



**PROSPECTS FOR DI-HIGGS BOSON SEARCHES WITH
MULTILEPTON FINAL STATES IN THE $q\bar{q}HH$ PRODUCTION MODE**

Master's Thesis

Student: Nalong-Norman Seeba

Supervisor: Torben Lange, NICPB Research Fellow

Co-Supervisor: Christian Veelken, NICPB Senior Research Fellow

Study program: YAFM

Tallinn 2023

Declaration

Hereby I declare that I have compiled the paper independently and all works, important standpoints and data by other authors have been properly referenced and the same paper has not been previously been presented for grading.

Author: Nalong-Norman Seeba

Signed digitally, 04.01.2023

The paper conforms to requirements in force.

Supervisor: Torben Lange

Co-supervisor: Christian Veelken

Signed digitally, 04.01.2023

Permitted to the defence.

Chairman of the Defence Committee: Jüri Elken

Signed digitally

Abstract

CMS recently presented its first search for Higgs boson pair (HH) production in the $WWWW$, $WW_{\tau\tau}$, and $\tau\tau\tau\tau$ decay modes using the full LHC Run2 dataset with 138 fb^{-1} of proton-proton collision data recorded by the CMS experiment at a center of mass energy of $\sqrt{s} = 13 \text{ TeV}$. The analyzed events contain two, three, or four reconstructed leptons, including electrons, muons, and hadronically decaying tau leptons. While the original analysis mostly focuses on constraints on the trilinear Higgs self-coupling λ and the HH signal component in the $ggHH$ production mode, the presented work focuses on the feasibility and prospects of extending the analysis to the sub-dominant $qqHH$ production mode. While at a factor of about 5-10 lower cross section than $ggHH$, the study of $qqHH$ allows to study the so far unmeasured $HHVV$ Higgs boson (H) coupling to Vector bosons (V). The presented results for a final state with two same charge leptons (electrons or muons) and one hadronically decaying tau are quite promising, motivating a more dedicated study of $qqHH$ in future analysis iterations. The projected sensitivity is now in the range of other HH analyses from CMS and ATLAS.

List of abbreviations and terms

ALICE	A Large Ion Collider Experiment
APD	Avalanche photodiodes
AR	Application region
ATLAS	A Toroidal LHC ApparatuS
BDT	Boosted decision tree
BM	Benchmark
BPIX	Barrel pixel detector
BSM	Beyond Standard Model
CERN	European Organization for Nuclear Research
CL	Confidence level
CMS	Compact Muon Solenoid
CP	Charge conjugation parity
CSC	Cathode strip chambers
DNN	Deep neural network
DT	Drift tube chambers
EB	ECAL barrel
ECAL	Electromagnetic calorimeter
EE	ECAL endcap
EFT	Effective field theory
FF	Fake factor
FPIX	Forward pixel detector
GCT	Global calorimeter trigger
ggF	Gluon-gluon fusion
GMT	Global muon trigger
GSF	Gaussian sum filter
GT	Global trigger
HB	Hadron barrel calorimeter
HCAL	Hadron calorimeter
HE	Hadron endcap calorimeter
HEP	High energy physics
HF	Hadron forward calorimeter
HH	Higgs boson pair
HLT	High level trigger
HO	Hadron outer calorimeter
HPS	Hadrons plus strips
IP	Interaction point
L1	Level 1 trigger

LEP	Large Electron–Positron Collider
LHC	Large Hadron Collider
LHCb	LHC beauty
LO	Leading-order
MB	Muon barrel
MC	Monte Carlo simulation
ME	Muon endcap
MET	Missing transverse energy
MIP	Minimum ionizing particle
ML	Machine learning
MR	Measurement region
MVA	Multivariate analysis
NLO	Next-to-LO
NNLO	Next-to-next-to-LO
N3LO	Next-to-NNLO
PF	Particle flow
POG	Physics object group
PSB	Proton Synchrotron Booster
QED	Quantum electrodynamics
QFT	Quantum field theory
QCD	Quantum chromodynamics
RCT	Regional calorimeter trigger
ROC	Receiver operating characteristic
ROCs	Readout chips
RPC	Resistive plate chambers
SC	Supercluster
SFOS	Same flavour opposite sign charge
SL	Superlayer
SM	Standard Model
SPS	Super Proton Synchrotron
SR	Signal region
TEC	Tracker endcap
TIB	Tracker inner barrel
TID	Tracker inner disc
TOB	Tracker outer barrel
VBF	Vector boson fusion
VPT	Vacuum photodiodes
WLS	Wavelength shifting
WP	Working point

Table of Contents

1. Introduction	1
2. The Standard Model of particle physics and the Higgs boson	2
2.1. The Standard Model of particle physics	2
2.2. The Higgs boson	4
2.3. Di-Higgs production and Higgs boson interactions	5
2.3.1. Trilinear self-coupling and couplings to vector bosons	6
2.3.2. Gluon-gluon fusion like ggHH production	6
2.3.3. Vector boson fusion like qqHH production	7
3. CMS experiment at the LHC	8
3.1. The Large Hadron Collider	8
3.2. The Compact Muon Solenoid	10
3.2.1. Coordinate system	10
3.2.2. Inner tracking system	11
3.2.3. Calorimeters	13
3.2.4. Muon system	16
3.2.5. Triggers	19
3.2.6. The Particle-Flow algorithm	20
4. Machine learning methods	22
4.1. Boosted decision trees	23
5. The qqHH analysis in multilepton final states	24
5.1. CMS HH \rightarrow multilepton analysis	24
5.2. Analysis overview	26
5.3. Datasets and Monte Carlo simulation	27
5.4. Object reconstruction	27
5.4.1. Electrons	28
5.4.2. Muons	29
5.4.3. Jets	30
5.4.4. Hadronic τ decays	30
5.4.5. Reconstruction of event level quantities	32
5.5. Data driven background estimation	32
5.6. Event selection	33
5.7. BDT discriminants	34

5.7.1. BDT input variables and performance	35
5.7.2. 2D histograms	40
6. Results	44
6.1. Systematic uncertainties	44
6.2. Signal extraction	45
6.3. Results on qqHH production	47
6.4. Comparison to other analyses	50
7. Summary and conclusions	52
8. Acknowledgements	53
Bibliography	54

1. Introduction

Particle physics aims at unraveling the mysteries of the fundamental laws of nature through studying elementary particles and the forces between them. The combined effort of our understanding of particle physics is contained in the theory known as the Standard Model. The final discovered particle predicted by the Standard Model is the Higgs boson, which was found by the CMS and ATLAS collaborations [1, 2] in 2012. The Higgs boson is the first known elementary scalar particle and a key part of the Higgs mechanism, necessary for a self consistent formulation of the Standard Model. Physics beyond Standard Model (BSM) is expected, since the SM is not powerful enough to explain everything. Hence it is of great interest to study the properties of the Higgs boson. Many of these have already been measured with great precision, however a few important ones still remain uncertain, such as the Higgs boson self-coupling (λ) and the quartic Higgs boson coupling to vector bosons (HHVV). CMS has recently published results on Higgs boson pair production with multilepton final states [3] focusing on the trilinear Higgs boson self-coupling and the signal component in the dominant ggHH production mode. This thesis focuses on the feasibility and prospects of extending this CMS HH \rightarrow multilepton analysis [3] to the sub-dominant qqHH production mode, giving insight about the Higgs boson HHVV coupling. The focus towards qqHH is inspired by the work done in reference [4] proposing a strategy to extract the HHVV quartic coupling via focusing on qqHH production in $b\bar{b}b\bar{b}$ final state. In this thesis the $2lss + 0/1\tau_h$ channel from the CMS HH \rightarrow multilepton analysis is chosen for the extension.

This thesis is structured as follows. Chapter 2 introduces the Standard Model and its structure, together with the Higgs boson, its self-interaction and interactions to vector bosons, and di-Higgs production processes. The overview of the experimental setup is given in chapter 3. Chapter 4 describes the machine learning methods used in this analysis and chapter 5 gives a brief introduction to the CMS HH \rightarrow multilepton analysis, and contains the qqHH focused analysis in multilepton final states. The results of the analysis are presented in chapter 6 and the thesis is concluded with a summary in chapter 7.

2. The Standard Model of particle physics and the Higgs boson

The Standard Model of particle physics (SM) is the quantum field theory (QFT) that describes the combined knowledge of the known fundamental forces and particles in the universe, excluding dark matter, dark energy and gravity [5]. It has been confirmed by numerous dedicated experiments probing various physical properties magnificently described by the SM. The Higgs boson was the last elementary particle of the SM yet to be observed, and was finally discovered in 2012 at CERN (European Organization for Nuclear Research) by the CMS and ATLAS collaborations [1, 2]. The discovery of the Higgs boson supports the existence of the Higgs mechanism providing a self consistent way for fundamental particles to gain mass. [6] The sections that follow, aim to give a brief description of the SM and the Higgs boson. Section 2.1 describes the structure and particle content of the SM, section 2.2 introduces the Higgs boson and gives a brief overview of the Higgs mechanism and the Higgs potential, and section 2.3 discusses di-Higgs production, the Higgs boson self interaction and coupling to vector bosons, and the leading and sub-leading (di-) Higgs production mechanisms.

2.1. The Standard Model of particle physics

The SM describes matter as a small number of different fundamental half-integer spin (1/2) particles called fermions, which can be divided into 6 quarks and 6 leptons. The SM explains the interactions of these fermions by the exchange of integer spin bosons. Three types of fundamental interactions can be explained by the SM: the strong interaction, weak interaction, and the electromagnetic interaction. [7] The up-quark, down-quark, together with the electron and the electron neutrino make up the first generation of fermions, and represent the building blocks of the low-energy universe. At higher energy scales, each of the four first generation particles has exactly two copies, which vary only in mass. These extra particles are known as the second and third generations of fermions. With the exception of the differences in mass, the properties of the corresponding particles in different generations are the same, as the particles possess exactly the same fundamental interactions as their first generation counterparts. [8] The three generations of fermions can be seen in Table 1.

Table 1. The three generations of fermions

Fermions				
Generation		I	II	III
Quarks	up-type	up (u)	charm (c)	top (t)
	down-type	down (d)	strange (s)	bottom (b)
Leptons	charged	electron (e^-)	muon (μ^-)	tau (τ^-)
	neutrinos	ν_e	ν_μ	ν_τ

The properties of the twelve fundamental fermions are classified by the types of interaction

they experience. Most fermions experience the weak force and undergo weak interactions. [8] Weak interactions are for example known by the slow process of nuclear β -decay, involving the emission of an electron and neutrino by a radioactive nucleus, and nuclear fusion. The mediators of weak interactions are the massive spin-1 bosons: the W^\pm and Z^0 . [7] Omitting the neutrinos, which are electrically neutral, the other fermions are electrically charged and participate in the electromagnetic interaction of quantum electrodynamics (QED), which is a QFT that explains the electromagnetic force, where the interactions between charged particles are mediated by photon exchange. Quarks are the only fermions that carry color charge, which is the quantum chromodynamics (QCD) equivalent of the electric charge, and participate in strong interactions. [8] Strong interactions are responsible for binding the quarks in the neutron and proton, and the neutrons and protons within atomic nuclei. The interquark force is mediated by a massless boson, named the gluon. [7] Quarks are always confined to color neutral bound states called hadrons, such as the proton and neutron, and never observed as free particles. [8] The mediators of the three fundamental forces are all spin-1 particles known as gauge bosons and can be seen in Table 2 together with their corresponding forces.

Table 2. Gauge bosons

Boson	Force
Gluon (g)	Strong
Photon (γ)	Electromagnetism
W and Z bosons (W^\pm, Z^0)	Weak

The final element of the SM is the Higgs boson, which is a spin-0 scalar boson discovered by the ATLAS and CMS experiments at the Large Hadron Collider in 2012 [1, 2]. The Higgs boson is discussed in the section that follows. The complete structure of the SM together with the particle content can be seen in Figure 1

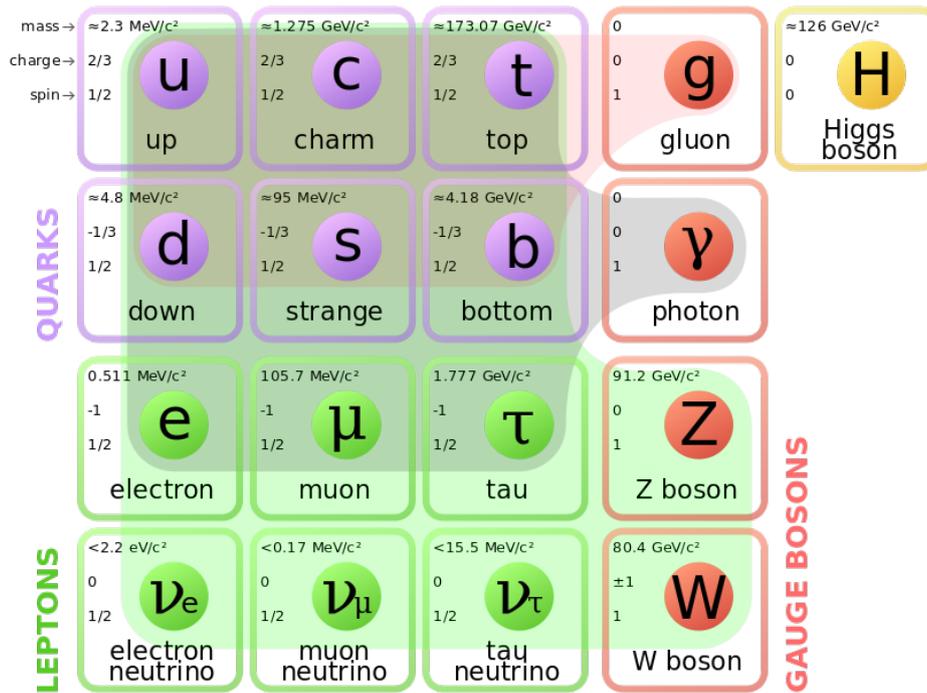


Figure 1. The particle content of the Standard Model. [9]

2.2. The Higgs boson

The Higgs boson, with a mass of about 125 GeV differs from all other SM particles, being a neutrally charged spin-0 scalar particle and the only fundamental scalar discovered to date. The Higgs boson has a particular function in the SM, as it provides the mechanism by which other particles can acquire mass through interacting with the Higgs field, which has a non-zero vacuum expectation value, as opposed to the fields of the other fundamental particles. The higgs mechanism breaks the $SU(2) \times U(1)$ local gauge symmetry of the electroweak sector of the SM. [8]

The SM Higgs potential is given by:

$$V = \frac{1}{2}m_H^2\Phi_H^2 + \lambda v\Phi_H^3 + \frac{\tilde{\lambda}}{4}\Phi_H^4, \quad (2.1)$$

where v is vacuum expectation value of the Higgs field, with $v = 246 \text{ GeV}$, m_H is the mass of the Higgs boson, and the cubic and quartic terms are proportional to the Higgs boson self-couplings, with λ being the trilinear coupling constant, and $\tilde{\lambda}$ denoting quartic Higgs self-coupling. Measuring the Higgs boson self couplings therefore allows the probing of the higher order terms in the Higgs potential, thus providing valuable information concerning the mechanism of electroweak symmetry breaking [10] and the shape of the Higgs potential. [3] An illustration of the shape of the Higgs potential can be seen in Figure 2. Particles couple to the Higgs boson proportional to their mass - the stronger the interaction to the Higgs, the bigger the mass of the particle.

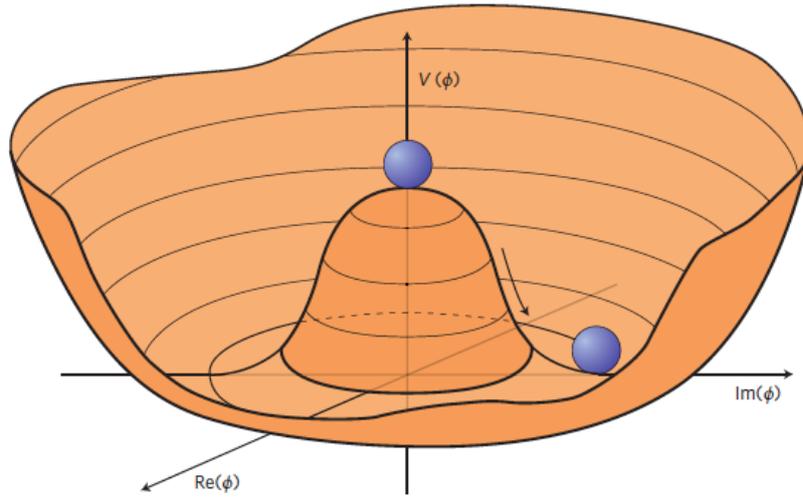


Figure 2. An illustration of the Higgs potential [11]

2.3. Di-Higgs production and Higgs boson interactions

The search for Higgs boson pair (HH) production can be used to study the shape of the Higgs potential and verifying the electroweak gauge symmetry breaking mechanism. It is done through probing Higgs boson couplings, such as the trilinear self-coupling, which is predicted by the SM. [3]

The main single Higgs boson production modes at the Large Hadron Collider (LHC) are gluon-gluon fusion (ggF or ggH), vector boson fusion (VBF or qqH), and associated production with top quarks ($t\bar{t}H$) or with vector bosons (WH, ZH). The dominant Higgs production mechanism is ggF and the sub-dominant one is VBF [12]. Analogously to single Higgs boson production, di-Higgs production can be divided into multiple production modes, such as ggF-like (ggHH) or VBF-like (qqHH). In addition to this, it is also possible for HH production to occur in associated production with top quark pairs ($t\bar{t}HH$) or with vector bosons (WHH, ZHH), but the cross sections of these processes are significantly lower than the ones of ggHH and qqHH. [13] The HH production cross sections as a function of λ/λ_{SM} can be seen in Figure 3. Sections 2.3.2 and 2.3.3 discuss ggHH and qqHH more in detail.

