in Monte Carlo simulations,” *Journal of High Energy Physics*, vol. 03, p. 015, (2013).
- [83] J. ALWALL *ET AL.*, “The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations,” *Journal of High Energy Physics*, vol. 07, p. 079, (2014).
- [84] P. NASON, “A New method for combining NLO QCD with shower Monte Carlo algorithms,” *Journal of High Energy Physics*, vol. 11, p. 040, (2004).
- [85] S. FRIXIONE, P. NASON, AND C. OLEARI, “Matching NLO QCD computations with Parton Shower simulations: the POWHEG method,” *Journal of High Energy Physics*, vol. 11, p. 070, (2007).

-
- [86] S. ALIOLI, P. NASON, C. OLEARI, AND E. RE, “A general framework for implementing NLO calculations in shower Monte Carlo programs: the POWHEG BOX,” *Journal of High Energy Physics*, vol. 06, p. 043, (2010).
- [87] T. SJÖSTRAND AND P. Z. SKANDS, “Transverse-momentum-ordered showers and interleaved multiple interactions,” *European Physical Journal C*, vol. 39, pp. 129–154, (2005).
- [88] CMS COLLABORATION, “Study of the underlying event, b-quark fragmentation and hadronization properties in $t\bar{t}$ events,” (2013). CMS-PAS-TOP-13-007.
- [89] CMS COLLABORATION, “Event generator tunes obtained from underlying event and multiparton scattering measurements,” *European Physical Journal C*, vol. 76, no. 3, p. 155, (2016).
- [90] T. SJÖSTRAND, S. MRENNNA, AND P. Z. SKANDS, “PYTHIA 6.4 Physics and Manual,” *Journal of High Energy Physics*, vol. 05, p. 026, (2006).
- [91] T. SJÖSTRAND, S. MRENNNA, AND P. Z. SKANDS, “A Brief Introduction to PYTHIA 8.1,” *Computer Physics Communications*, vol. 178, pp. 852–867, (2008).
- [92] N. DAVIDSON *ET AL.*, “Universal Interface of TAUOLA Technical and Physics Documentation,” *Computer Physics Communications*, vol. 183, pp. 821–843, (2012).
- [93] S. AGOSTINELLI *ET AL.*, “GEANT4: A Simulation toolkit,” *Nuclear Instruments and Methods in Physics*, vol. A506, pp. 250–303, (2003).
- [94] J. ALLISON *ET AL.*, “Geant4 developments and applications,” *IEEE Transactions on Nuclear Science*, vol. 53, p. 270, (2006).
- [95] CMS COLLABORATION, “Particle-Flow Event Reconstruction in CMS and Performance for Jets, Taus, and MET,” (2009). CMS-PAS-PFT-09-001.
- [96] CMS COLLABORATION, “Commissioning of the Particle-flow Event Reconstruction with the first LHC collisions recorded in the CMS detector,” (2010). CMS-PAS-PFT-10-001.
- [97] T. SPEER *ET AL.*, “Vertex Fitting in the CMS Tracker,” (2006). CMS-NOTE-2006-032.
- [98] CMS COLLABORATION, “Tracking and Primary Vertex Results in First 7 TeV Collisions,” (2010). CMS-PAS-TRK-10-005.
- [99] K. ROSE, “Deterministic annealing for clustering, compression, classification, regression and related optimization problems,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2210–2238, (1998).
- [100] R. FRÜHWIRTH, W. WALTENBERGER, AND P. VANLAER, “Adaptive Vertex Fitting,” (2007). CMS-NOTE-2007-008.
- [101] W. ADAM, R. FRÜHWIRTH, A. STRANDLIE, AND T. TODOROV, “Reconstruction of electrons with the Gaussian-sum filter in the CMS tracker at the LHC,” *Journal of Physics G: Nuclear and Particle Physics*, vol. 31, no. 9, p. N9, (2005).

- [102] M. CACCIARI, G. P. SALAM, AND G. SOYEZ, “The anti- k_t jet clustering algorithm,” *Journal of High Energy Physics*, vol. 04, p. 063, (2008).
- [103] S. CATANI, Y. L. DOKSHITZER, M. H. SEYMOUR, AND B. R. WEBBER, “Longitudinally invariant K_t clustering algorithms for hadron hadron collisions,” *Nuclear Physics B*, vol. 406, pp. 187–224, (1993).
- [104] Y. L. DOKSHITZER, G. D. LEDER, S. MORETTI, AND B. R. WEBBER, “Better jet clustering algorithms,” *Journal of High Energy Physics*, vol. 08, p. 001, (1997).
- [105] M. CACCIARI AND G. P. SALAM, “Dispelling the N^3 myth for the k_t jet-finder,” *Physics Letters B*, vol. 641, pp. 57–61, (2006).
- [106] M. CACCIARI, G. P. SALAM, AND G. SOYEZ, “FastJet User Manual,” *European Physical Journal C*, vol. 72, p. 1896, (2012).
- [107] CMS COLLABORATION, “Determination of jet energy calibration and transverse momentum resolution in CMS,” *Journal of Instrumentation*, vol. 6, no. 11, p. P11002, (2011).
- [108] CMS COLLABORATION, “Introduction to Jet Energy Corrections at CMS.” <https://twiki.cern.ch/twiki/bin/view/CMS/IntroToJEC>, last accessed on 14.04.2016, last accessed on 17.05.2016.
- [109] CMS COLLABORATION, “CMS JEC Run I legacy performance plots,” (2015). CMS-DP-2015-044.
- [110] CMS COLLABORATION, “Determination of jet energy calibration and transverse momentum resolution in CMS,” *Journal of Instrumentation*, vol. 6, p. 11002, (2011).
- [111] CMS COLLABORATION, “Identification of b-quark jets with the CMS experiment,” *Journal of Instrumentation*, vol. 8, p. P04013, (2013).
- [112] CMS COLLABORATION, “Performance of b tagging at $\sqrt{s} = 8$ TeV in multijet, $t\bar{t}$ and boosted topology events,” (2013). CMS-PAS-BTV-13-001.
- [113] CMS COLLABORATION, “Performance of b-tagging algorithms at 13 TeV,” (2016). CMS-AN-2016/036.
- [114] CMS COLLABORATION, “Performance of b tagging in boosted topology events.” <https://twiki.cern.ch/twiki/bin/view/CMSPublic/BoostedBTaggingPlots2015>, last accessed on 14.04.2016.
- [115] CMS COLLABORATION, “MET Analysis.” https://twiki.cern.ch/twiki/bin/view/CMSPublic/WorkBookMetAnalysis\#7_7_6_MET_Corrections, last accessed on 14.04.2016.
- [116] CMS COLLABORATION, “METType1Type2Formulae.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/METType1Type2Formulae>, last accessed on 14.04.2016.