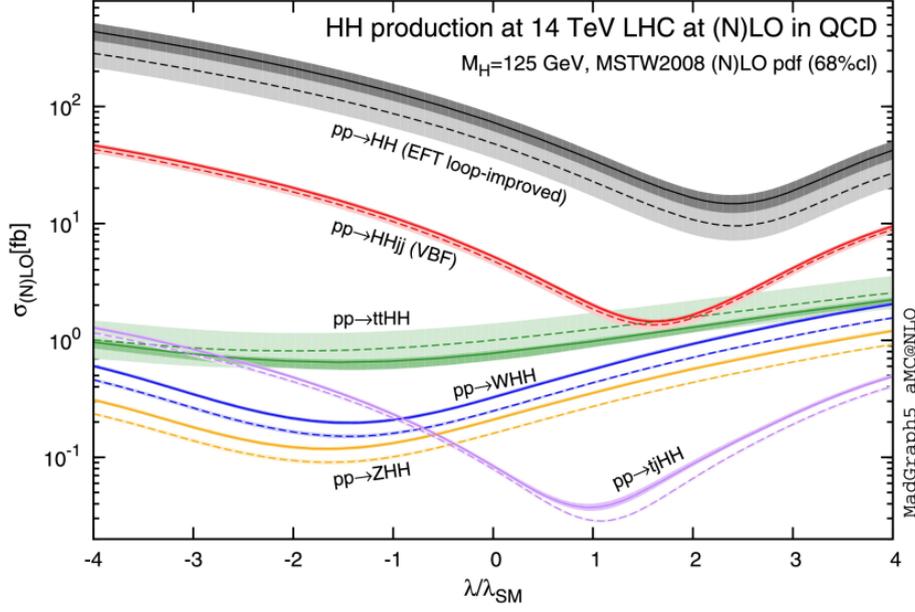


Figure 3. Cross sections for HH production channels, at the $\sqrt{s} = 14$ TeV LHC as a function of λ/λ_{SM} . Here ggHH is drawn in grey, qqHH in red, and the associated production modes in green, blue, yellow, and violet. [14]

2.3.1. Trilinear self-coupling and couplings to vector bosons

The SM trilinear self-coupling of the Higgs boson is defined by the coupling strength $\lambda = \frac{m_H^2}{2v}$, where m_H is the mass of the Higgs boson and v is the electroweak symmetry breaking vacuum expectation value [15]. The ratios of couplings to their SM predictions are referred to as coupling strength modifiers and are defined at 1 in the SM. The coupling strength modifier for the trilinear self-coupling λ is denoted as κ_λ . [3]

In addition to probing the trilinear self-coupling of the Higgs boson, HH production can also be used to probe the strength of the Higgs couplings to vector bosons at high energies in the VBF production mode. There exists a quartic coupling between two Higgs bosons and two vector bosons denoted as HHVV, and a coupling between one Higgs boson and two vector bosons denoted as HVV. The coupling strength modifiers for HHVV and HVV are denoted as C_{2V} and C_V respectively, with the SM value 1. [4]

2.3.2. Gluon-gluon fusion like ggHH production

Gluon-gluon fusion is the dominant Higgs production mechanism at the LHC and is a process that involves two gluons fusing via a loop of a heavy colored particle. The loop facilitates an indirect coupling, despite there being no direct coupling between the Higgs and gluons in the SM. The biggest contribution in the loop comes from the top-quark, as it has the highest mass and thus the biggest coupling constant to the Higgs, and is supplemented by a smaller contribution of bottom-quark loops. The ggHH process is governed by two parameters: the Higgs boson self-coupling (λ) and the Yukawa coupling of the top-quark (Y_t), and can be

seen in Figure 4. There are two destructively interfering diagrams known as the triangle diagram, which is sensitive to the Higgs boson self-coupling, and the box diagram, depending only on the top Yukawa coupling. The destructive interference of these diagrams results in a low total cross section for HH production. [3, 13] The SM ggHH cross section at NNLO (next-to-next-to-leading-order) accuracy in QCD at 13 TeV has been calculated to be $31.1^{+2.1}_{-7.2}$ fb [16].

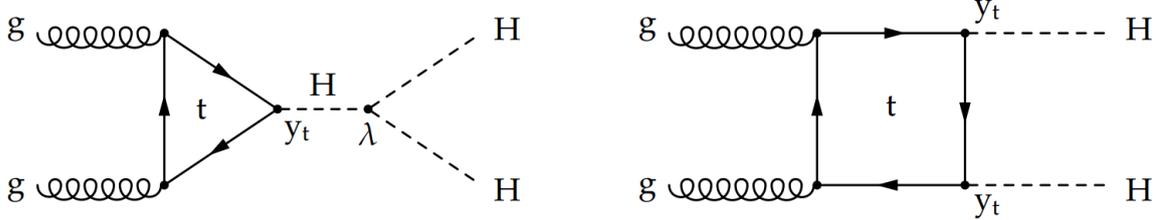


Figure 4. Leading order Feynman diagrams contributing to ggHH production. The left diagram is referred to as the triangle diagram and is sensitive to the Higgs boson trilinear coupling λ and the Yukawa coupling of the top-quark Y_t . The right diagram is the box diagram, which is insensitive to λ , and scales as Y_t^2 . [3]

2.3.3. Vector boson fusion like qqHH production

Vector boson fusion is the second largest Higgs production mechanism at the LHC and is a process where two initial-state quarks both radiate off a virtual electroweak gauge boson (namely a W or Z boson), that fuse to produce a particle, such as the Higgs at the LHC. The initial-state quarks are detected as high energy jets by the detector. [13]

VBF induced di-Higgs production can be seen in Figure 5. Starting from left, there are three destructively interfering leading-order diagrams contributing to qqHH production, depending on the HHVV, HVV, and HVV in combination with λ couplings respectively. The coupling strength modifiers are C_{2V} for HHVV and C_V for HVV. [4] The SM cross section for qqHH is about an order of magnitude smaller than the cross section for qqHH, and has been computed at N³LO (next-to-NNLO) accuracy in QCD at 14 TeV to be $2.055^{+0.001}_{-0.001}$ fb [17].

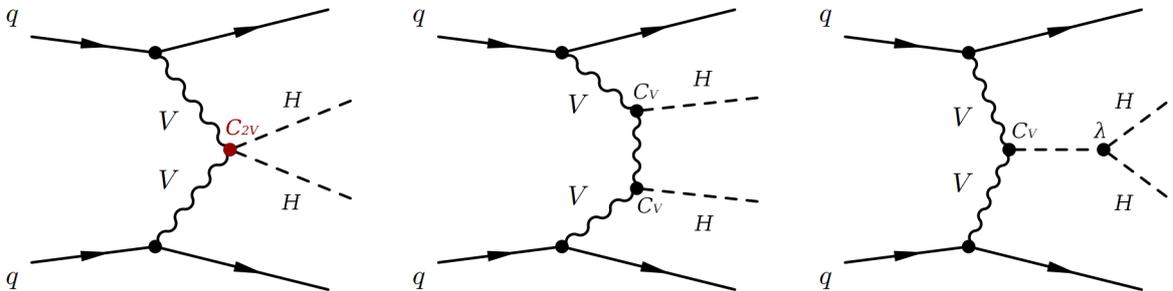


Figure 5. Leading-order Feynman diagrams contributing to qqHH production. The left, center, and right diagrams scale as C_{2V} , C_v^2 , and $C_v\lambda$ respectively. [4]

3. CMS experiment at the LHC

This chapter gives a brief overview of the world's largest particle collider, named the Large Hadron Collider (LHC) at CERN (European Organization for Nuclear Research), and one of its four particle detector experiments, named the Compact Muon Solenoid (CMS).

3.1. The Large Hadron Collider

The Large Hadron Collider is a hadron accelerator and collider that was built at CERN between 1998-2008 and installed in the existing tunnel of the previously used LEP (Large Electron-Positron Collider) near Geneva, Switzerland. The accelerator ring has a circumference of 27 km and is located about 100 m below ground level. The aim of the LHC and its experiments is to discover new physics beyond the SM with proton-proton and heavy ion collisions at centre of mass energies of up to 14 TeV. [18]

The LHC has four main detector experiments (CMS, ATLAS, ALICE and LHCb) located at four different interaction points. There are two general-purpose detectors - CMS and ATLAS (A Toroidal LHC Apparatus) [19], which both aim to study different particle physics interactions, with the main goal of discovering and studying the Higgs boson, and clarifying the nature of electroweak symmetry breaking [20]. The remaining two detectors are used for more specialized research. The first one is ALICE (A Large Ion Collider Experiment), studying heavy ion collisions, focuses on the strong-interaction sector of the SM. It is designed to study the physics of strongly interacting matter and quark-gluon plasma.[21] LHCb (LHC beauty), is dedicated to flavour physics, precision measurements of CP (charge conjugation parity) violation and rare decays of B hadrons [22].

The LHC itself belongs to a larger accelerator complex as seen in Figure 6. A pre-accelerator chain is used to supply the LHC with pre-accelerated protons. The chain consists of the Linac2 linear accelerator, Proton Synchrotron Booster (PSB), Proton Synchrotron and the Super Proton Synchrotron (SPS). These accelerators bring the proton beam energies up to 450 GeV before they reach the LHC, where the beams can get accelerated to 7 TeV. [18]

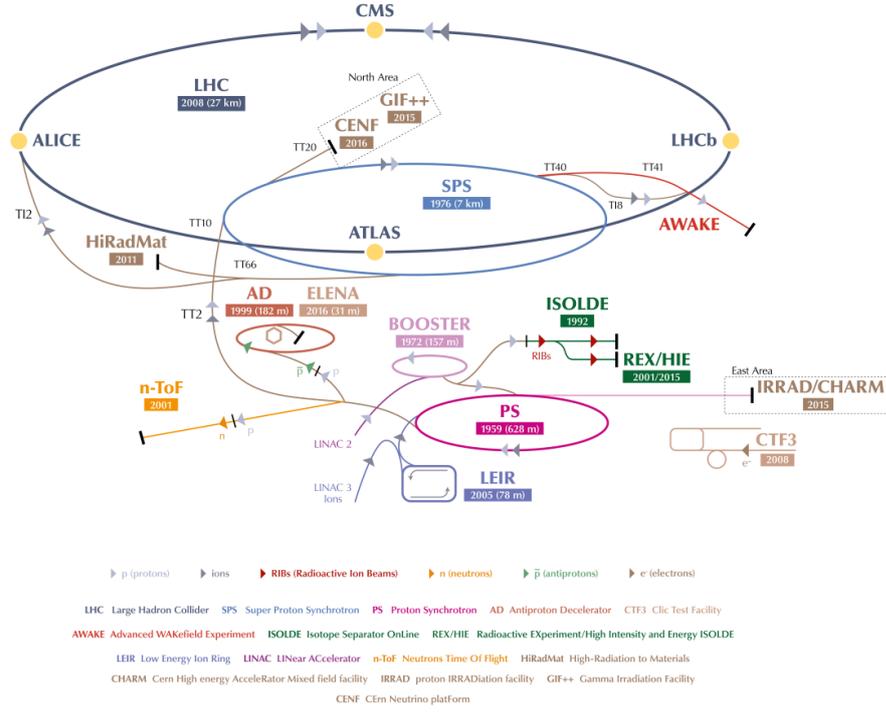


Figure 6. The CERN accelerator complex. The LHC is depicted as the dark blue line in a complex chain of particle accelerators, shown with its four main experiments. The smaller machines (pre-accelerators) are used in a chain to help boost the particles to their final energies before they get fed into the LHC. They are also used for smaller experiments.[23]

The two key parameters for a collider are the collision energy, and the luminosity L which is a measure of the number of collisions. The number of events generated in the LHC collisions is given by:

$$N = \sigma \times L \quad (3.1)$$

where σ is the production cross section (the probability that two particles will collide and react a certain way) for the event we are interested in and L is the integrated luminosity of the machine, defined as:

$$L = \int_{t_0}^t \mathcal{L} dt \quad (3.2)$$

[24] The machine luminosity depends only on beam parameters and can be written for a Gaussian beam distribution as:

$$\mathcal{L} = \frac{N_b^2 n_b f_{rev} \gamma_r}{4\pi \epsilon_n \beta^*} F, \quad (3.3)$$

where N_b is the number of particles per bunch, n_b is the number of bunches per beam, f_{rev} is the revolution frequency, γ_r is the relativistic Lorentz factor of the proton beams, ϵ_n the normalized transverse beam emittance, β^* the beta function at the collision point, and F the geometric luminosity reduction factor due to the crossing angle at the interaction point (IP), given by:

$$F = \left(1 + \left(\frac{\theta_c \sigma_z}{2\sigma^*} \right)^2 \right)^{-1}, \quad (3.4)$$

where θ_c is the full crossing angle at the IP, σ_z is the RMS bunch length, and σ^* the transverse RMS beam size at the IP. [18] Integrated luminosity presents a measurement of the size of collected data, and is usually given in inverse femtobarns (fb^{-1}) - the units of inverse cross section. The LHC design parameters are given in Table 3 [25].

Table 3. LHC design parameters

\sqrt{s} (TeV)	N_b	n_b	β^* (m)	ϵ_n (mm μ rad)	f_{rev} (kHz)	θ_c (μ rad)	σ_z (cm)	σ^* (μ m)
14	2808	1.15×10^{11}	0.55	3.75	11.245	285	7.55	16.7

3.2. The Compact Muon Solenoid

The Compact Muon Solenoid is located at the LHC interaction point 5, focusing on a wide range of Standard Model and beyond Standard Model physics, with one of its main goals being the discovery and study of the Higgs boson. The 13 m long, 4 T superconducting solenoid with an inner diameter of 6 m at the centre of CMS is the main feature of the detector, giving CMS its name. The solenoid is used to bend the trajectory of high-energy charged particles in order to precisely measure their momentum and identify the charge of the particle. The detector components are placed cylindrically in layers around the interaction point. The innermost main component is the silicon tracker, followed by calorimeters, then the superconducting solenoid magnet and ending with muon detectors. All of the components are confined within a steel return yoke of the magnet. The CMS detector is 21.6-m long and has a diameter of 14.6 m. It has a total weight of 12500 t. The overall layout of CMS is shown in Figure 7. [20]

3.2.1. Coordinate system

CMS uses a right-handed Cartesian coordinate system with the point of origin centered at the nominal collision point inside the experiment, the x -axis pointing radially inward towards the centre of the LHC ring, the y -axis pointing vertically upward, and the z -axis pointing along the beam direction. Since the detector is radially symmetric around the z -axis we can use a cylindrical coordinate system to describe the direction of particles travelling through the detector. We define ϕ as the azimuthal angle which is measured from the x -axis in the $x - y$ plane, the radial coordinate in this plane is denoted by r , and θ as the polar angle measured from the z -axis. For the sake of convenience pseudorapidity η , which is defined as

$$\eta = -\ln \tan\left(\frac{\theta}{2}\right) \quad (3.5)$$

is used instead of the polar angle θ , since differences in pseudorapidity are Lorentz invariant. Momentum and energy transverse to the beam line are denoted by p_T and E_T and are computed from the x and y components. [20] Angular separation between particles is denoted as ΔR and is defined as

$$\Delta R = \sqrt{(\Delta\eta)^2 + (\Delta\phi)^2} \quad (3.6)$$

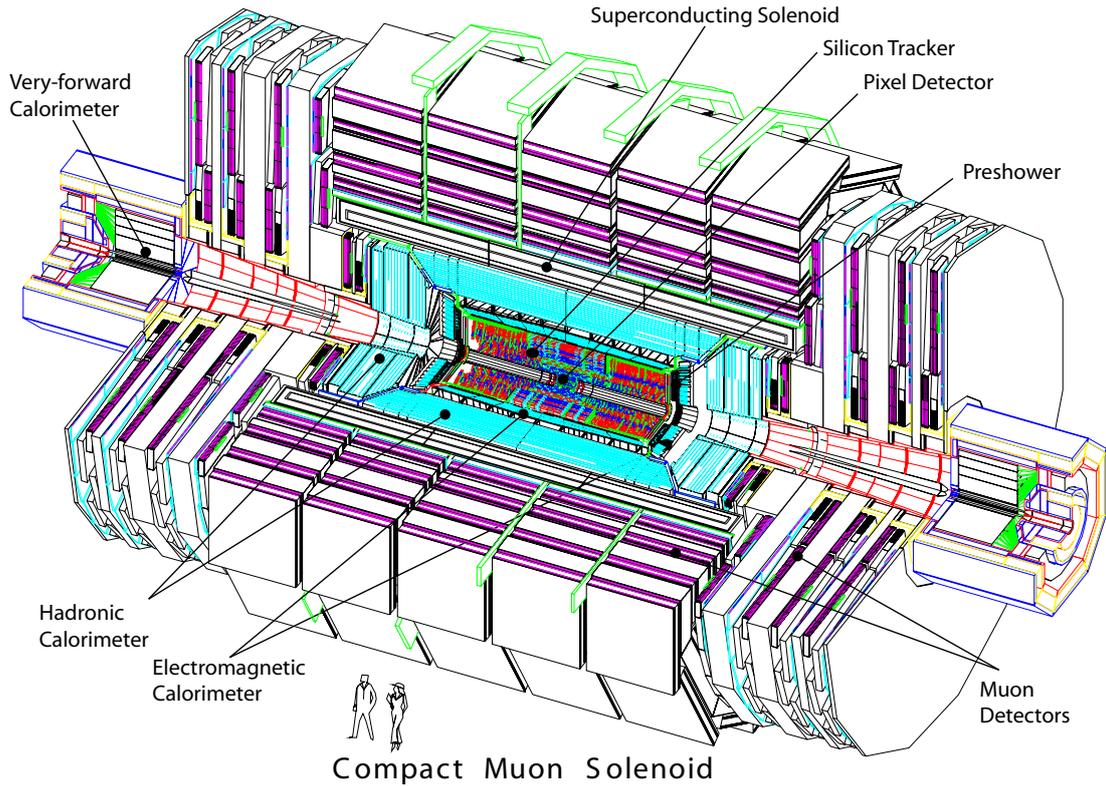


Figure 7. Schematic layout of the CMS detector. [26]

where $\Delta\eta$ is the difference in pseudorapidity and $\Delta\phi$ the difference in azimuthal angle of particles. The CMS coordinate system together with a visualisation of pseudorapidity η can be seen in Figure 8.

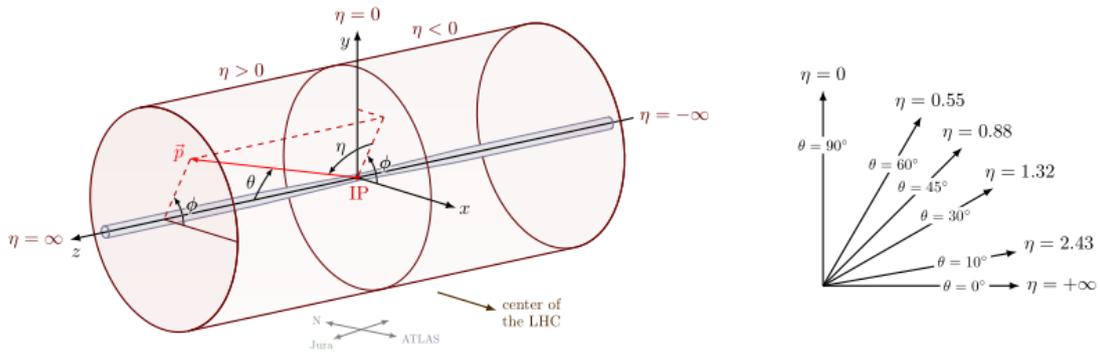


Figure 8. CMS coordinate system and pseudorapidity η . Here IP is the interaction point, which is the place where particles collide, \vec{p} describes the path a particle is travelling after the collision, ϕ is the azimuthal angle, and θ is the polar angle. [27, 28]

3.2.2. Inner tracking system

The purpose of CMS's inner tracking system is to provide a precise and efficient measurement of the trajectories of charged particles that emerge from the collisions. The tracking system surrounds the interaction point, and is 5.8 m long with a diameter of 2.5 m. It consists of a

pixel detector and a silicon strip tracker, and can be seen in Figure 9. [20]

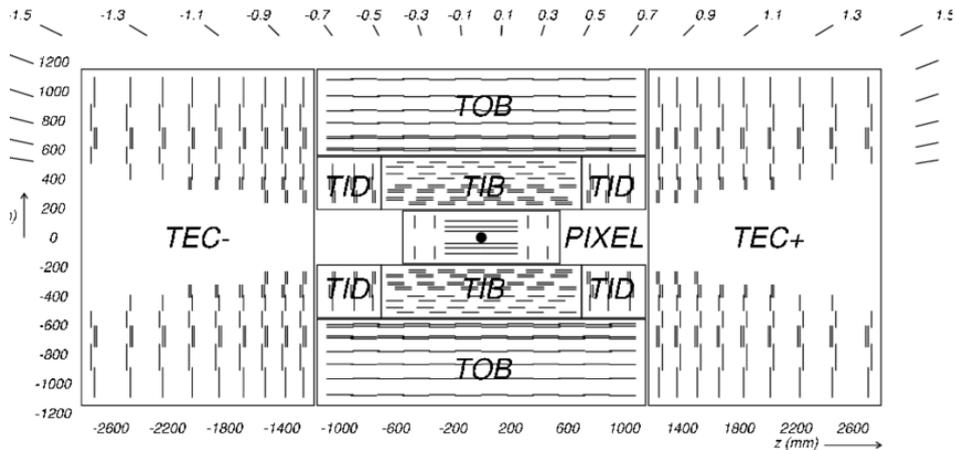


Figure 9. Schematic cross section through the CMS tracker. The silicon pixel detector in the center is surrounded by the silicon strip tracker parts: Tracker Inner Barrel (TIB) and Tracker Outer Barrel (TOB) detectors in the barrel region, Tracker Inner Disk (TID) and the Tracker End Caps (TEC) in the endcap regions. [20]

The pixel detector is in the region closest to the interaction point, and provides 3-dimensional space points, that enable high precision, charged particle tracking and vertex reconstruction. As it is the innermost detector, it is located in a challenging environment with especially strong radiation, characterized by a flux of approximately 1000 charged particles traversing the inner tracker every 25 ns. It consists of a cylindrical barrel (BPIX) detector, covering a pseudorapidity range of $|\eta| < 1.3$ and a forward-facing (FPix) detector, covering $|\eta| < 2.5$. At the end of 2016, the phase-0 CMS pixel detector was replaced with an upgraded version, the CMS Phase-1 pixel detector. The original pixel detector consisted of three barrel layers (BPIX) at radii 4.4 cm, 7.3 cm and 10.2 cm and two endcap disks (FPix) on each end at distances of 34.5 cm and 46.5 cm from the interaction point. The 3-layer barrel, 2-disk endcap system was replaced with a 4-layer barrel, 3-disk endcap system, and has layers at radii 2.9, 6.8, 10.9, and 16.0 cm, and three disks on each end at distances of 29.1, 39.6, and 51.6 cm from the center of the detector. The layout of the CMS Phase-1 pixel detector compared to the original pixel detector can be seen in Figure 10. [29]

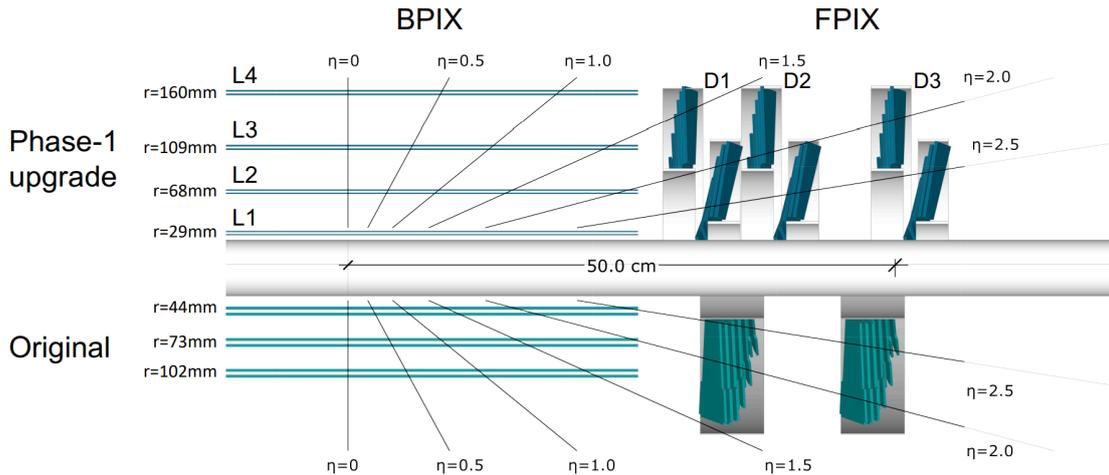


Figure 10. Comparison of the layout of the CMS Phase-1 pixel detector to the original detector. [29]

In total the CMS Phase-1 pixel detector uses 1856 segmented silicon sensor modules: 1184 modules in the barrel pixel detector (BPIX) and 672 modules in the forward disks (FPIX). Each module consists of a sensor with 160×416 pixels connected to 16 readout chips (ROCs), which are divided into 4160 readout channels and read out the analog pulse height information for each pixel, when they come into contact with charged particles. The analog signal is converted to light and transmitted to the readout electronics via optical fibers. The standard pixel size is $100 \times 150 \mu\text{m}^2$. [29, 30]

The silicon strip tracker, which surrounds the pixel detector, is made out of 10 detection layers in the barrel region and 12 disks in the endcap regions. Four of the 10 detection layers in the barrel region make up the Tracker Inner Barrel (TIB), which is accompanied by three discs on both sides, making up the Tracker Inner Disc (TID). Both the TIB and TID are surrounded by the remaining six detection layers, which make up the Tracker Outer Barrel (TOB), and the remaining nine disks on each end of the endcap regions make up the Tracker End Cap (TEC). The silicon strip detector is installed at a radial distance between 20 and 116 cm from the beam line. Similar to the pixel detector, the pseudorapidity range of $|\eta| < 2.5$ is covered. [20, 30] The different sub detectors making up the silicon strip tracker can be seen in Figure 9. In total 15158 modules with 24244 sensors cover about 9.3 million strips, making up an active area of 198 m^2 [31]. In order to reduce radiation damage, both the pixel and strip detectors are cooled to an operating temperature of -10°C .

3.2.3. Calorimeters

The energy of charged particles is determined from their momentum, measured using the tracking system. Calorimeters also allow to find neutral particles, measure their energy and distinguish them. Together with the tracking system, CMS uses two kinds of calorimeters measuring the energy of from charged particles by absorbing them completely: sampling calorimeters consist of a passive absorber material and an active detector material. When particles interact with the passive material either through electromagnetic or strong processes,

particle showers are induced. The energy deposited by the particles of the shower in the active part of the calorimeter can be detected in the form of scintillation light, and serves as a measurement of the energy of the incident particle [32]. In a homogeneous calorimeter, the active detector material also serves as the shower inducing material. The energy of electrons and photons is measured by the Electromagnetic Calorimeter (ECAL), which is the inner of the two calorimeters. Hadrons, which are composite particles made out of quarks and gluons, pass through the ECAL and are absorbed by the outer layer named the Hadron Calorimeter (HCAL).

The ECAL, a homogeneous calorimeter, has a barrel region and two endcap regions. The barrel region is composed of 61200 scintillating lead tungstate (PbWO_4) crystals, and both endcap regions of 7324 crystals. PbWO_4 is a material with high density, therefore it has short radiation length¹ X_0 and a small Moliere radius of 22 mm, providing a fine granularity and a compact calorimeter, which fits into the design of CMS. In order to detect the scintillation light the ECAL uses two types of photodetectors that are installed to the end of each crystal. In the barrel region avalanche photodiodes (APD) are used, which can operate in strong transverse magnetic fields. In the endcaps vacuum phototriodes (VPT) are used in order to cope with higher levels of radiation. Photodetectors convert the scintillation light into an electrical signal that is amplified and sent for analysis. [30, 33] The barrel part of the ECAL (EB) covers the pseudorapidity range $|\eta| < 1.479$ and the endcaps (EE) cover $1.479 < |\eta| < 3.0$. The PbWO_4 crystals in the EB form 36 supermodules, each holding 1700 crystals. The endcap regions are divided into 2 halves, called Dees, each holding 3662 crystals. The ECAL also contains preshower detectors installed in front of both endcap. The preshower detectors aim to identify neutral pions in the endcaps, help with the identification of electrons against minimum ionizing particles (MIP), and improve the position determination of electrons and photons with high granularity. [20] A 3-D view of the barrel and endcap electromagnetic calorimeter is shown in Figure 11.

¹Radiation length is a characteristic property of a material that shows energy loss of particles interacting with the material electromagnetically with respect to distance.

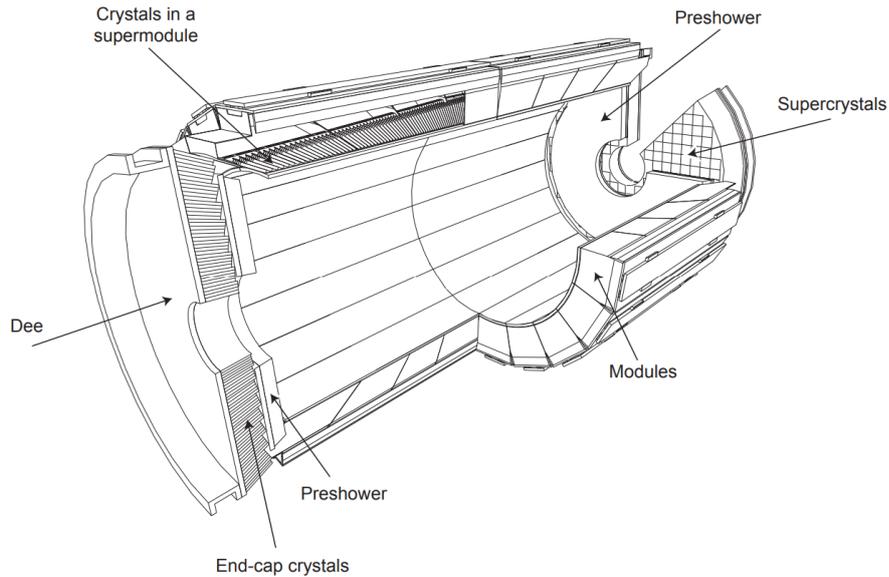


Figure 11. A 3-D view of the electromagnetic calorimeter [20]

The HCAL consists of four parts: The Hadron Barrel (HB), the Hadron Endcap (HE), the Hadron Outer (HO) and the Hadron forward (HF) calorimeter. It is a sampling calorimeter, which means it is constructed of passive absorber layers alternating with active detector layers. The absorber layers force showering of hadrons, and the active layers detect the energy of particles in the showers. The HB and HE are placed inside the solenoid magnet, and the HO and HF are positioned outside. The positioning of the sub calorimeters can be seen in Figure 12.

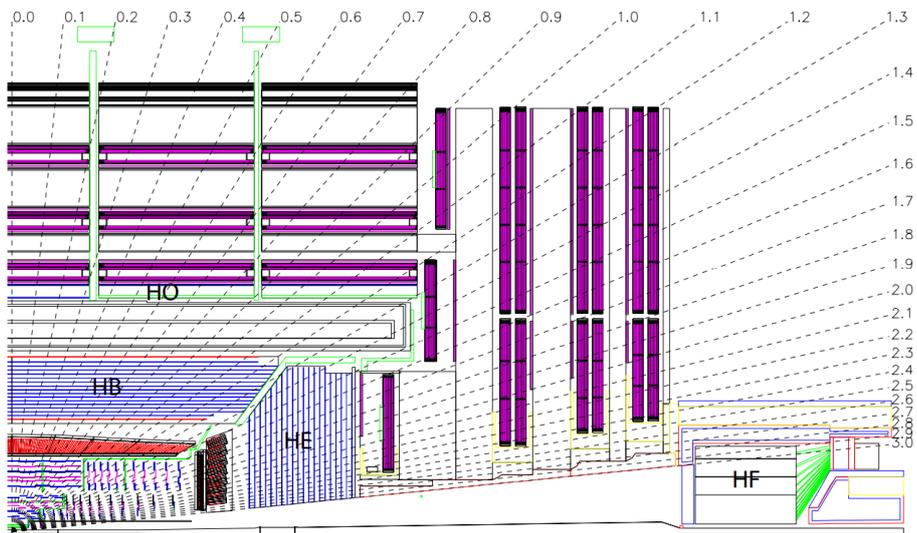


Figure 12. Longitudinal view of the CMS detector showing the locations of the hadron barrel (HB), endcap (HE), outer (HO) and forward (HF) calorimeters. [20]

The HB is placed between the EB (at $r = 1.77$ m) and the CMS solenoid magnet's coil (at $r = 2.95$ m), and consists of alternating flat brass absorber plates and plastic scintillator layers, that are divided into 36 wedges in ϕ and 16 towers in η , covering the pseudorapidity range of

$|\eta| < 1.3$. The HE covers the pseudorapidity range $1.3 < |\eta| < 3$, and is located at the ends of the solenoidal magnet, thus requiring the absorbing layers be made from a non-magnetic material. For the absorbing layers brass plates are chosen, which are 79 mm thick with 9 mm gaps to accommodate the scintillators. In total there are 18 layers of scintillators, 17 made of Kuraray SCSN81 and an outer layer made of Bicron BC408. The scintillation light is collected by wavelength shifting (WLS) fibres that are inserted in the grooves of the scintillators. The HO is placed after the coil and uses the coil itself as the absorber. It is used to identify late starting showers and to measure the shower energy deposited after the HB. A large iron yoke, designed in the form of five 2.536 m wide rings along the z -axis, covers the solenoid magnet. The HO is placed as the first sensitive layer in each of these five rings, followed by muon detectors. The HF extends the angular coverage of the detector up to $|\eta| < 5$. It is exposed to very high levels of radiation, thus quartz fibres were chosen as the active medium and steel is used for the absorber material. The quartz fibers provide a fast collection of Cherenkov radiation using photomultipliers. Charged particles entering the HF produce particle showers, where only electrons and positrons are fast enough to produce Cherenkov light, making the calorimeter mainly sensitive to the electromagnetic component of particle showers. [20, 33]

3.2.4. Muon system

The muon system forms the outermost layer of CMS and has three functions - muon identification, momentum measurement and triggering. Due to the shape of the iron yoke of the solenoid magnet, the system has a cylindrical barrel section (MB) and two planar endcap regions (ME). For the task of identifying muons, the only measurable SM particle crossing the whole detector, the CMS uses three types of gaseous particle detectors. In the barrel region, where the neutron induced background is insignificant and both the muon rate and the magnetic field are low, Drift Tube (DT) chambers are used, which cover the pseudorapidity region of $|\eta| < 1.2$. In the two endcap regions, where the magnetic field is large and non-uniform, and both the muon rate and background levels are high, Cathode Strip Chambers (CSC) are utilized due to their high rate and precision capabilities, increasing the pseudorapidity coverage up to $|\eta| < 2.4$. Lastly, Resistive Plate Chambers (RPC) are installed in both regions and provide high-precision timing information for triggering. [20] The operation procedure of the three gaseous detectors is similar - traversing muons ionize gas atoms in the detectors, a high voltage disperses the charges, and a stream of electrons induces mirror charges on the cathode, which are amplified and measured. [30]

Drift tube chambers are divided into 12 sectors in the $r - \phi$ plane, forming 4 stations (MB1-MB4) at different radii distributed among the layers of the magnet flux return yoke, as can be seen in Figure 13.

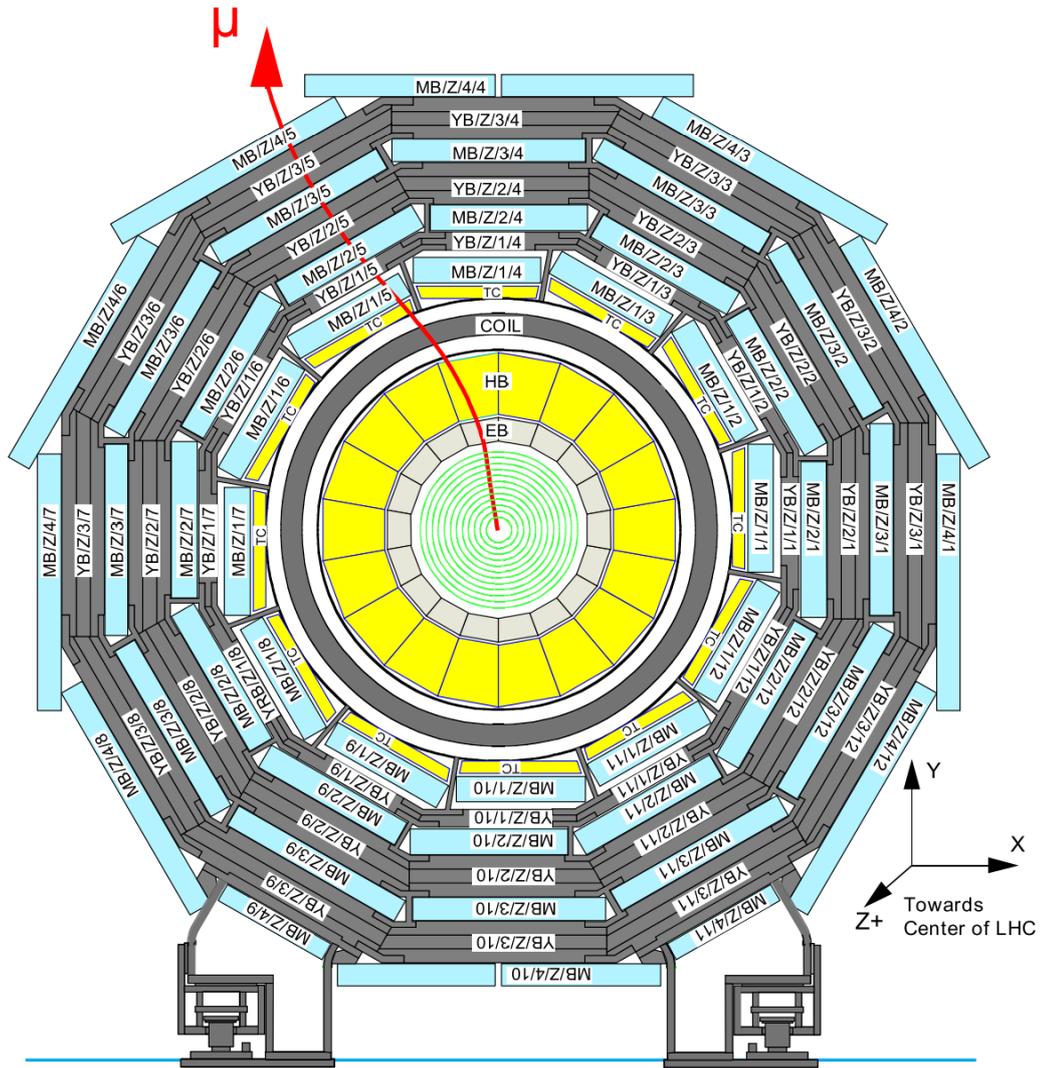


Figure 13. Layout of the CMS MB Drift tube chambers for one wheel of the iron yoke in $r - \phi$ plane. The MB is divided into 12 sectors, which form four stations interleaved by the return yoke [30]

There are five such wheels in z -direction, making up the barrel muon system, consisting of 240 chambers in total. The inner three stations (MB1-MB3) contain three superlayers (SL), with the lower one (SL1) and upper one (SL3) measuring the muon coordinates in the $r - \phi$ plane, and the middle one (SL2) measuring coordinates in the $r - z$ plane. MB4 has only two superlayers, measuring coordinates in the $r - \phi$ plane. [30] Each SL is built of four layers of rectangular drift cells staggered by half a cell. The cells have a width of 42 mm and a height of 13 mm, and are made out of a central anode wire, two cathode strips at the edges, and field shaping electrodes at the top and bottom. The cross section of a drift cell can be seen in Figure 14. The cells are flushed with a gas mixture of 85% Ar and 15% CO₂, kept at atmospheric pressure, which provides a drift velocity of about 55 $\mu\text{m}/\text{ns}$, resulting in a maximum drift time of 380 ns, corresponding to about 16 bunch crossings. [33]

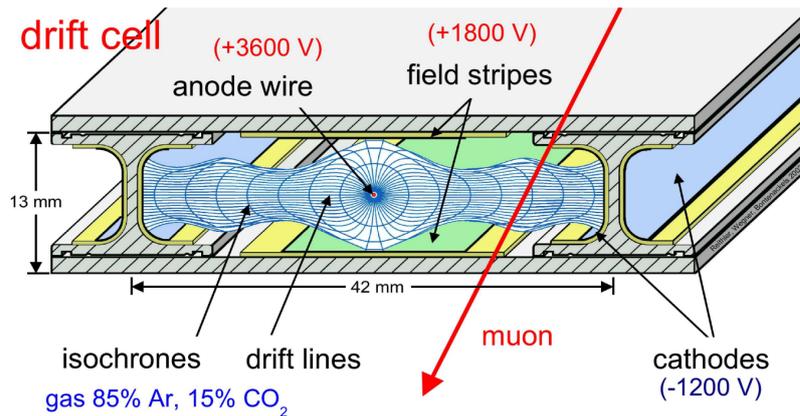


Figure 14. Cross section of a CMS drift cell with drift lines of electrons and isochrones [33]