-
- [117] A. HOECKER *ET AL.*, “TMVA-Toolkit for multivariate data analysis,” *arXiv preprint physics/0703039*, (2007).
- [118] I. ANTICHEVA *ET AL.*, “ROOT – A C++ framework for petabyte data storage, statistical analysis and visualization,” *Computer Physics Communications*, vol. 182, no. 6, pp. 1384–1385, (2011).
- [119] C. G. BROYDEN, “The convergence of a class of double-rank minimization algorithms 1. general considerations,” *IMA Journal of Applied Mathematics*, vol. 6, no. 1, pp. 76–90, (1970).
- [120] R. FLETCHER, “A new approach to variable metric algorithms,” *The Computer Journal*, vol. 13, no. 3, pp. 317–322, (1970).
- [121] D. GOLDFARB, “A family of variable-metric methods derived by variational means,” *Mathematics of Computation*, vol. 24, no. 109, pp. 23–26, (1970).
- [122] D. F. SHANNO, “Conditioning of quasi-Newton methods for function minimization,” *Mathematics of Computation*, vol. 24, no. 111, pp. 647–656, (1970).
- [123] L. BREIMAN, J. FRIEDMAN, C. J. STONE, AND R. A. OLSHEN, *Classification and regression trees*. CRC press, (1984).
- [124] Y. FREUND AND R. E. SCHAPIRE, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of computer and system sciences*, vol. 55, no. 1, pp. 119–139, (1997).
- [125] V. BLOBEL AND E. LOHRMANN, *Statistische und numerische Methoden der Datenanalyse*. Springer-Verlag, (2013).
- [126] E. GROSS, “Search and Discovery Statistics in HEP.” ESHEP, Lecture (2015), retrieved from https://indico.cern.ch/event/381289/contributions/1807986/attachments/1149483/1655277/Statistics_CERN_2015.pdf.
- [127] W. VERKERKE, “Systematic uncertainties and profiling.” <https://indico.in2p3.fr/event/9742/contribution/16/material/1/0.pdf>, Lecture, last accessed on 26.05.2016.
- [128] G. COWAN, “Statistics for Searches at the LHC,” *Proceedings, 69th Scottish Universities Summer School in Physics : LHC Phenomenology (SUSSP69)*, (2013).
- [129] M. BAAK, S. GADATSCH, R. HARRINGTON, AND W. VERKERKE, “Interpolation between multi-dimensional histograms using a new non-linear moment morphing method,” *Nuclear Instruments and Methods in Physics A*, vol. 771, pp. 39–48, (2015).
- [130] J. NEYMAN AND E. S. PEARSON, *On the problem of the most efficient tests of statistical hypotheses*. Springer-Verlag, (1992).
- [131] A. L. READ, “Modified frequentist analysis of search results (the CL_s method),” (2000). CERN-OPEN-2000-205.

- [132] G. COWAN, K. CRANMER, E. GROSS, AND O. VITELLS, “Asymptotic formulae for likelihood-based tests of new physics,” *European Physical Journal C*, vol. 71, p. 1554, (2011).
- [133] S. S. WILKS, “The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses,” *Annals of Mathematical Statistics*, vol. 9, pp. 60–62, (03, 1938).
- [134] A. WALD, “Tests of statistical hypotheses concerning several parameters when the number of observations is large,” *Transactions of the American Mathematical society*, vol. 54, no. 3, pp. 426–482, (1943).
- [135] CMS COLLABORATION, “Documentation of the RooStats-based statistics tools for Higgs PAG.” <https://twiki.cern.ch/twiki/bin/view/CMS/SWGuideHiggsAnalysisCombinedLimit>, last accessed on 26.05.2016.
- [136] W. VERKERKE AND D. KIRKBY, “The RooFit toolkit for data modeling,” *arXiv preprint physics*, (2003).
- [137] CMS COLLABORATION, “Search for H to bbbar in association with single top quarks as a test of Higgs couplings,” (2014). CMS-PAS-HIG-14-015.
- [138] CMS COLLABORATION, “Search for the Associated Production of a Higgs Boson with a Single Top Quark in Proton-Proton Collisions at $\sqrt{s} = 8$ TeV,” (2015). CMS-PAS-HIG-14-027.
- [139] N. FALTERMANN, *Search for Standard Model Higgs boson production in association with a single top quark with the CMS experiment*. (2015). Diploma thesis, KIT, EKP-2015-00020.
- [140] CMS COLLABORATION, “Measurement of the t -channel single-top-quark production cross section and of the $|V_{tb}|$ CKM matrix element in pp collisions at $\sqrt{s} = 8$ TeV,” *Journal of High Energy Physics*, vol. 06, p. 090, (2014).
- [141] CMS COLLABORATION, “Measurement of the inclusive cross section of single top-quark production in the t -channel at 13 TeV,” (2016). CMS-PAS-TOP-16-003.
- [142] CMS COLLABORATION, “Search for $t\bar{t}H$ production in the $H \rightarrow b\bar{b}$ decay channel with $\sqrt{s} = 13$ TeV pp collisions at the CMS experiment,” (2016). CMS-PAS-HIG-16-004.
- [143] CMS COLLABORATION, “Search for the standard model Higgs boson produced in association with top quarks in multilepton final states,” (2013). CMS-PAS-HIG-13-020.
- [144] C. BÖSER, *Search for the Higgs Boson in WH and tHq Production Modes with the CMS Experiment*. (2015). PhD thesis, KIT, IEKP-KA/2015-06.
- [145] DATA CERTIFICATION TEAM, “JSON files for 22 January ReReco (2012ABCD): Runs 190456 - 208686 .” <https://hypernews.cern.ch/HyperNews/CMS/get/physics-validation/2065.html>, last accessed on 27.04.2016.
- [146] CMS COLLABORATION, “Measurement of the cross section ratio $\sigma_{t\bar{t}b\bar{b}}/\sigma_{t\bar{t}jj}$ in pp collisions at $\sqrt{s} = 8$ TeV,” *Physics Letters B*, vol. 746, pp. 132–153, (2015).