Cathode Strip Chambers are arranged in four stations per endcap, and are perpendicular to the beam axis. Each station, arranged as a disc, consists of two rings of chambers. The inner ring is segmented into 18 trapezoidal chambers and the outer one into 36, giving full ϕ coverage. A single chamber is composed of six layers of anode wires and seven layers of cathode strips in an interleaved manner, and are separated by a 9.5 mm gas gap, that is filled with a mixture of 30% Ar, 50% CO₂ and 20% CF₄. The cathode strips are aligned perpendicular to the wires in radial direction, allowing for measurements in ϕ , whereas the anode wires run along ϕ , allowing for measurements in r . A sketch of a cathode strip chamber can be seen in Figure 15. [33]

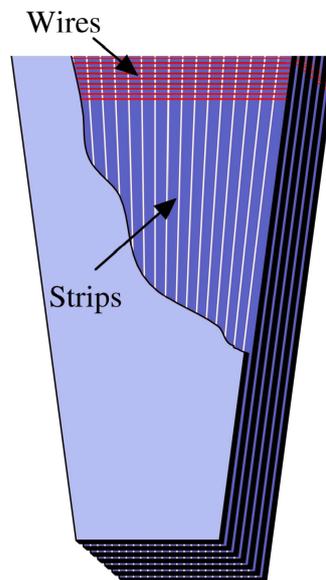


Figure 15. Sketch of a muon cathode strip chamber [33]

Resistive Plate Chambers are gaseous parallel-plate detectors, able to provide very precise timing information. With time resolution of a few nanoseconds, they are mainly designed for trigger purposes. In the barrel region, the RPCs are installed on both sides of the two innermost DT stations, while only one RPC is attached to each of the outer two stations. In the endcap region, there is one RPC per muon station. RPC chambers are made of a pair of

parallel bakelite plates, separated by a 2 mm small gap, that is filled with a gas mixture of 96% $C_2H_2F_4$, 3.5% iC_4H_{10} and 0.5% SF_4 . The plates are coated with graphite electrodes, and a high voltage of 9.5 kV is applied to them. CMS uses double gap RPC units, that consist of insulated aluminium readout strips, placed between two single gap pairs. The layout of a double gap RPC can be seen in Figure 16. [33]

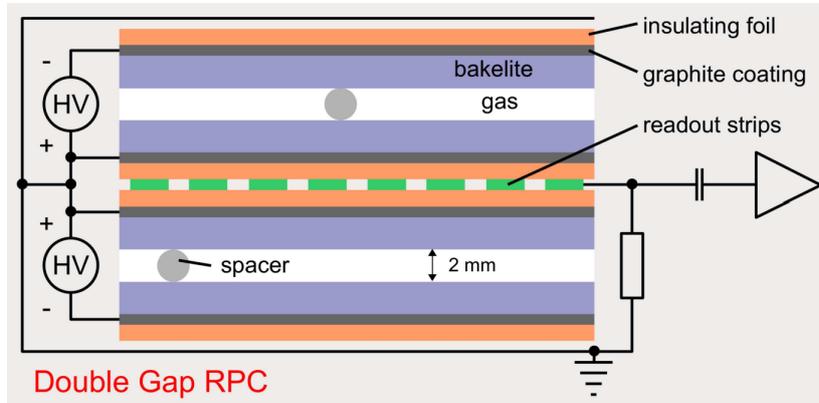


Figure 16. Cross section of a double gap resistive plate chamber [33]

3.2.5. Triggers

To be able to measure even the rarest types of events, the LHC provides proton-proton collisions at high interaction rates, having an average of 25 interactions per bunch crossing at a rate of 40 MHz. Only a small fraction of events can be stored for later analysis. It is the job of the trigger system to select potentially interesting events for offline storage. The selection is done in two consecutive steps called Level-1 (L1) trigger and High Level-Trigger (HLT). [20] The CMS level-1 trigger is designed to reduce the initial bunch crossing rate of 40 MHz to 100 kHz. It is a hardware system with a fixed latency, that takes information from the calorimeters and muon detectors to decide within $4 \mu s$ of a collision if an event should be accepted or rejected. The L1 trigger system consists of multiple stages. First a regional calorimeter trigger (RCT) receives transverse energies and quality flags from the ECAL and HCAL, to process and form electron and photon candidates. The RCT sends the information to the global calorimeter trigger (GCT), that sorts the candidates further, finds jets using transverse energy sums, and calculates global quantities such as missing transverse energy (E_T^{miss}). Regarding muons, each of the three muon detector systems transmit information to the global muon trigger (GMT), that builds muon candidates. Each candidate is assigned a quality code, transverse momentum p_T , and a position (η, ϕ) in the muon system. The GMT combines muon candidates found by more than one system, in order to eliminate a single candidate passing multiple muon triggers. Both the information from the GCT and GMT is forwarded to the final step of the L1 trigger system, which is called the global trigger (GT), that finally takes the decision whether to reject the event or accept it for subsequent evaluation by the HLT. [34]

If an event is accepted by the L1 trigger, the full detector information is read out and sent to the HLT, where a simplified event reconstruction is performed. Physics objects such as electrons, muons and jets are reconstructed and identification criteria are applied in an effort

to select events with possible interest for data analysis. The HLT hardware consists of a single processor farm composed of about 1000 computers, that run algorithms to determine type and multiplicity of particles in events. The data processing is built around HLT paths, which are a set of algorithmic processing steps, that are run in a predetermined order that reconstruct physics objects and make selections on them. Accepted events are sent to the storage manager, where the data is stored locally on disks and eventually transferred to the CMS computing center for offline processing and permanent storage. [34]

3.2.6. The Particle-Flow algorithm

The CMS detector allows to identify and reconstruct electrons, muons, photons, charged hadrons and neutral hadrons using the Particle-Flow (PF) algorithm, determining the direction, energy and type of the particles. These particles are then used to build jets, determine the missing transverse energy (MET), reconstruct and identify τ leptons, and quantify lepton isolation with respect to other particles. [35] After the collisions, particles first enter the inner tracking system, where the charged particle trajectories and origins are reconstructed from hits in the sensitive layers. The tracking system is immersed in the magnetic field of the solenoid, which bends the trajectories and allows the measurement of the electric charges and momenta of charged particles. Electrons and photons are absorbed in the ECAL, where the electromagnetic particle showers are detected as clusters of energy, recorded in the cells of the calorimeter, from which the energy and direction of the particles is determined. Charged and neutral hadrons initiate hadronic particle showers by interacting with the absorber layers of the HCAL. The showers are detected as clusters of energy and are used to estimate the energy and direction of the corresponding particles. Only muons and neutrinos pass through the calorimeters with hardly any interactions. Neutrinos escape undetected and are measured indirectly from momentum imbalance, while muons produce hits in muon detectors, allowing for their identification and giving further information for momentum reconstruction. Jets consist of hadrons and photons, the energy of which can be inclusively measured by the calorimeters without any attempt to separate individual jet particles. [36] Figure 17 summarizes the specific particle interactions in a transverse slice of the CMS detector. The PF algorithm reconstructs and identifies each individual particle in an event, using the information from the CMS detector elements.

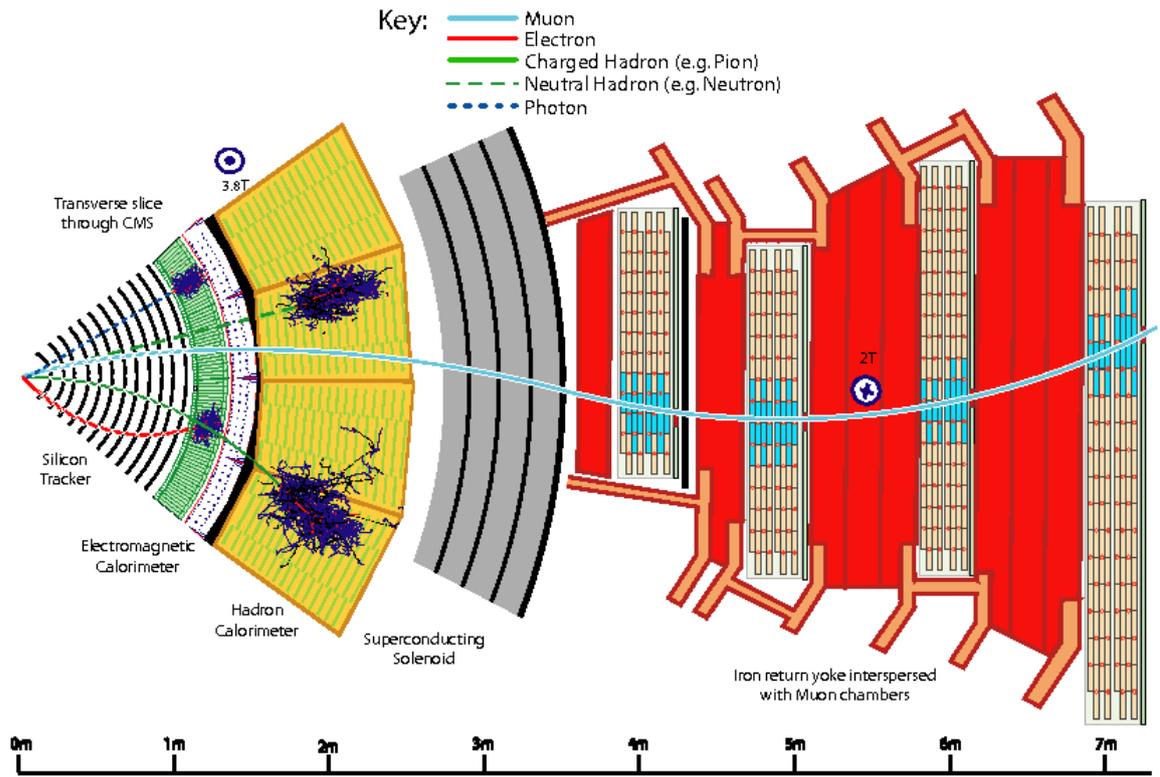


Figure 17. Schematic of the specific particle interactions in a transverse slice of the CMS detector. [36]

4. Machine learning methods

Experimental high-energy physics (HEP) has two main objectives - probing the SM with increasing precision and searching for physics beyond the SM. Both of these tasks require correct and effective identification of rare signals hidden in very large backgrounds. [37] Machine learning (ML) algorithms are the convenient solution for satisfying these requirements. Hence, the use of multivariate analysis (MVA) techniques in the analysis of data from colliders is growing. Recent analyses have given a clear example of improving the signal significance through the use of decision tree-based tools and neural networks, that identify the optimal way to analyse and separate signal events from background events, such as boosted decision trees (BDT). [38]

BDTs are used in HEP for separating different physics processes by assigning them different functional values based on a set of discriminating input variables. This is done by applying a boosting algorithm on a supervised learning method called decision tree learning. The boosting algorithm is a method where many “weak” classifiers are combined to achieve a final powerful classifier [39]. In this thesis XGBoost [40] is used for the BDT implementation.

The aim of a supervised learning problem for classification is to approximate an existing mapping f , that relates elements of a dataset \mathcal{D}

$$y = f(x), \quad \mathcal{D} = \{x_i, y_i\}, \quad (4.1)$$

where x_i is input data with multiple features, and y_i is the final prediction or target variable. Using the characteristics of the features of the input data, and the classification as target, we get an approximation to f :

$$\hat{f} = \mathcal{A}(\mathcal{D}), \quad (4.2)$$

which is used to classify or categorise inputs. In a physics event classification task, events are chosen as the input data x_i , and y_i are the given classification labels for the events, such as a label 1 for signal events and a label 0 for background events. The input data features are discriminating variables that characterise each event, for instance the number or transverse momentum of specific particles in events. The mapping is found by solving an optimisation task, which measures the quality of the prediction for a single x_i with a loss function $L(y_i, \hat{f}(x_i))$, which in order to get the best model should be minimized. The data is divided into training, validation, and testing sets, which allow assessing the quality of the model. Training data is used to fit and tune the model, validation data is used to measure the accuracy of the model during training, and testing data is used to measure the performance of the model. It is important to have a model that is well tuned, advanced enough to interpret and express underlying patterns in data, but also simple enough to produce genuine results. [38] If the model is too complex, it will learn too much about the specific details of the training data, which do not represent the underlying patterns. When applying the model on testing

data, the performance will be worse than that of a model less complex, because it does not generalise well. This issue is called overtraining. [41] The section that follows gives a brief overview of the concept of boosted decision trees, used in the analysis for separating different di-Higgs production mechanisms from each other, and signal events from background.

4.1. Boosted decision trees

Decision trees are rooted binary trees, that learn a tree-like model of decision rules from input data to make predictions. A given dataset is split recursively into smaller groups, called nodes, each representing a decision point. The tree branches out from the node according to the decisions that are made. The final prediction is represented by the leaves of the tree, which are assigned a score. [42] For a physics classification task, a decision tree sorts input events into leaves, which are assigned a score, giving each event a functional value.

In the case of BDTs, an ensemble of trees is used. Instead of training the model with a single tree, boosting algorithms iteratively add new trees in an additive manner one at a time, and within each tree one split is added to the pre-existing structure of trees from the previous iteration, thus minimizing (boosting) the loss function, resulting in a smaller error rate and better performance. The final prediction score for a single input is obtained by the sum of all leaf scores assigned to the input. The training of the BDT is controlled by hyperparameters, which are set before the training. [43] The XGBoost hyperparameters used in BDT training in this thesis are:

- `n_estimators`: Maximum number of trees.
- `subsample`: Controls the number of training events that are used to grow each tree to a fraction of the full training sample [44].
- `colsample_bytree`: Specifies the number of different features that are used in a tree.
- `gamma`: Minimum loss reduction required to make a further partition on a leaf node of the tree.
- `learning_rate`: Controls the effect that trees added at a later stages of boosting iterations have on the output of the BDT in comparison with the effect of trees added at an earlier stage. Small values decrease the effect of trees added during the boosting iterations. [44]
- `max_depth`: Maximum depth of a tree. A large depth will make the model more complex and lead to overtraining.
- `min_child_weight`: Minimum sum of weighted inputs in each leaf.

The choice of potential input variables or discriminating features in the input dataset needs to be optimized as well. In this thesis, it is done by finding physics inspired variables and selecting the most significant ones. A variable ranking feature implemented in XGBoost evaluates the gain of each variable. The gain is measured as the improvement in accuracy of a branch before the split compared to the accuracy of the two branches created after the split [43].

5. The qqHH analysis in multilepton final states

This chapter gives a brief overview of the CMS HH→multilepton analysis [3] and describes the extension of this multilepton analysis to the sub-dominant qqHH production mode in events containing two same-sign electric charge electrons or muons with the possibility of an additional hadronically decaying τ (denoted as $2lss + 0/1 \tau_h$) in final state.

5.1. CMS HH → multilepton analysis

The CMS HH→multilepton analysis searches for Higgs boson pair (HH) production in final states with multiple electrons (e), muons (μ), or hadronically decaying τ leptons (τ_h), using data recorded by the CMS experiment in proton-proton collisions with center-of-mass energy of 13 TeV during LHC Run 2 (2016-2018), corresponding to an integrated luminosity of 138 fb^{-1} . The analysis focuses on both non-resonant and resonant HH production, where resonant HH production involves the decay of a possible new heavy resonance X into a Higgs boson pair. Seven different final states or "search categories", distinguished by the final decay products, are included in the analysis:

- $0l + 4\tau_h$
- $1l + 3\tau_h$
- $2l + 2\tau_h$
- $3l + 1\tau_h$
- $3l + 0\tau_h$
- $2lss + 0/1\tau_h$
- $4l + 0\tau_h$

where " l " refers to an electron or muon, τ_h refers to a hadronically decaying τ lepton, and "ss" refers to same sign electrical charge electrons or muons. These final states target HH signal events in which the Higgs boson pair decays either into WWWW, WW $\tau\tau$, or $\tau\tau\tau\tau$. Monte Carlo event simulation is used for the modeling of the HH signal and most backgrounds, and BDT classifiers are used to distinguish the HH signal from backgrounds. The inputs to the BDT classifiers include various physics inspired variables of reconstructed particles, such as p_T , η , invariant mass, angular separation ΔR , visible mass of the Higgs boson pair, scalar p_T sum of all reconstructed particles. Additional inputs are used, which allow the BDTs to learn that distributions in the input observables for HH signal events change as a function of the model parameters. For non-resonant HH production these parameters are a set of five effective field theory (EFT) couplings, and for resonant HH production the mass of a heavy particle X is used. In order to test new theories and hypothesis effectively, EFT couplings are studied instead of the whole model, as they represent a variety of "full" theory models with the same phenomenology at once. The values for the EFT couplings for the CMS HH→multilepton analysis are chosen according to 12 EFT benchmark (BM) scenarios, taken from [45]. BDTs for non-resonant HH production are trained on HH samples corresponding both to the SM prediction and to the 12 BM scenarios, indicated by 13 binary inputs. The BDTs for resonant HH production are instead trained on a full set of 19 resonant masses as an input. For the extraction of the HH signal, a binned maximum likelihood fit is applied

to the distributions in the output of the BDT classifiers, in the seven search categories as well as kinematic distributions in two background control regions. The data from each of the three years of LHC Run 2 are fit separately. Background contributions are classified as either “reducible” or “irreducible”. Three types of reducible backgrounds are considered in the analysis: misidentified l or τ_h candidates (“fakes”), mismeasurement of the electron charge (“flips”), and electrons from photon conversions. Irreducible backgrounds arise from events in which all selected l and τ_h candidates come from W, Z or single H boson decays, and are reconstructed with the correct charge. The “fakes” and “flips” backgrounds are both determined from data, while electron conversions and irreducible backgrounds are modeled using MC simulation.

The SM signal strength parameter is defined as $\mu_{SM} = \sigma_{HH}/\sigma_{HH}^{SM}$, where σ_{HH} is the measured HH production cross section, and σ_{HH}^{SM} is its predicted value in the SM. The 95% confidence level (CL) upper limit on μ is calculated using approximation of the CL_s method [46, 47], and limits on $\mu = \sigma_{HH}/\sigma_{HH}^{theo}$ are used to calculate limits on HH signal cross sections and coupling strength modifiers ($\kappa_\lambda, C_2, C_V, C_{2V}$). Observed and expected 95% CL upper limits on the signal strength modifier μ for non-resonant HH production, obtained for both individual search categories and the combination can be seen in Figure 18. [3]

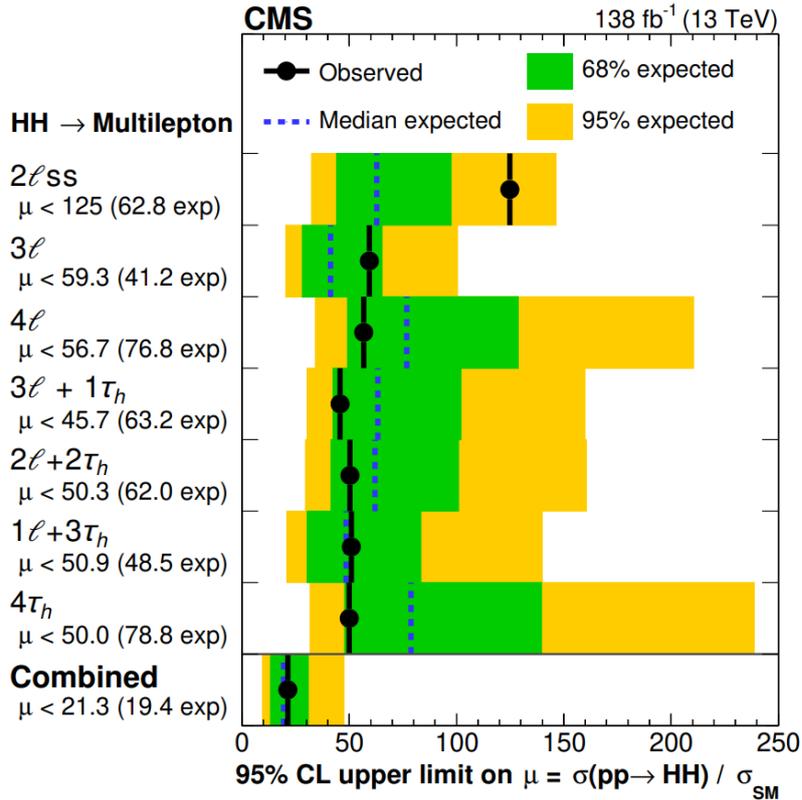


Figure 18. CMS $H \rightarrow$ multilepton analysis observed and expected 95% CL upper limits on the signal strength modifier μ for non-resonant HH production. [3]

The observed (expected) 95% CL interval for the Higgs boson trilinear self-coupling strength modifier κ_λ is measured to be $-6.98 < \kappa_\lambda < 11.17$ ($-6.98 < \kappa_\lambda < 11.73$). Due to the low expected qqHH sensitivity, the original CMS $HH \rightarrow$ multilepton κ_λ analysis does not contain any

dedicated result on the qqHH rate or constraints on C_{2V} . Still, the analysis results allow to extract these results as shown in another thesis [43]. The observed (expected) HHVV coupling strength modifier C_{2V} interval is set at $-3.42 < C_{2V} < 5.56$ ($-2.73 < C_{2V} < 4.83$). The 95% confidence level (CL) upper limit on non-resonant HH production cross section as a function of κ_λ , and C_{2V} can be seen in Figure 19. [43]

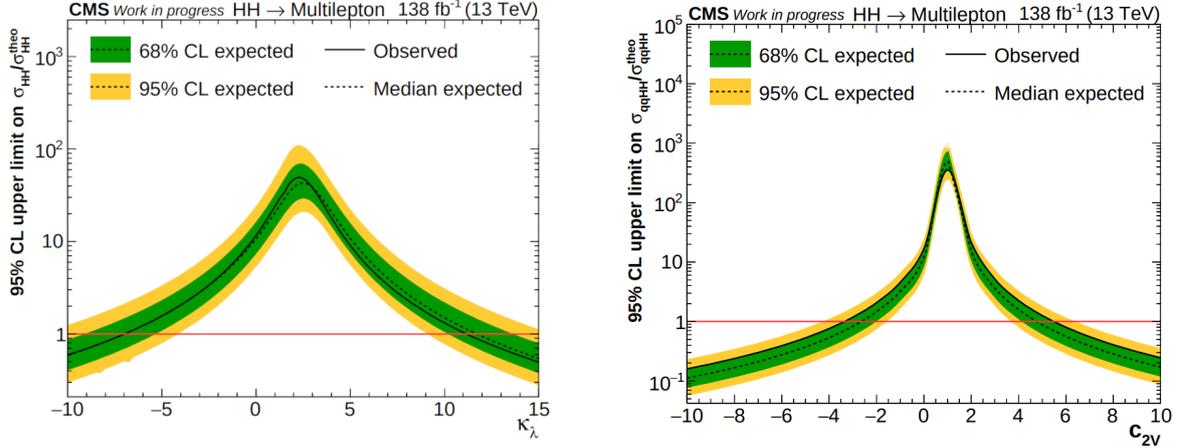


Figure 19. Observed and expected 95% CL upper limits on the Higgs boson self-coupling strength modifier κ_λ on the left, and Higgs boson coupling to vector bosons C_{2V} on the right. [43]

5.2. Analysis overview

For the extension of the CMS HH \rightarrow multilepton analysis to better suit the challenges of a qqHH focused search, the strategy for separating signal and background in the signal extraction has been expanded. In addition to the signal/background separating BDT, a second BDT separating ggHH-like events from qqHH-like events is trained, and the signal/background BDT reoptimized, allowing the use of a 2-D MVA distribution in the signal extraction. With both BDTs including new variables targeting the qqHH signature, based on the additional jets expected in qqHH. For identifying these jets, multiple tagging strategies have been tried.

The $2lss + 0/1\tau_h$ search category is chosen due to having the biggest HH signal from the seven search categories. The channel is characterized by two leptons, which are either electrons or muons, and "ss" (same-sign) indicates that the electric charge of the leptons is the same. The ss requirement is needed in order to suppress $t\bar{t}$ background. This channel also takes into account events up to a single τ_h in addition to the two leptons in the event. The main backgrounds for this channel come from WZ production, $W\gamma$ events in which the photon decays to an electron-positron pair (e^-e^+) and one of them fails to get reconstructed, and from events in which one or both reconstructed leptons come from misidentified jets. Additionally other backgrounds, such as same-sign W boson pairs and WWW production contribute. HH signal events selected for the $2lss$ category are made up of about 70% HH \rightarrow WWWW decays and 30% from HH \rightarrow WW $\tau\tau$ decays. Data driven background estimation is discussed in 5.5, event selection is discussed in section 5.6 and BDT training in section 5.7.

5.3. Datasets and Monte Carlo simulation

The data and simulation used both in this, as well as in the original multilepton analysis correspond to the full LHC Run 2 dataset with 138 fb^{-1} of proton-proton collision data recorded at a center-of-mass energy of $\sqrt{s} = 13 \text{ TeV}$ and 25 ns bunch crossing period between 2016 and 2018. The events have been recorded using a combination of single, double and triple lepton triggers and by triggers that are based on two hadronically decaying τ leptons (denoted as τ_h) or a lepton and a single τ_h in the event. Only data that was taken when all detector systems were fully operational is included in the analysis. On average approximately 30 inelastic pp interactions occurred per bunch crossing. [3]

Monte Carlo (MC) simulation is used to create samples of HH signals and of background events, for the purpose of estimating signal and background contributions in the signal extraction and to train machine learning algorithms.

Samples produced from LO matrix elements are generated with MADGRAPH5 aMC@NL [48], and samples at NLO matrix element level are generated with MADGRAPH5 aMC@NLO and POWHEG v2 [49–51]. Parton showers, hadronization processes and decays of τ leptons are modelled using the generator PYTHIA [52] with the tunes CUETP8M1, CUETP8M2 or CUETP8M2T4 [53] in 2016 era MC samples, and with tune CP5 [54] in 2017 and 2018 samples. The event generator tunes are based on the Monash tune [55]. Samples produced with CUETP8M* tune use the NNPDF3.0 set, and samples produced with CP5 tune use the NNPDF3.1 set [56]. [3]

HH signal samples are generated at LO (leading-order) and at NLO (next-to-LO) accuracy in QCD to simulate nonresonant HH production in ggHH and qqHH production modes. Both Higgs bosons are required to decay into either WW, ZZ or $\tau\tau$. The ggHH LO samples are used in BDT training, while the ggHH NLO samples are used in the HH signal extraction from the data. Separate ggHH samples are produced for SM HH production and for 12 benchmark scenarios in the EFT approach [45].

Background MC samples include processes producing a single boson (W, Z), two bosons (WW, WZ, ZZ, $W\gamma$, and $Z\gamma$), three bosons ((WWW, WWZ, WZZ, ZZZ, $WZ\gamma$), one or two top quarks, top quarks associated with one or more bosons (t , $t\bar{t}$, $t\bar{t}W$, $t\bar{t}Z$, $t\bar{t}H$, tHq , tHW), and a single Higgs boson (ggH, qqH, WH, ZH). All samples that include a Higgs boson are produced for a Higgs boson mass of 125 GeV. Most samples, including the dominant WZ and ZZ backgrounds, are generated at NLO and scaled to cross sections computed at NNLO in QCD. [3]

5.4. Object reconstruction

The CMS detector allows reconstruction of electrons, muons, taus, and jets. The general reconstruction is described in section 3.2.6.

5.4.1. Electrons

Electron reconstruction is based on the combined information from the inner tracking system and the calorimeters. Because of the large amount of material in the inner tracking system, most electrons emit a sizeable fraction of their energy in the form of bremsstrahlung photons, which convert into an electron-positron pair, before getting absorbed by the ECAL. This means that by the time the electron or photon reaches the ECAL, it may no longer be a single particle, but rather consist of a shower of multiple electrons and photons. An algorithm is used for combining the clusters from the individual particles into a single object, making it possible to recover the energy of the primary electron. The emittance of bremsstrahlung photons causes electrons to lose momentum, thus changing the curvature of their trajectory in the tracker. For the estimation of the track parameters, a dedicated tracking algorithm, based on the Gaussian sum filter (GSF) is used. A dedicated group of specialists in CMS known as the EGamma physics object group (POG) has developed means of reconstruction and identification for electrons. The reconstruction is done in the following steps:

- The formation of clusters is started by the energy reconstruction algorithm by grouping together cells of the ECAL with energies exceeding a predefined threshold. The cluster containing most of the energy is defined as the seed cluster, forming a supercluster (SC) when combined with nearby cluster. SCs account for photon conversions and losses from bremsstrahlung.
- Suitable trajectory seeds in the pixel detector, matching the position of the SC and the trajectory of an electron, are used for seeding GSF tracks.
- At the same time with the two steps described above, all reconstructed tracks in the event are tested for electron trajectories, and are also used for seeding GSF tracks if they are suitable. Tracks with $p_T > 2$ GeV, reconstructed from hits in the inner tracking system, are denoted as generic tracks.
- ECAL clusters, SCs, GSF tracks, and generic tracks are imported into the PF algorithm that groups the elements together into blocks of particles.
- The blocks are fixed into electron and photon objects, starting from either a GSF track or a SC respectively. Objects with an associated GSF track, passing loose selection requirements, are labeled as electrons.

[57]

The separation of prompt leptons from non-prompt ones is done with a BDT-based algorithm, developed for the ttH multilepton search [58–60]. The BDT outputs are values between 0 and 1, where non-prompt leptons are assigned lower scores than prompt leptons. The EGamma POG has developed an algorithm for separating electrons from jets, making use of MVA methods [61]. Three working points have been defined: WP-loose, WP-90 and WP-80, corresponding to 98%, 90% and 80% signal efficiency respectively. In this analysis electrons are required to pass WP-loose. For this analysis three different electron identification IDs have been defined, summarized in Table 4. The tight ID electrons are used in the signal event selection. For the estimation of background from misidentified leptons, the fakable ID

together with the tight ID are used. The loose electron ID is used for low level event vetoes and for cleaning collections of jets and τ_h from overlapping electron candidates. [43] This analysis targets electrons with p_T as low as 10 GeV, and uses a modified p_T variable, defined as cone- p_T , in order to avoid potential biases with the fakeable lepton ID definition and the data driven background estimation for fake leptons. Isolated leptons are separated from leptons in jets by means of the scalar sum p_T values of charged particles, neutral hadrons and photons reconstructed within a narrow cone centered on the direction of the lepton. The size of the cone shrinks inversely proportionally with the p_T of the lepton. The cone size is referred to as mini isolation and is denoted as I_e for electrons, and I_μ for muons. More information about electron isolation and cone- p_T can be found in reference [3].

Table 4. Loose, fakeable and tight selection criteria for electrons

Observable	Tight	Fakable	Loose
cone- p_T	$> 10 \text{ GeV}$	$> 10 \text{ GeV}$	$> 7 \text{ GeV}$
$ \eta $	< 2.5	< 2.5	< 2.5
I_e	$< 0.4 \times p_T$	$< 0.4 \times p_T$	$< 0.4 \times p_T$
EGamma POG MVA	$> \text{WP-loose}^2$	$> \text{WP-90} (> \text{WP-loose}^2)^\dagger$	$> \text{WP-loose}^2$
Prompt-e MVA	> 0.30	$< 0.30 (> 0.30)^\dagger$	-

† denotes electrons failing (passing) the requirement prompt-e MVA > 0.30 .

5.4.2. Muons

Muons are reconstructed in various ways, with the final collection being composed of three different types:

- Standalone muon: Hits within each DT or CSC detector are clustered to form track segments, which are used as seeds for the pattern recognition algorithm that gathers all DT, CSC and RPC hits along the muon trajectory, resulting in a standalone muon track.
- Global muon: Each standalone muon track is matched to a track in the inner tracking system (inner track) if the parameters of the two tracks are compatible propagated onto a common surface. Hits from the inner track and the standalone muon track are combined and fit to form a global muon track.
- Tracker muon: Inner tracks with $p_T > 0.5 \text{ GeV}$ and a total momentum $p > 2.5 \text{ GeV}$ are extrapolated to the muon system. If at least one muon segment matches with the extrapolated track, the inner track qualifies as a tracker muon track.

Global muons and tracker muons that share the same inner track are merged into a single muon candidate by the PF algorithm, that associates muon energy deposits in the ECAL and HCAL to the muon.

The CMS collaboration has developed muon identification and isolation criteria to separate muons from jets. The muon candidates used in the analysis are required to pass the loose PF muon identification criteria. Three muon IDs are defined for this analysis: tight, fakeable, and loose. Tight muons are used for the signal event selection, fakeable muons are used for fake

background estimation (background coming from misidentified leptons), and loose muons are used for separating prompt leptons from non-prompt, and for low level event vetoes. The three muon IDs are summarized in Table 5. [36]

Table 5. Loose, fakeable and tight selection criteria for muons

Observable	Tight	Fakable	Loose
p_T	$> 10 \text{ GeV}$	$> 10 \text{ GeV}$	$> 5 \text{ GeV}$
$ \eta $	< 2.4	< 2.4	< 2.4
I_μ	$< 0.4 \times p_T$	$< 0.4 \times p_T$	$< 0.4 \times p_T$
PF muon	$> \text{WP-loose}$	$> \text{WP-loose}$	$> \text{WP-medium}$

5.4.3. Jets

Color charged particles cannot exist by themselves in nature because of QCD confinement only allowing for colorless states. The quarks and gluons hadronize in the detector by creating other colored objects around them in order to form colorless objects. A collection of these colorless objects travelling in the same direction form a narrow jet of particles. Jets are studied in order to measure the properties of quarks and are reconstructed using the anti- k_t algorithm [62] with a distance parameter R , that clusters together PF objects. The parameter R regulates the approximate size of the jets built by the algorithm. Jets reconstructed with $R = 0.4$ are known as AK4 jets and are used as standard jets, while jets reconstructed with $R = 0.8$ are known as AK8 jets and are in this analysis used for boosted W boson decays. The AK8 jets in this analysis are reconstructed from all PF candidates except leptons, and are known as lepton subtracted AK8 jets

The analysis uses three types of jets: AK4 jets, AK8 jets and b-tagged jets. AK4 and AK8 jets are used to reconstruct hadronically decaying W bosons decaying either into two AK4 jets or a single AK8 jet. AK4 jets are required to pass $p_T > 25 \text{ GeV}$ and be within the geometric acceptance of $|\eta| < 2.4$. AK8 jets are required to have $p_T > 100 \text{ GeV}$, $|\eta| < 2.4$, subjettiness parameter [63] $\tau_2/\tau_1 < 0.75$, contain two subjets with $p_T > 20 \text{ GeV}$ and $|\eta| < 2.4$, and at least one tight ID lepton within $\Delta R < 1.2$ around the jet. Events containing b-tagged jets coming from top quark or from single Higgs boson decay are vetoed in the analysis, in order to reduce background events. B-tagged jets are AK4 jets coming from heavy b-quarks, that pass the DeepJet algorithm [64] DeepJet-L or DeepJet-M working points. The DeepJet algorithm is a neural network-based algorithm, that uses low-level features from as many jet constituents as possible, in order to classify different jet flavours. AK4 jets with high p_T are used as possible VBF jet candidates.

5.4.4. Hadronic τ decays

Because of the short lifetime of a τ lepton, they can only be seen indirectly by their decay products. τ leptons decay either leptonically into an electron or muon with 2 neutrinos or into a hadronically decaying W boson with a neutrino. Therefore only hadronically decaying τ are referred to and reconstructed as τ_h , as leptonically decaying τ are already reconstructed as electrons or muons. Hadronic τ decays (τ_h) are reconstructed by the hadrons plus strips

(HPS) algorithm [65]. The algorithm aims at reconstructing individual hadronic τ decay modes seen in Table 6 [65], where h^\pm denotes a charged hadron, π^0 a neutral pion, and branching ratios show the fraction of particles decaying in a certain mode.

Table 6. Possible hadronic decay modes of τ

Decay mode	Branching ratio $\mathcal{B}(\%)$
$\tau^- \rightarrow h^\pm \nu_\tau$	11.5
$\tau^- \rightarrow h^\pm \pi^0 \nu_\tau$	26.0
$\tau^- \rightarrow h^\pm \pi^0 \pi^0 \nu_\tau$	9.5
$\tau^- \rightarrow h^- h^+ h^- \nu_\tau$	9.8
$\tau^- \rightarrow h^- h^+ h^- \pi^0 \nu_\tau$	4.8

The hadronic τ candidates are built by combining charged hadrons constructed by the PF algorithm with neutral pions. The neutral pions are reconstructed by clustering photons and electrons reconstructed by the PF algorithm within rectangular strips. The size of the strip is narrow in η -direction but wide in ϕ -direction, depending on the p_T of its constituents. Neutral pions mostly decay into a photon pair, where the photons have a high probability of decaying into an electron-positron pair. The bigger the p_T of a charged particle, the smaller the bending of its track by the magnetic field, which results in the strip having a smaller size. The charged hadrons are required to have $p_T > 0.5$ GeV and must originate from the primary vertex of the event.

In order to separate hadronic τ decays from jets, muons and electrons, a multivariate convolutional deep neural network (DNN), referred to as DeepTau [66] algorithm, is used. 42 high-level observables together with low-level information, taken from the silicon tracker, ECAL, HCAL, and the muon detectors, are used as input. The high-level observables include the p_T , η , ϕ , mass of the τ_h candidate, the reconstructed τ_h decay mode, the isolation of τ_h with respect to charged and neutral particles, and the estimated distance covered by τ_h between its production and decay. The low-level information comes from particle activity within two $\eta \times \phi$ grids which are centered on the direction of the τ_h candidate. Different working points (WPs) for the DNN output scores are defined: eight in the discriminator against jets and electrons (VVVLoose, VVLoose, VLoose, Loose, Medium, Tight, VTight, VVTight), and four in the discriminator against muons (VLoose, Loose, Medium, Tight). This analysis uses the Medium WP against jets, VLoose against muons, and VVVLoose and VLoose WPs against electrons. These WPs are chosen to ensure a high signal efficiency while ensuring a low enough τ_h misidentification rate. Two levels of τ_h IDs are defined: the tight ID and fakable ID. Tight τ_h ID is used for the selection of events in the signal region (SR), and the fakable τ_h ID is used to obtain data driven estimation of the misidentified ("fake") τ_h background. The selection criteria applied to the τ_h IDs are summarized in Table 7.