-
- [147] ATLAS COLLABORATION, “Study of heavy-flavor quarks produced in association with top-quark pairs at $\sqrt{s}=7$ TeV using the ATLAS detector,” *Physical Review D*, vol. 89, no. 7, p. 072012, (2014).
- [148] CMS COLLABORATION, “First Measurement of the Cross Section Ratio $\sigma(t\bar{t}b\bar{b})/\sigma(t\bar{t}jj)$ in pp Collisions at $\sqrt{s} = 7$ TeV,” (2012). CMS-PAS-TOP-12-024.
- [149] CMS COLLABORATION, “Search for Higgs Boson Production in Association with a Top-Quark Pair and Decaying to Bottom Quarks or Tau Leptons,” (2013). CMS-PAS-HIG-13-019.
- [150] CMS COLLABORATION, “Search for $H \rightarrow b\bar{b}$ in association with single top quarks as a test of Higgs boson couplings,” (2016). CMS-AN-2013/113.
- [151] H. KIRSCHENMANN, “Jet performance in CMS,” (2013). CMS-CR-2013-325.
- [152] CMS COLLABORATION, “Baseline muon selections for Run-I.” <https://twiki.cern.ch/twiki/bin/view/CMSPublic/SWGGuideMuonId>, last accessed on 27.04.2016.
- [153] CMS COLLABORATION, “Multivariate Electron Identification.” https://twiki.cern.ch/twiki/bin/view/CMS/MultivariateElectronIdentification#Recommended_Working_Points_With, last accessed on 07.05.2016.
- [154] CMS COLLABORATION, “Performance of the Particle-Flow jet identification criteria using proton-proton collisions at $\sqrt{s} = 8$ TeV,” (2014). CMS-AN-14/227.
- [155] CMS COLLABORATION, “Recommended Jet Energy Corrections and Uncertainties For Data and MC.” <https://twiki.cern.ch/twiki/bin/view/CMS/JECDataMC>, last accessed on 27.04.2016.
- [156] T. CHWALEK, *Measurement of W-boson helicity fractions in top quark decays with the CDF II experiment and prospects for an early $t\bar{t}$ cross-section measurement with the CMS experiment.* (2010). PhD thesis, KIT, CERN-THESIS-2010-255.
- [157] CMS COLLABORATION, “Reference muon id and isolation efficiencies,” (2014). <https://twiki.cern.ch/twiki/bin/view/CMS/MuonReferenceEffs>, last accessed on 27.04.2016.
- [158] N. KIDONAKIS, “NNLL threshold resummation for top-pair and single-top production,” *Physics of Particles and Nuclei*, vol. 45, no. 4, pp. 714–722, (2014).
- [159] CMS COLLABORATION, “b Tag & Vertexing Physics Object Group.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/BtagPOG>, last accessed on 10.05.2016.
- [160] CMS COLLABORATION, “Methods to apply b-tagging efficiency scale factors.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/BTagSFMethods>, last accessed on 14.04.2016.
- [161] D. KROHN, M. D. SCHWARTZ, T. LIN, AND W. J. WAALEWIJN, “Jet Charge at the LHC,” *Physical Review Letters*, vol. 110, no. 21, p. 212001, (2013).

- [162] CMS COLLABORATION, “CMS Luminosity Based on Pixel Cluster Counting - Summer 2013 Update,” (2013). CMS-PAS-LUM-13-001.
- [163] CMS COLLABORATION, “Estimating Systematic Errors Due to Pileup Modeling.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/PileupSystematicErrors>, last accessed on 26.05.2016.
- [164] JETMET POG, “Jet Energy Resolution,” (2014). https://twiki.cern.ch/twiki/bin/view/CMS/JetResolution#JER_Uncertainty, last accessed on 27.04.2016.
- [165] CMS COLLABORATION, “Jet energy scale uncertainty sources.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/JECUncertaintySources>, last accessed on 27.04.2016.
- [166] CMS COLLABORATION, “b-quark JES and Hadronization Modeling.” https://twiki.cern.ch/twiki/bin/viewauth/CMS/TopMassSystematics#b_quark_JES_and_Hadronization_Mo, last accessed on 27.04.2016.
- [167] R. J. BARLOW AND C. BEESTON, “Fitting using finite Monte Carlo samples,” *Computer Physics Communication*, vol. 77, pp. 219–228, (1993).
- [168] J. S. CONWAY, “Incorporating Nuisance Parameters in Likelihoods for Multisource Spectra,” in *Proceedings, PHYSTAT 2011 Workshop on Statistical Issues Related to Discovery Claims in Search Experiments and Unfolding, CERN, Geneva, Switzerland 17-20 January 2011*, (2011).
- [169] CMS COLLABORATION, “Search for associated production of a single top quark and a Higgs boson in events where the Higgs boson decays to two photons at $\sqrt{s} = 8$ TeV,” (2014). CMS-PAS-HIG-14-001.
- [170] CMS COLLABORATION, “Search for Associated Production of a Single Top Quark and a Higgs Boson in Leptonic Channels,” (2015). CMS-PAS-HIG-14-026.
- [171] B. MAIER, *Search for the associated production of a single top quark and a Higgs boson decay channel in the $H \rightarrow b\bar{b}$ at 8 and 13 TeV with the CMS experiment*. (2016). PhD thesis, KIT, IEKP-KA/2016-8.
- [172] MADGRAPH, “MadGraph Wiki - Reweight.” <https://cp3.irmp.ucl.ac.be/projects/madgraph/wiki/Reweight>, last accessed on 27.04.2016.
- [173] CMS COLLABORATION, “PdmV2015Analysis,” (2015). https://twiki.cern.ch/twiki/bin/viewauth/CMS/PdmV2015Analysis#ReReco_at_25_ns, last accessed 11.05.2016.
- [174] CMS COLLABORATION, “Hadron based origin identification of heavy flavour jets at generator level.” <https://twiki.cern.ch/twiki/bin/view/CMSPublic/GenHFHadronMatcher#Introduction>, last accessed on 17.05.2016.