Table 7. Selection criteria for the tight ID τ_h and fakable ID τ_h

Observable	Tight ID	Fakable ID
p_T	$> 20 \text{ GeV}$	$> 20 \text{ GeV}$
$ \eta $	< 2.3	< 2.3
DeepTau vs jets	WP-Medium	$> \text{WP-VVLoose}$
DeepTau vs electrons	$> \text{WP-VVVLoose}$	$> \text{WP-VVVLoose}$
DeepTau vs muons	$> \text{WP-VLoose}$	$> \text{WP-VLoose}$

5.4.5. Reconstruction of event level quantities

Particle detectors are unable to reconstruct only weakly interacting particles such as neutrinos. The contribution of particles that can not be reconstructed are seen in momentum imbalance in the $x - y$ plane, where energy and momentum conservation is required. Neutrinos can be taken into account as MET E_T^{miss} , which is the magnitude of the negative sum of the transverse momentum vectors of all PF candidates reconstructed in the event.

5.5. Data driven background estimation

The effects leading to misidentified l/τ_h and lepton charge misidentification are not well modelled in simulation, thus the corresponding backgrounds are estimated from data [43]. The l/τ_h misidentification (fakes) background, which includes non-prompt leptons is the largest reducible background, and is estimated using the fake factor (FF) method from [67]. The FF method selects events that satisfy the selection criteria discussed in section 5.6, except l and τ_h are required to pass the fakable instead of the tight object selection criteria. The sample of the selected events is referred to as the application region (AR) of the FF method. In order to avoid overlap with the SR, events where l and τ_h pass the tight object selection criteria are vetoed. An estimate of the fake background in the SR is then obtained by applying correctly chosen weights w to the events in the AR. The weights depend on the probability $f_i(p_T, \eta)$ for a single misidentified l or τ_h to pass both the fakable and tight selection criteria. The probabilities are determined in a measurement region (MR) with no real l/τ_h , and measured separately for e, μ, τ_h , parameterized as functions of cone- p_T and η . The weights w are given by the expression

$$w = (-1)^{n+1} \prod_{i=1}^n \frac{f_i(p_T, \eta)}{1 - f_i(p_T, \eta)}, \quad (5.1)$$

where the product extends over all l/τ_h , that pass the fakable but fail the tight selection criteria, and n is the number of such l/τ_h . The alternating sign is used to avoid double counting coming from events with multiple misidentified l/τ_h . [3, 43]

The electron charge misidentification (flips) background is estimated in a similar manner. A sample consisting of events with two opposite sign electric charge electrons, passing all selection criteria of the SR, is selected and assigned adequately chosen weights. The weights are computed by summing charge mismeasurement probabilities. The probability of mismeasuring the charge of either electron is determined using $Z \rightarrow ee$ events, and is parameterized as a function of p_T and η of the electron. For muons, the charge misidentification rate is negligible, thus making the flips background for events containing electrons notably higher

than events with muons. [3]

5.6. Event selection

The $2lss + 0/1\tau_h$ category is focused on selecting HH signal events in the WWWW decay mode, where two W bosons decay leptonically into an electron or muon with its corresponding neutrino, and the remaining W bosons decay hadronically into jets. To better include events from the $2W2\tau$ and $\tau\tau\tau\tau$ HH decay modes, up to one additional τ_h is allowed as well. Events with two or more τ_h are already covered by other search categories of the $HH \rightarrow$ multilepton analysis. The leptons are required to pass tight lepton selection and have $p_T > 25/15$ GeV for the leading (biggest p_T) and sub-leading (second biggest p_T) lepton respectively. The selected leptons are also required to pass tight charge requirements to reduce charge misidentification. Events with more than one τ_h , and events with invariant mass of lepton pairs, passing loose lepton selection criteria, less than 12 GeV, are excluded, in order to avoid any overlap with other channels and eliminate low mass resonances, as they are not covered well by the simulation. The τ_h is also required to pass tight selection criteria. Events containing 2 pairs of same flavoured opposite sign electric charge (SFOS) leptons, passing loose lepton criteria, with a 4-lepton invariant mass < 140 GeV, are also cut in order to guarantee rejection of $H \rightarrow ZZ \rightarrow 4l$ events. [3]

In order to reduce background from the Drell-Yan process, events with either SFOS leptons passing loose lepton selection criteria or electron pair passing tight electron selection criteria are vetoed. Events in this category are also required to contain jets targeting hadronic W boson decays. The W bosons are reconstructed using either AK8 or AK4 jets. Selected events must have a minimum of 2 AK4 jets, and are sub-categorized into boosted, semi-boosted and resolved W-jets categories if the number of AK8 jets are ≥ 2 , 1 and 0 respectively. Since we are searching for vector boson fusion induced di-higgs production, we look for two high transverse momentum jets in addition to the jets that are used for W boson reconstruction. A selection of possible VBF jet candidates is made by choosing jets with high p_T and cleaning them of jets used for W boson reconstruction. The VBF jets are required to have $p_T > 20$ GeV and $|\eta| < 4.7$, and the jet pair with the biggest mass is chosen as the VBF jet pair. This VBF jet selection is verified using a technique known as MC Truth Matching, where generator level particles in simulated signal events are matched to reconstructed objects through angular separation. It is important to note that this technique is only used for verifying the jet finding, not the final analysis. For measuring background from lepton charge misidentification a dedicated sideband, where the same-sign charge requirement for the two leptons is inverted, is used. In an effort to reduce background contamination from processes with top quarks, events involving b-tagged jets are rejected. The selection conditions used for the signal region are summarized in Table 8.

Table 8. The selection criteria for the signal region of the $2lss + 0/1\tau_h$ channel

Cut	$2lss + 0/1\tau_h$
nLeptons (tight)	2
Lepton p_T	25/15 GeV (leading/subleading lepton)
Product of lepton charges	>0
nJets	≥ 1 AK8 or ≥ 2 AK4
nTau (tight)	≤ 1
tauID (vs Jets)	deepVsjMedium
low mass veto	$m_{ll}^{loose} > 12$ GeV
bjet veto	$N_{bjet}^{loose} < 2 \ \&\& \ N_{bjet}^{medium} == 0$
$E_T^{miss} LD$	> 30/0 GeV for events without/with tight ID muons
VBF jet p_T	> 20 GeV
VBF jet $ \eta $	<4.7

5.7. BDT discriminants

In the qqHH focused search of the $HH \rightarrow$ multilepton analysis two BDTs are trained for the $2lss + 0/1\tau_h$ channel. The general training procedure is analogous to the BDT training in the original CMS $HH \rightarrow$ multilepton analysis. The dataset used to train the BDTs consists of simulated events only and is produced using a looser event selection compared to the one used in the SR, in order to increase the amount of events for the training. For this purpose, the lepton ID requirements are changed to the loose lepton definition and the DeepTau discriminator WP is relaxed to VVVLoose. Each simulated background event is replicated 13 times in the training sample, for the different EFT benchmark scenarios and the SM. This oversampling guarantees constant background statistics for a fixed benchmark point. Since the model is trained using looser object selection, the background events are reweighted according to their expected yields in the SR, making sure that all backgrounds are considered in the BDT training. To ensure the same importance of signal and background events in the BDT training, the sum of all signal events and all background events are both normalized to 10^5 events. The dataset is split into two samples of equal size, based on even and odd events, to guarantee that the model is able to generalize well and to avoid potential biases. Two separate sub-BDTs are trained, the one trained on even events is evaluated on odd events, and vice versa, making certain that the BDT is not trained and evaluated on the same events. XGBoost algorithm is used for the training.

The first BDT is aimed at separating ggHH-like events from qqHH-like events, and will be referred to as VBF/ggF BDT. The second BDT separates HH signal events from background events, and will be referred to as s/bkg BDT. Both BDTs are optimized by the following steps:

- Physics inspired variables are constructed in the analysis code and chosen as possible BDT input variables
- In order to reduce overtraining and make the model generalize better, the optimization of input variables is done. A selection of the provided variables is made, starting with removing one variable from a pair that has a correlation of 80% or more. After the highly correlated variables are removed, an iterative process starts in which a new model is

trained and the least performing variables are removed. The process continues until 10 variables are left.

The training of the VBF/ggF BDT is analogous to the s/bkg BDT, with qqHH samples chosen as background. Several samples with different coupling values ($C_V, C_{2V}, \kappa_\lambda$) including SM (1,1,1) and BSM samples (0.5,1,1), (1,0,1), (1,1,0), (1,2,1), (1.5,1,1) are used for the training.

5.7.1. BDT input variables and performance

To separate the ggHH signal from the qqHH signal, and signal events from background, numerous physics inspired variables are studied. The dominant signal contribution in this channel comes from $HH \rightarrow WWWW$ decays resulting in two same-sign leptons, each originating from separate Higgs bosons. The $HH \rightarrow WW\tau\tau$ decay also contributes to this channel to a smaller extent. Both Higgs bosons are identified in a similar way. The leading and sub-leading same sign leptons are paired with the reconstructed W -jets through differences in angular separation ΔR (a lepton is paired with a W -jet pair that has the smallest ΔR with the corresponding lepton). If the event has a τ_h , one of the Higgs bosons is reconstructed with a lepton/ τ_h pair, which is again matched with having the smallest ΔR , and the other Higgs boson with the remaining lepton and W -jet pair. With the focus on qqHH production two VBF jets are expected, which are chosen as described in section 5.6. VBF jets are characterized as high p_T jets, meaning that the invariant mass of the VBF jet pair should be relatively large in comparison with central jets. Several possible VBF jet selections are tested using a generator study:

- The jet pair with the highest invariant mass from a collection of high p_T jets is chosen as the VBF jet pair.
- The VBF jet candidates, chosen from the collection of high p_T jets, are cleaned using differences in angular separation with other jets (if ΔR is zero, then the potential VBF jet candidate is already used for other purposes), and the jet pair with the highest invariant mass is chosen as the VBF jet pair.
- The VBF jet candidates from the collection of high p_T jets are cleaned, and are required to pass selection criteria $|\eta| < 4.7, p_T > 20$. The jet pair with the highest invariant mass is chosen as the VBF jet pair.

In the end the third method is chosen, as this VBF jet selection is more reliable than the other two.

For training the VBF/ggF BDT, the invariant masses of di-lepton pairs m_{ll} and VBF jet pairs m_{jj}^{VBF} are used as input variables. In addition to this, angular separation between different reconstructed objects, such as the VBF jet pair ΔR_{jj}^{VBF} , HH pair $\Delta R_{H_1 H_2}$, and lepton VBF jet pair $\Delta R_{lj}^{Max_{VBF}}$ are used. The number of jets, the di-Higgs mass with E_T^{miss} , and the p_T of the sub-leading VBF jet are also included. Finally the scalar p_T sum of all reconstructed final state objects, denoted as HT is also used, since the non-resonant case offers a wide di-Higgs mass spectrum resulting in a relatively high energy content of signal events. The

main inputs to the VBF/ggF BDT classifier are summarized in Table 9.

Table 9. Input variables for $2lss + 0/1\tau_h$ VBF/ggF BDT

Variable	Description
vbv_m_jj	Invariant mass of the VBF jet pair.
HT	Scalar p_T sum of all reconstructed final state objects
maxJetPt_vbf	Maximum VBF jet p_T .
diHiggsMass_wMet_sel	Invariant mass of the ll pair and W-jets with E_T^{miss} .
nJet	Number of jets.
m_ll	Invariant mass of the ll pair.
maxdR_vbfjet_lep	Maximum opening angle between VBF jet and lepton.
vbv_pt_sublead	sub-leading VBF jet p_T .
dR_h1h2	Opening angle between reconstructed Higgs bosons.
vbv_dR_jj	Opening angle between the VBF jets.

The s/bkg BDT also uses the invariant mass of di-lepton pairs m_{ll} , the di-Higgs mass with E_T^{miss} , and number of jets as input variables. In addition, the number of VBF jets, the scalar p_T sum of all reconstructed final state objects with E_T^{miss} , denoted as $STMET$, and angular separations between reconstructed W-jets and leptons are used. The input variables used can be seen in Table 10.

Table 10. Input variables for $2lss + 0/1\tau_h$ s/bkg BDT

Variable	Description
dR_l_Wjets_min	Minimum opening angle between lepton and W-jets.
mindr_lep1_jet	Minimum opening angle between leading lepton and jet.
nJet_vbf	Number of VBF jets.
STMET	Scalar p_T sum of all final state objects and E_T^{miss} .
mindr_lep2_jet	Minimum opening angle between sub-leading l and jet.
dihiggsMass_wMet_sel	Invariant mass of the ll pair and W-jets with E_T^{miss} .
dR_l_leadWjet_min	Minimum opening angle between l and the highest p_T W-jets.
m_ll	Invariant mass of the di-lepton pair.
nJet	Number of jets.
dR_2j_fromW1	Opening angle of the $W \rightarrow jj$ jet pair.

In the case of both BDTs, the optimization of input variables omitted variables containing τ_h , such as the invariant mass of the τ_h and l pair $m_{l\tau_h}$, the angular separation between the τ_h and l pair $\Delta R_{l\tau_h}$. Other variables of interest that were studied but not chosen by the optimization include the differences in η and ϕ for the vbf jet pair ($\Delta\eta_{jj}^{VBF}$, $\Delta\phi_{jj}^{VBF}$), angular separations between the reconstructed Higgs bosons and the VBF jets (ΔR_{Hj}^{VBF}), between the VBF jet pair and lepton pair (ΔR_{lljj}^{VBF}). The hyperparameters used in the BDT training can be seen in Table 11. The hyperparameter values chosen for the s/bkg BDT come from the original HH multilepton analysis, and the VBF/ggF values are chosen arbitrarily and tweaked accordingly until the effects of overtraining are minimal. Distributions in some of the observables used as inputs to the BDT classifiers are shown in Figure 20

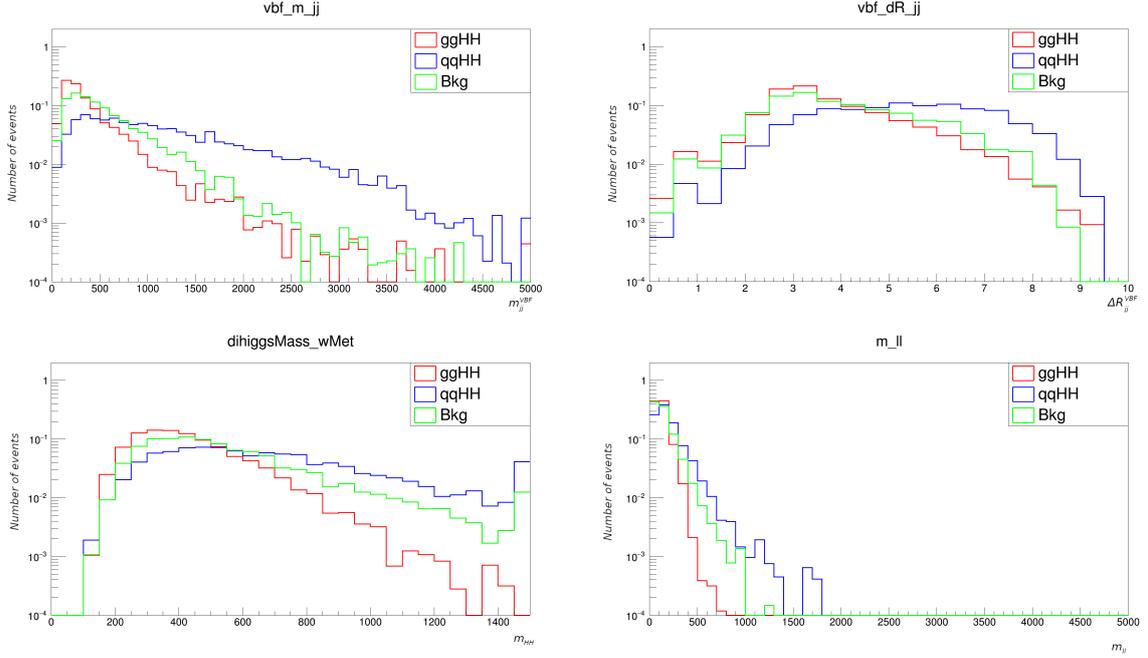


Figure 20. Distributions of m_{jj}^{VBF} , ΔR_{jj}^{VBF} , diHiggsMass_wMet_sel, and m_{ll} . Blue lines correspond to qqHH events, red to ggHH and green to background events.

Table 11. Hyperparameters of the $2lss + 0/1\tau_h$ BDT trainings

Parameter	VBF/ggF BDT	s/bkg BDT
n_estimators	150	496
subsample	0.7	0.910
colsample_bytree	0.358	0.416
gamma	4.734	3.295
learning_rate	0.2	0.174
max_depth	2	2
min_child_weight	62.426	485.384

In the VBF/ggF BDT training ggHH events are chosen as signal and qqHH events as background, labelled 1 and 0 respectively. The training is done on 5340257 signal events and 106414 background events, and achieves a good performance, as can be seen in Figure 21 showing the receiver operating characteristic (ROC) curve for the BDT output. The ROC curve is a graph showing the performance of the classification model. It is constructed by plotting the true positive rate, which is the proportion of observations classified correctly against the false positive rate, which is the proportion of observations that are classified incorrectly. ROC curves that are closer to the top left corner indicate classifiers with better performance. The amount of overtraining is kept at an acceptable level by the chosen hyperparameters. The BDT output for EFT scenarios, showing separation of ggHH and qqHH, and agreement between training and test events, can be seen in figure Figure 22.

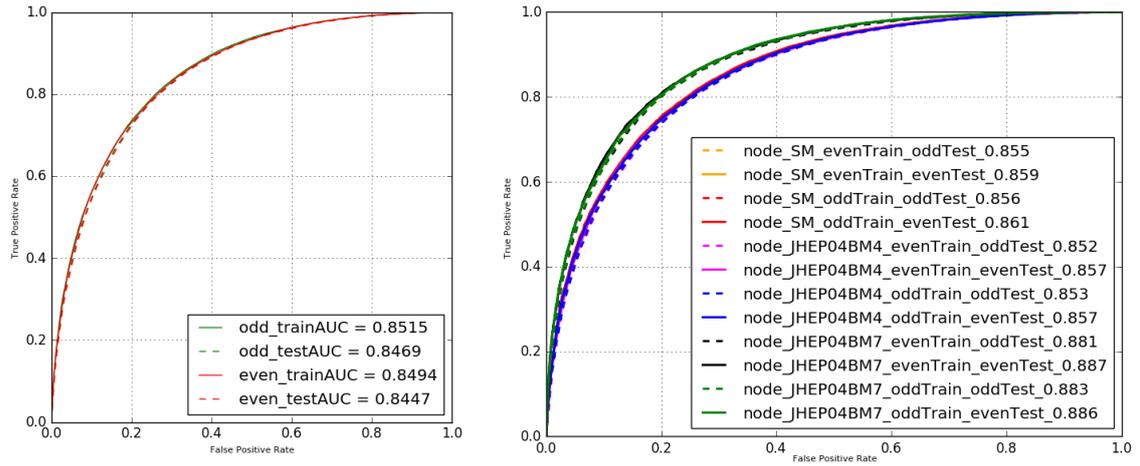


Figure 21. ROC curves for VBF/ggF BDT, showing the classification ability of the trained BDT on the training dataset. The training dataset is split into even and odd numbered halves. Two sub-BDTs are trained on the two halves and tested on the corresponding other half. The curve on the right shows the separation potential for the EFT scenarios (SM, BM4, BM7).