- [175] CMS COLLABORATION, “Baseline muon selections for Run-II - Tight Muon.” https://twiki.cern.ch/twiki/bin/view/CMS/SWGuideMuonIdRun2#Tight_Muon, last accessed on 27.04.2016.
- [176] CMS COLLABORATION, “Triggering electron MVA details and working points.” https://twiki.cern.ch/twiki/bin/view/CMS/MultivariateElectronIdentificationRun2#Triggering_electron_MVA_details, last accessed on 07.05.2016.
- [177] CMS COLLABORATION, “Jet Identification - Recommendations for 13 TeV data analysis.” https://twiki.cern.ch/twiki/bin/view/CMS/JetID#Recommendations_for_13_TeV_data, last accessed on 27.04.2016.
- [178] CMS COLLABORATION, “Jet Energy Resolution.” https://twiki.cern.ch/twiki/bin/viewauth/CMS/JetResolution#JER_Scaling_factors_and_Uncertai, last accessed on 27.04.2016.
- [179] MUONPOG, “Reference muon id, isolation and trigger efficiencies for Run-II,” (2016). <https://twiki.cern.ch/twiki/bin/view/CMS/MuonReferenceEffsRun2>, last accessed on 27.04.2016.
- [180] MUONPOG, “Muon T&P Instructions for Run-II,” (2016). <https://twiki.cern.ch/twiki/bin/view/CMS/MuonTagAndProbeTreesRun2>, last accessed on 27.04.2016.
- [181] EGAMMAPOG, “Instructions for applying electron and photon ID,” (2016). <https://twiki.cern.ch/twiki/bin/view/CMS/EgammaIDRecipesRun2>, last accessed on 27.04.2016.
- [182] EGAMMAPOG, “Instructions for applying electron and photon ID,” (2016). <https://twiki.cern.ch/twiki/bin/view/CMS/ElectronScaleFactorsRun2>, last accessed on 27.04.2016.
- [183] BTagPOG, “Event reweighting using scale factors calculated with a tag and probe method,” (2016). <https://twiki.cern.ch/twiki/bin/view/CMS/BTagShapeCalibration>, last accessed on 17.05.2016.
- [184] CMS COLLABORATION, “Calibration of the Combined Secondary Vertex b-Tagging discriminant using dileptonic $t\bar{t}$ and Drell-Yan events,” (2016). CMS Note 2013/313.
- [185] M. NICOLAS BULTE, EFE YAZGAN, “Towards a better description of extra jet emission in $t\bar{t}$ events.” <https://indico.cern.ch/event/517453/contributions/2031156/>, last accessed on 07.05.2016.
- [186] V. D. BARGER, J. OHNEMUS, AND R. J. N. PHILLIPS, “Event shape criteria for single lepton top signals,” *Physical Review D*, vol. 48, pp. 3953–3956, (1993).
- [187] G. C. FOX AND S. WOLFRAM, “Observables for the analysis of event shapes in e^+e^- annihilation and other processes,” *Physical Review Letters*, vol. 41, no. 23, p. 1581, (1978).

- [188] CMS COLLABORATION, “CMS Luminosity Measurement for the 2015 Data Taking Period,” (2016). CMS-PAS-LUM-15-001.
- [189] CMS COLLABORATION, “NNLO+NNLL top-quark-pair cross sections.” <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/TtbarNNLO>, last accessed on 27.04.2016.
- [190] CMS COLLABORATION, “NLO single-top channel cross sections.” <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/SingleTopRefXsec>, last accessed on 27.04.2016.
- [191] F. MALTONI, D. PAGANI, AND I. TSINIKOS, “Associated production of a top-quark pair with vector bosons at NLO in QCD: impact on $t\bar{t}H$ searches at the LHC,” *Journal of High Energy Physics*, vol. 02, p. 113, (2016).
- [192] CMS COLLABORATION, “Standard model cross sections for CMS at 13 TeV.” <https://twiki.cern.ch/twiki/bin/view/CMS/StandardModelCrossSectionsat13TeVInclusive>, last accessed on 27.04.2016.
- [193] CMS COLLABORATION, “CMS Phase II Upgrade Scope Document,” (2015). CERN-LHCC-2015-019.
- [194] LHC HIGGS CROSS SECTION WORKING GROUP, “Handbook of LHC Higgs Cross Sections: 3. Higgs Properties,” (2013). arXiv:1307.1347.
- [195] M. CZAKON, P. FIEDLER, AND A. MITOV, “The total top quark pair production cross-section at hadron colliders through $\mathcal{O}(\alpha_S^4)$,” *Physical Review Letters*, vol. 110, p. 252004, (2013).
- [196] N. KIDONAKIS, “Differential and total cross sections for top pair and single top production,” pp. 831–834, (2012). arXiv:1205.3453.
- [197] CMS COLLABORATION, “Standard Model Cross Sections for CMS at 8 TeV.” <https://twiki.cern.ch/twiki/bin/view/CMS/StandardModelCrossSectionsat8TeV>, last accessed on 31.05.2016.
- [198] CMS COLLABORATION, “Summary table of samples produced for the 1 Billion campaign, with 25ns bunch-crossing.” <https://twiki.cern.ch/twiki/bin/viewauth/CMS/SummaryTable1G25ns>, last accessed on 27.05.2016.

Danksagung

Eine Doktorarbeit ist nicht nur der Verdienst einer einzelnen Person. In dieser Danksagung möchte ich deshalb all denen danken, die diese Doktorarbeit auf die ein oder andere Weise beeinflusst haben.