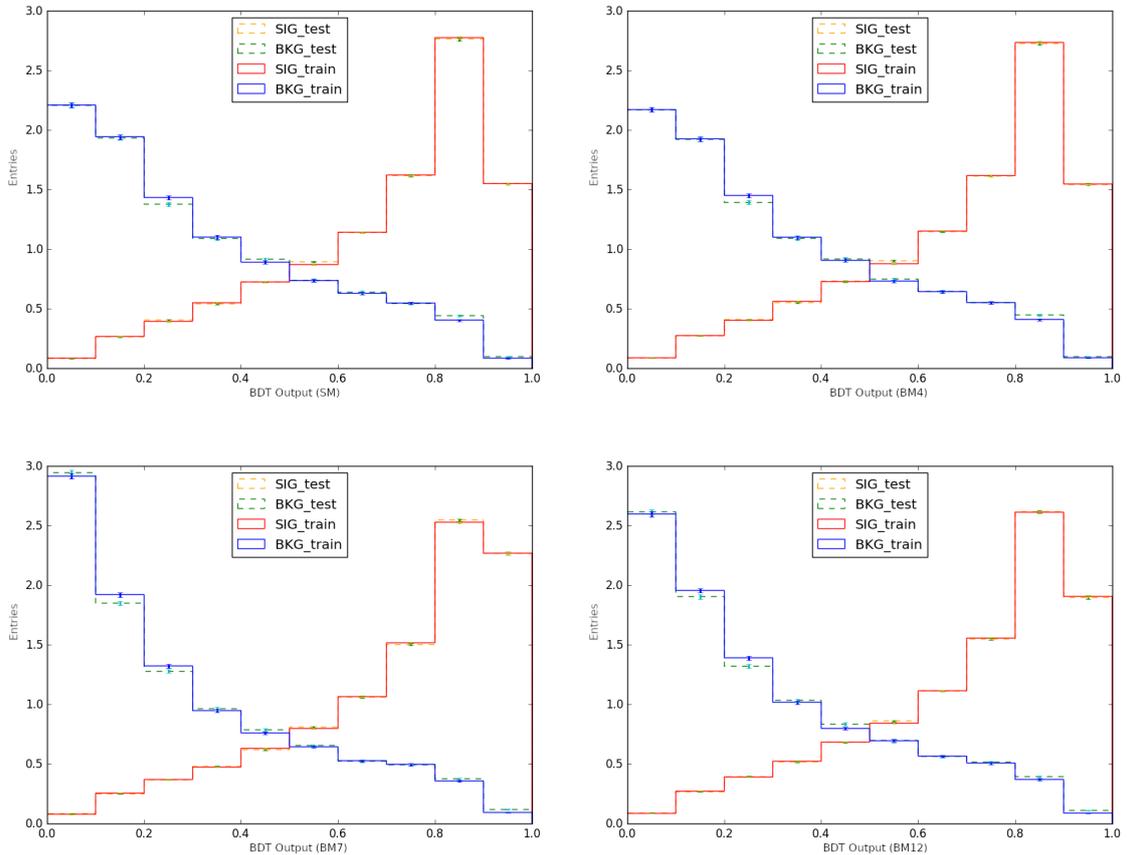


Figure 22. BDT outputs on the BDT training dataset, showing the separation of signal and background events (ggHH is drawn with a red line and denoted as SIG, qqHH is drawn with a blue line and is denoted as BKG). Two sub-BDTs are trained on even and odd numbered halves of the training dataset, and tested on the corresponding other half. The training events are drawn with a solid line, test events with a dashed line, allowing to see the agreement between them. The chosen EFT scenarios are SM, BM4, BM7 and BM12.

For the s/bkg BDT training ggHH events are chosen as signal and all other events as

background, again labelled 1 and 0 respectively. The training is done on 5345899 signal events and 486124 background events, and similar to the VBF/ggF BDT, reaches a good performance, as can be seen in Figure 23 showing the ROC curve for the BDT output. The BDT output for EFT scenarios, showing separation of signal and background, and agreement between training and test events, can be seen in figure Figure 24.

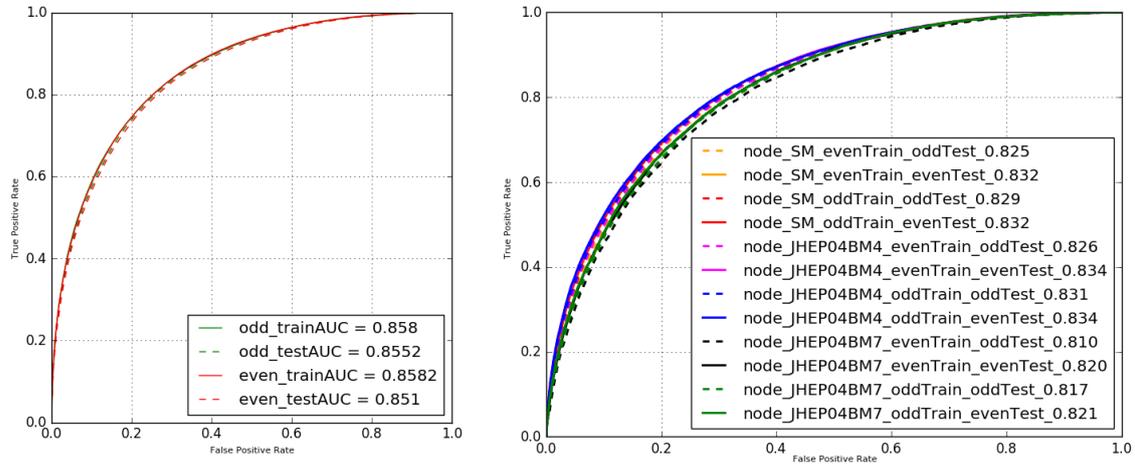


Figure 23. ROC curves for s/bkg BDT, showing the classification ability of the trained BDT on the training dataset. The training dataset is split into even and odd numbered halves. Two sub-BDTs are trained on the two halves and tested on the corresponding other half. The curve on the right shows the separation potential for the EFT scenarios (SM, BM4, BM7).

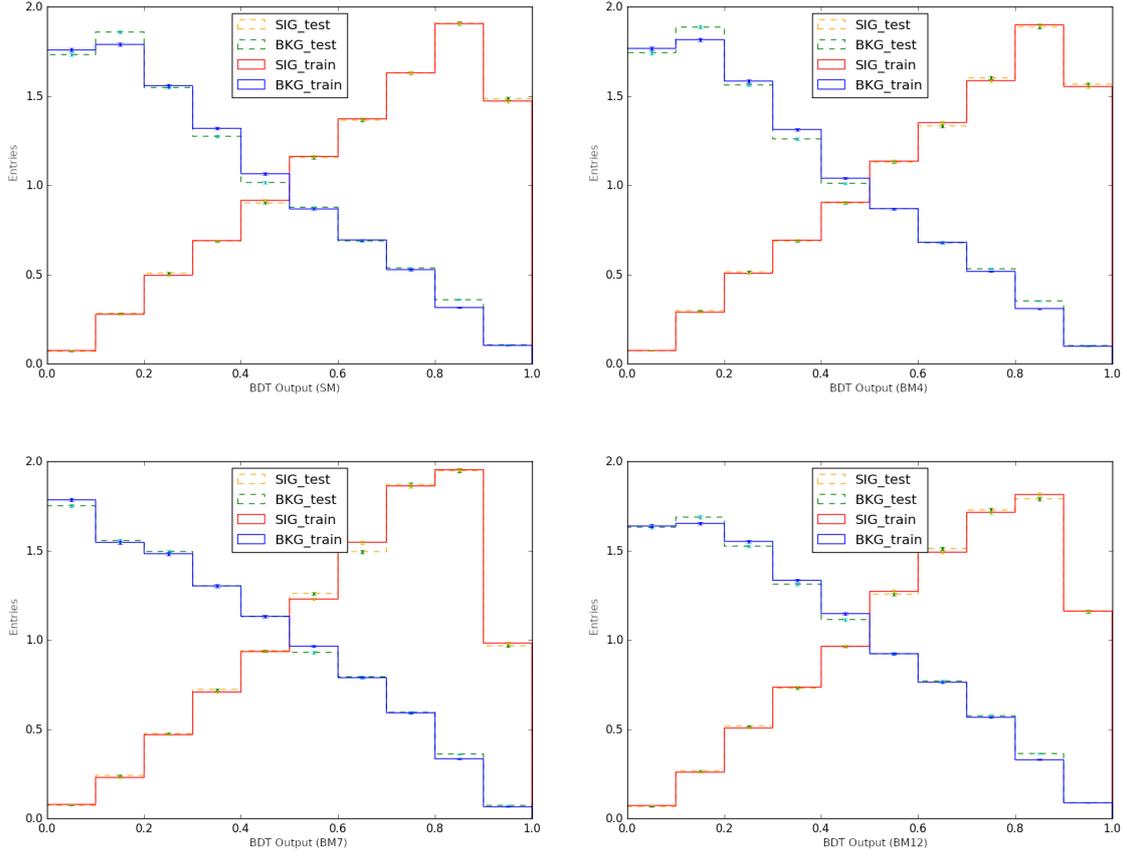


Figure 24. BDT outputs on the BDT training dataset, showing the separation of signal and background events (signal is drawn with a red line and denoted as SIG, background is drawn with a blue line and is denoted as BKG). Two sub-BDTs are trained on even and odd numbered halves of the training dataset, and tested on the corresponding other half. The training events are drawn with a solid line, test events with a dashed line, allowing to see the agreement between them. The chosen EFT scenarios are SM, BM4, BM7 and BM12.

5.7.2. 2D histograms

The two BDTs allow the use of a 2-D MVA distribution in the signal extraction. Figure 25 shows a schematic of the 2-D MVA distribution, where the output of the s/bkg BDT is on the x -axis, with 0 signifying background events and 1 signal events, and the VBF/ggF BDT is on the y -axis, with 0 indicating qqHH events (VBF) and 1 ggHH events (ggF). Here we can see a possible classification of events for a qqHH sample, where most of the events are gathered around the point (1;0), meaning that the events in the sample are classified as qqHH signal and the performance of the classifiers is excellent.

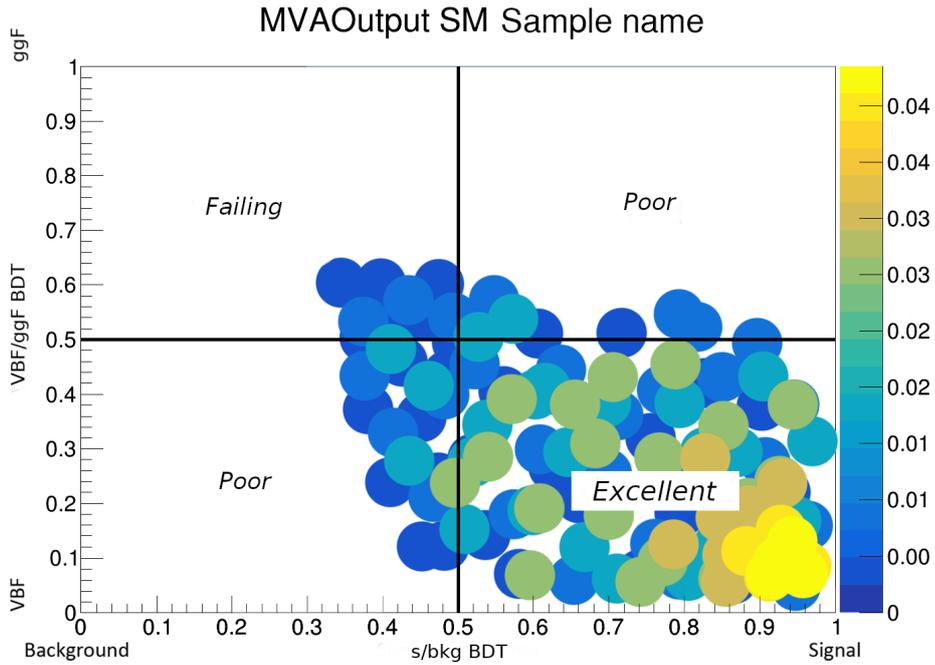


Figure 25. 2-D MVA distribution schematic for a qqHH sample, showing how the classification of the corresponding events should be if the two BDTs perform well. The canvas is split into 4 sections with coordinates (s/bkg; VBF/ggF). If in the case of a qqHH sample, the points are gathered around (0;1), the performance of both BDTs would be failing. Events gathered around (1;1) or (0;0) would mean that one BDT is performing well but the other one fails, making the combined performance of the BDTs poor.

The 2-D MVA distributions for the EFT SM scenario of the following 2018 era samples:

- SM VBF (SM-like qqHH production)
- BSM VBF (BSM-like qqHH production)
- SM ggF (SM-like ggHH production)
- SM WZTo3LNu (Background events from WZ to 3 leptons with neutrinos)
- SM ZZTo4L (Background events from ZZ to 4 leptons)

can be seen in Figure 26 (signal samples) and Figure 27 (background samples). The 2-D distributions for the 2017 and 2016 era samples look similar to 2018. For the signal extraction the 2-D distribution is unrolled into a 1-D distribution.

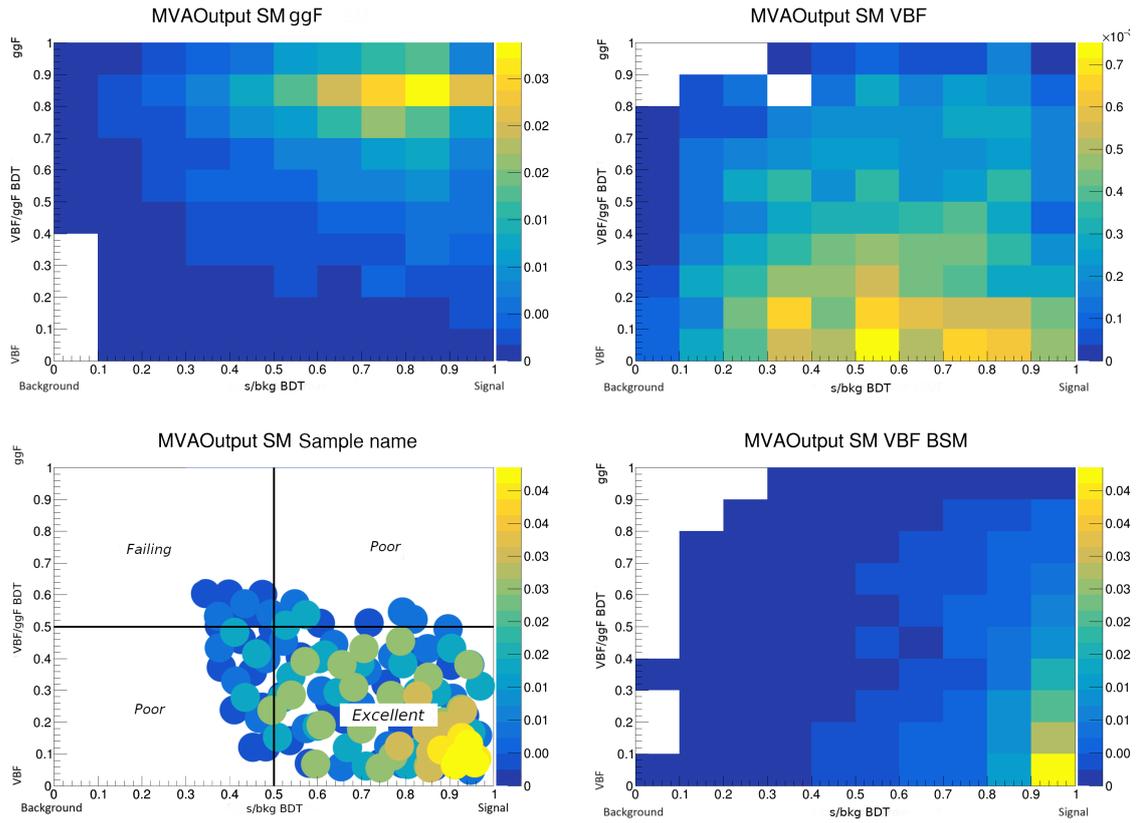


Figure 26. 2-D MVA distributions for the EFT SM scenario of the 2018 era signal events (SM ggF, SM VBF, BSM VBF) in comparison with the potential qqHH signal sample distribution. Here we can see how the BDTs classify different HH signal events. The top left plot has events gathered around (1;1), meaning that the events are classified as ggF signal events, which matches the sample, indicating that the BDTs are performing well. The same can be said about the bottom right BSM VBF sample plot, where the events are classified excellently as qqHH signal events. The top right SM VBF sample plot shows that the VBF/ggF BDT is able to classify the events as qqHH, but the signal and background separation, while the events being more on the signal side, is not perfect.

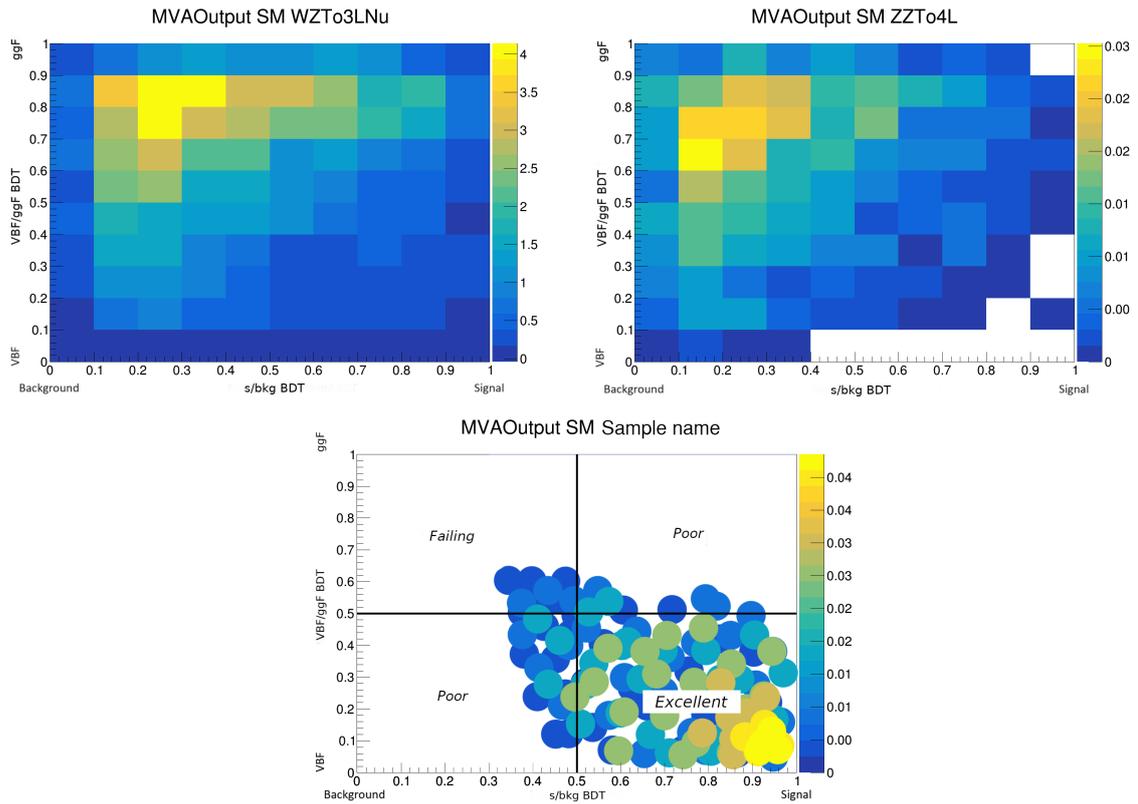


Figure 27. 2-D MVA distributions for the EFT SM scenario of the 2018 era background events (WZ to 3 leptons and neutrinos, ZZ to 4 leptons) in comparison with the potential qqHH signal sample distribution. Both plots show the classification of background events being classified well.

6. Results

6.1. Systematic uncertainties

The predicted event yields and distributions in the output of BDT classifiers are affected by various systematic uncertainties, which may be theoretical, affecting the predicted cross section or the kinematics of the collision process of decays, or experimental, coming from differences in object reconstruction or uncertainties on the estimation of the data driven lepton misidentification (fakes) and electron charge misidentification (flips) backgrounds. The uncertainties and their effects are treated as nuisance parameters in the signal extraction described in section 6.2. Theoretical uncertainties on signal and background rates, experimental uncertainties on trigger efficiency, B-tagging efficiency, luminosity and pileup are included in this analysis. The uncertainties are listed in Table 12 and are discussed more in detail in [3] and [43].