Zuallererst danke ich Prof. Dr. Thomas Müller, der mir schon vor über vier Jahren ermöglicht hat, bereits zu Beginn meiner Diplomarbeit mich seiner Arbeitsgruppe anzuschliessen. In dieser für Teilchenphysiker sehr spannenden Zeit war es mir durch ihn möglich Teil einer einmaligen Arbeitsgruppe zu sein und dort uneingeschränkt nach meinen Vorlieben forschen zu können.

Prof. Dr. Günter Quast möchte ich recht herzlich für die Übernahme des Korreferats und die äusserst hilfreichen Kommentare zu dieser Arbeit zu danken.

Weiterhin danke ich natürlich Dr. Thorsten Chwalek, der in den letzten Jahren mein direkter Betreuer war und in jeder noch so vertrackten Situation einen Ausweg sah und immer mit Rat und Tat zur Stelle war. Weiterhin möchte ich ihm für das ausführliche Korrekturlesen meiner Arbeit danken. Ebenso danke ich hier auch meinen ehemaligen Betreuern während der Diplomarbeit und der Anfangszeit meiner Promotion: Dr. Jeannine Wagner-Kuhr und Dr. Hauke Held.

Ein riesen Dankeschön geht an meine zwei engsten Kollegen, die hier mit mir zusammen an der tH Analyse gearbeitet haben: Dr. Benedikt Maier und Dr. Christian Böser. Diese Zusammenarbeit ging weit über eine einfache kollegiale Beziehung hinaus. Mit euch konnte man nicht nur richtig gut zusammen arbeiten, sondern auch besonders viel lachen.

Weiterhin geht ein Dank an alle, die in den letzten Jahren wichtigen Input zur Analyse beigetragen haben und in den letzten Monaten mitgeholfen haben die Analyse bis zur Veröffentlichung zu bringen. Vielen Dank an Nils Faltermann, Kevin Flöh, Denise Müller, Dr. Frank Roscher und Matthias Schnepf.

Vielen Dank geht auch an Dr. Steffen Röcker, der als Bürokollege vier Jahre lang zu einer einmaligen Büroatmosphäre beigetragen hat. Vielen Dank auch an Dr. Matthias Mozer, der mir oft bei Physik-Fragen weiterhelfen konnte. Auch weiteren Kollegen am EKP sei hier gedankt, herauszuheben wäre hier Michael Ziegler, der bei physikalischen, sowie nicht-physikalischen Fragen immer eine Antwort parat hatte.

Special thanks is owed to the people, who collaborated with me on the tH analysis within CMS over the years. Prof. Dr. Andrea Giammanco, Prof. Dr. Ken Bloom, Dr. Andrey Popov and Daniel Knowlton contributed a lot in shaping the tH analysis during its early days during Run-I, and

B. Danksagung

put in enormous effort culminating in the successful publication of the tH combination paper.

Vielen Dank auch an Frau Bräunling, die als Sekretärin des Instituts viel zur angenehmen Arbeitsatmosphäre beiträgt. Ebenso einen grossen Dank an alle Admins innerhalb des EKPs, die Ihre wertvolle Zeit für eine gut funktionierende IT-Infrastruktur aufwenden.

Ein ganz besonderer Dank geht natürlich auch an meine gesamte Familie: Ihr habt mich jahrelang unterstützt und mir so viel Vertrauen entgegen gebracht, dass ein einzelner Abschnitt hier niemals meine Dankbarkeit euch gegenüber genügend ausdrücken könnte. Danke.

Zuletzt möchte ich meiner Freundin Carolin danken, die mir immer ein immenser Rückhalt ist. Du hast mir vor allem in den letzten Monaten so viel Kraft gegeben und es mir so einfach wie möglich gemacht diese Arbeit zu verfassen. Ohne dich wäre diese Arbeit so nicht möglich gewesen. Ich bin der glücklichste Mensch, dass ich dich gefunden habe. Dankeschön.

Hiermit versichere ich, die vorliegende Arbeit selbstständig verfasst
und nur die angegebenen Hilfsmittel verwendet zu haben.

Simon Fink

Karlsruhe, den 22. Juni 2016