Table 12. Systematic uncertainties affecting the multilepton analysis

Uncertainty
Trigger efficiency
Electron, muon and τ_h reconstruction
Jet energy scale and resolution
B-tag efficiency and mistag rate
Luminosity
Pileup
L1 ECAL prefiring
Theory cross section
Data driven background estimation

In the case of the $2lss + 0/1\tau_h$ channel, results are mainly limited by the statistical uncertainty of the data with only partial contribution of the systematic uncertainties on data driven fake background estimation and other systematics to an even smaller extent. The fake background estimation uncertainties are associated with shape variations emerging from statistical uncertainties in the measurement region (MR) and the application region (AR). The statistical uncertainties can emerge from differences in background composition between MR (dominated by multijet background) and AR (dominated by WZ). The differences in background composition are determined by comparing the prediction for the fake background in the SR with an estimate of the fake background. The fake background prediction is obtained by MC simulation, where all events pass signal selection criteria, and is referred to as "nominal MC" shape templates, while the fake background estimate is acquired by scaling the simulated events selected for the AR according to the FF method, and is referred to as "MC closure" shape templates. A linear multivariable function $f(x, y)$ is used to fit the ratios of the nominal MC and MC closure shape templates, and the deviation of the slope from 0 is taken as an additional uncertainty, that is used in the BDT output shape. Examples can be seen in

Figure 28 where the data fakes shape uncertainties for electrons and muons is shown.

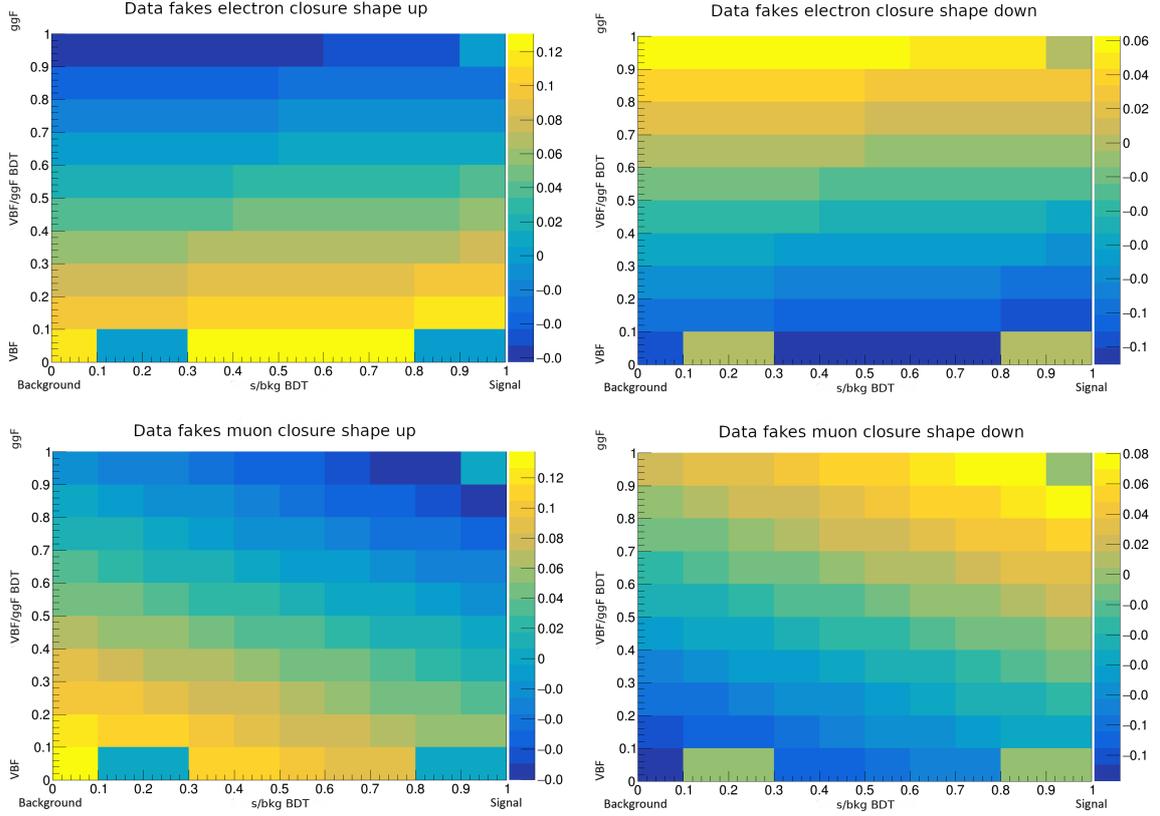


Figure 28. Up and down shape uncertainties for data fakes electrons and muons

The effect of statistical uncertainties in the measurement for the fake probabilities used for the FF method is evaluated by varying the probabilities within their uncertainties and using the varied probabilities for the calculation of the fake factors, thus resulting in shape uncertainties in the BDT output distribution. In addition, a 30% normalization uncertainty on the fake background estimation is added and is uncorrelated across all data taking eras. The uncertainty on the charge flip background in $2lss + 0/1\tau_h$ amounts to 30% and is correlated among all data taking eras. [3]

6.2. Signal extraction

For the signal extraction a binned maximum likelihood fit to BDT output distributions is performed, measuring the signal strength modifier defined as $\mu = \frac{\sigma_{qqHH}}{\sigma_{qqHH}^{theo}}$ with a profile likelihood test statistic [68]. The BDTs as in the original analysis are parametrized by a number of EFT benchmarks, but for simplicity the SM output node is used in this analysis. The null hypothesis (H_0) is the existence of a signal with strength μ (signal-plus-background) and the alternative hypothesis (H_1) is no signal (background-only). The aim is to exclude the null hypothesis above a certain μ . The likelihood test statistic is defined as:

$$\tilde{q}_\mu = -2 \ln \frac{\mathcal{L}(\text{data}|\mu = 0, \hat{\theta}_\mu)}{\mathcal{L}(\text{data}|\mu = \hat{\mu}, \hat{\theta})}, \quad (6.1)$$

where \mathcal{L} is the likelihood function, defined as:

$$\mathcal{L}(\text{data}|\mu, \theta) = \text{Poisson}(\text{data}|\mu \cdot s(\theta) + b(\theta)) \cdot p(\tilde{\theta}|\theta). \quad (6.2)$$

Here data represents actual experimental observation or pseudo-data, signal is denoted as s and background is denoted as b . Predictions for signal and background distributions are dependent on both systematic and statistical uncertainties, that are managed by nuisance parameters θ , such that the signal and background expectations become functions of the nuisance parameters. The distribution $p(\tilde{\theta}|\theta)$ characterizes nuisance parameters with default values $\tilde{\theta}$. In equation 6.1 $\hat{\theta}_\mu$ is the conditional maximum likelihood estimator of θ given μ , and $\hat{\mu}$ and $\hat{\theta}$ are parameter estimators that correspond to the global maximum of the likelihood. The test statistic has a lower constraint $0 \leq \hat{\mu}$, considering only models with positive signal rate, and an upper constraint $\hat{\mu} \leq \mu$ that guarantees a one-sided confidence interval. Upward fluctuations of data, such that $\hat{\mu} \geq \mu$ are not used as evidence against H_0 . The test statistic is used to calculate p-values for both hypotheses (p_μ for H_0 and p_b for H_1), which are used to calculate the 95% confidence level upper limit on μ by adjusting μ until

$$\text{CL}_s(\mu) = \frac{p_\mu}{1 - p_b} < 0.05. \quad (6.3)$$

The method described above is used to calculate observed limits. To calculate expected limits, a large set of background only pseudo-data is generated and 95% confidence level upper limits on μ are calculated for each of them, as if they were real data. A cumulative probability distribution of results is built by starting integration from the side corresponding to low event yields as can be seen in Figure 29 (right). The median expected value is the point where the cumulative probability distribution crosses the 0.5 quantile. The $\pm 1\sigma$ (68%) band is defined by the crossings of the 0.16 and 0.84 quantiles, and the $\pm 2\sigma$ (95%) band by crossings at 0.025 and 0.975 quantiles. [69] More about the binned likelihood fit and calculating upper limits can be found in [68, 69].

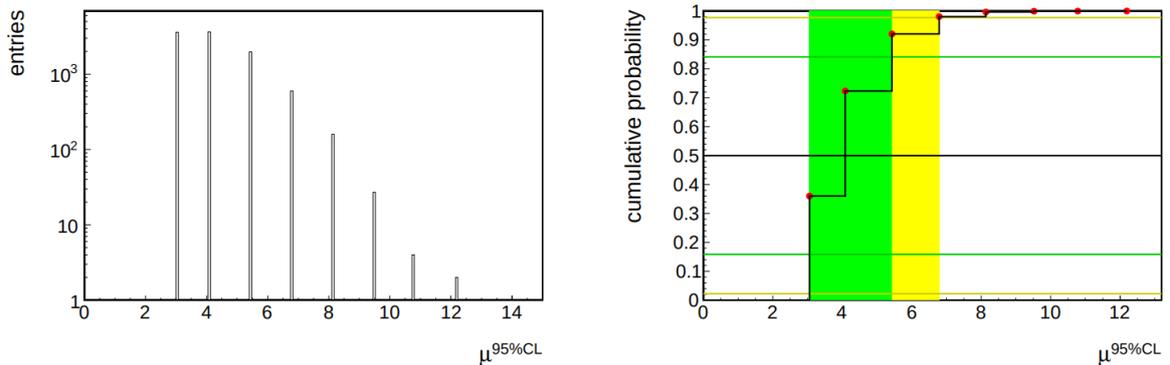


Figure 29. The plot on the left shows an example of a differential distribution of possible limits on μ for the background only hypothesis ($s = 1, b = 1$, no systematic errors). The plot on the right shows the cumulative probability distribution of the plot on the left with 0.025, 0.16, 0.5, 0.84, and 0.975 quantiles, denoted as horizontal lines, which define the median expected limit and the 68% and 95% bands for the expected value of μ [69].

All limits calculated in this analysis are extracted using the Higgs Combine toolkit [70]. Limits

on the signal strength modifier μ are used to calculate limits on signal cross section and coupling strength modifiers. Signal cross section limits are calculated by multiplying the upper limit on the signal strength modifier μ with the signal cross section the input signal model is normalized to, that being σ^{theo} or σ_{SM} . Limits on coupling strength modifiers such as the Higgs vector boson coupling C_{2V} are calculated by scanning the upper limit on μ as a function of the given coupling parameter, where the intersection of the limit on μ with a line at one gives the upper limit on the coupling strength modifier by excluding scenarios, where $\mu < 1$. [43] Upper limits on the signal strength modifier μ , coupling strength modifier C_{2V} , and signal cross section for qqHH production can be seen in the section that follows.

6.3. Results on qqHH production

Figure 30 shows the 95% confidence level (CL) upper limit on SM like ($C_{2V} = 1$) non-resonant qqHH production in the $2lss + 0/1\tau_h$ channel. The expected limit on SM qqHH production rate is $908 \times \sigma_{qqHH}^{SM}$. The comparison of the limits between the original CMS HH→multilepton analysis and the qqHH focused search shows the results on SM like qqHH production improving by a factor of 3.3.

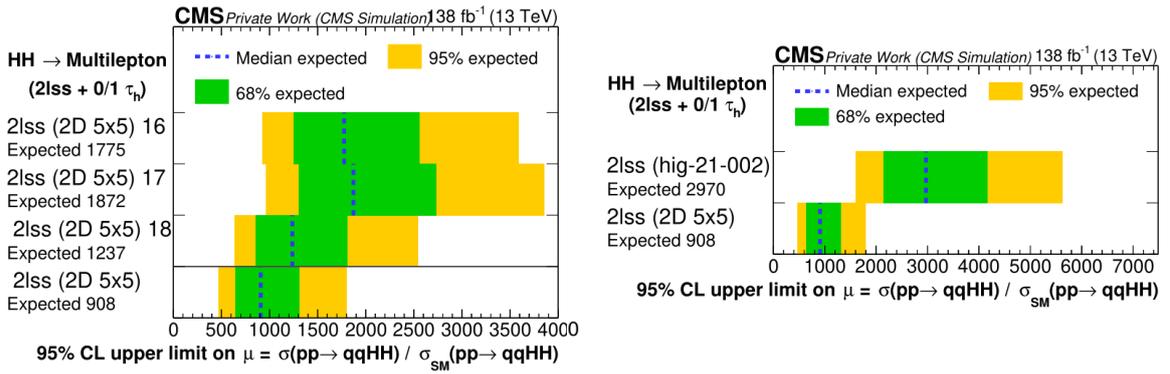


Figure 30. Limits on SM like qqHH production in the $2lss + 0/1\tau_h$ channel. The plot on the left shows the 95% confidence level (CL) upper limit on the signal strength modifier $\mu = \sigma(pp \rightarrow qqHH) / \sigma_{SM}(pp \rightarrow qqHH)$ for eras 2016, 2017 and 2018 as well as their combination. The plot on the right shows the comparison of the upper limits between the original HH→multilepton (denoted as hig-21-002) and the qqHH focused search (denoted as 2D 5x5).

The choice of the 5×5 2-D binning can be explained with Figure 31, which shows a comparison of different binnings (5×5 , 5×3 , 3×5 , and 3×3) of both BDT outputs with equal sized bins. A consistent background description is required, meaning that all backgrounds have entries in all bins, therefore finer binnings than 5×5 are not tested. The 5×5 2-D binning gives the best results from the different binnings satisfying this requirement. This rather simple binning choice offers room for further improvement in future iterations of the analysis.

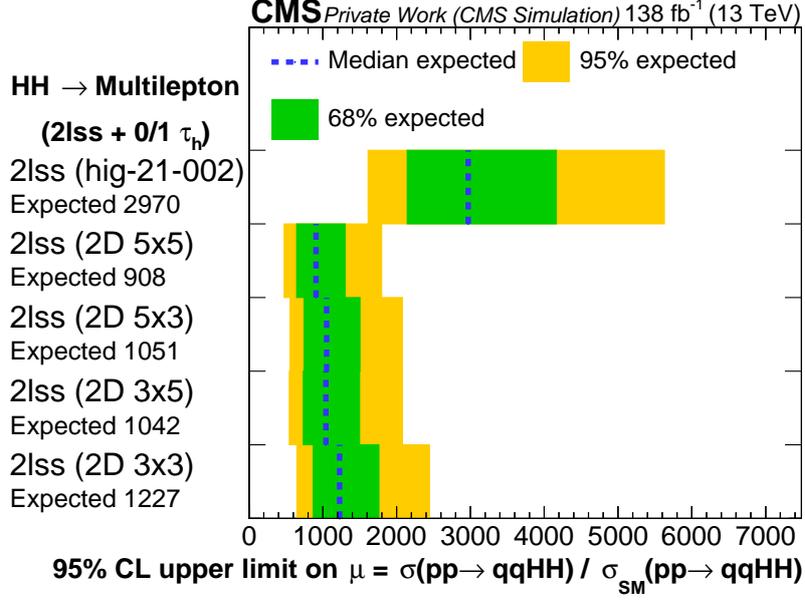


Figure 31. The 95% confidence level (CL) upper limit on the signal strength modifier $\mu = \sigma(\text{pp} \rightarrow \text{qqHH}) / \sigma_{SM}(\text{pp} \rightarrow \text{qqHH})$ for different 2D binnings.

The 95% CL upper limit on the coupling strength modifier for non-resonant qqHH production $\mu_{qqHH} = \sigma_{qqHH} / \sigma_{qqHH}^{theo}$ for different values of C_{2V} and for the cross section as a function of C_{2V} can be seen in Figure 32. The intersection of the limit on the coupling strength modifier with the red line at one corresponds to the expected limit on C_{2V} of $-2.66 < C_{2V} < 4.78$. For the cross section limit, theoretical uncertainties on σ_{theo} are frozen.

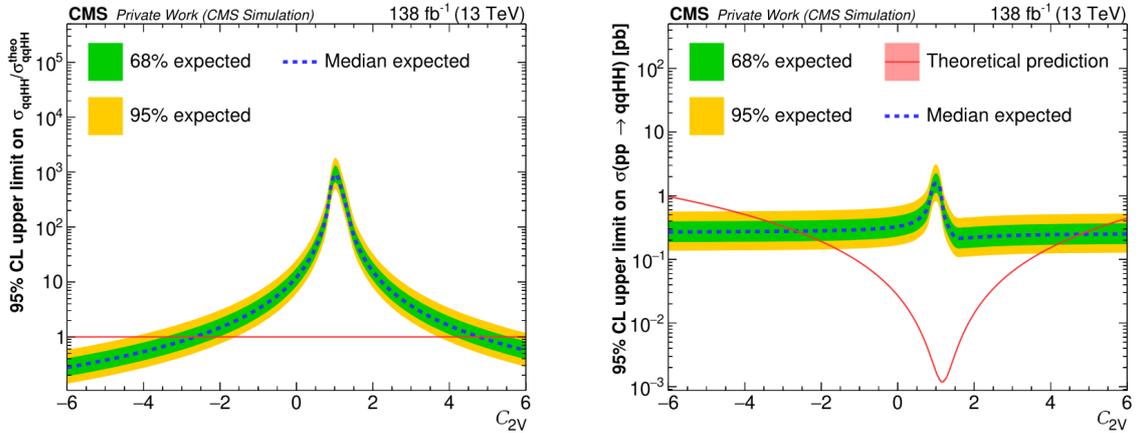


Figure 32. The plot on the left shows the 95% CL upper limit on the coupling strength modifier for non-resonant qqHH production in the 2lss + 0/1τ_h channel. The plot on the right shows the limit on the qqHH production cross section as a function of C_{2V} .

In the SM the qqHH production cross section depends on both the Higgs boson coupling to vector bosons C_{2V} and the coupling between the Higgs boson with a single vector boson C_V . A two dimensional likelihood profile in the $C_V - C_{2V}$ space can be seen in Figure 33, showing what regions of values for both couplings could be excluded.

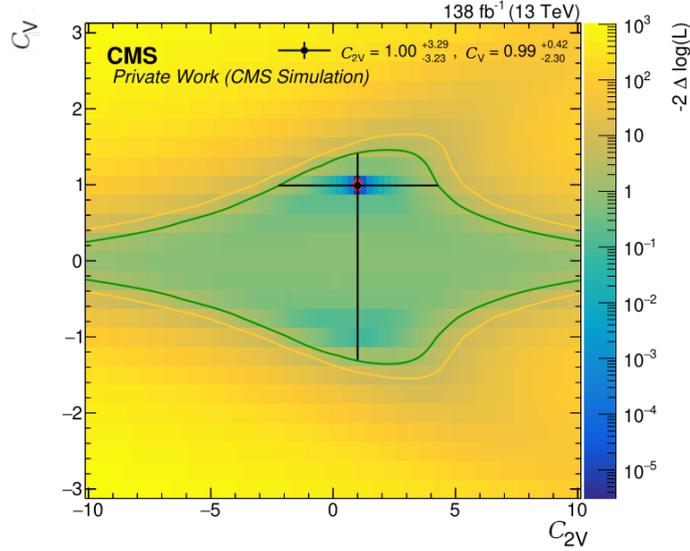


Figure 33. 2-D likelihood profile in the $C_V - C_{2V}$ space

Figure 34 shows the 95% upper limit on non-resonant BSM ($C_{2V} = 0$) qqHH production in the $2lss + 0/1\tau_h$ channel in comparison to the original $HH \rightarrow$ multilepton analysis. The expected limit on the BSM qqHH production rate is $12 \times \sigma_{theo}$, improving the result by a factor of 2.4. This is currently an interesting point, with the CMS collaboration being close to excluding this BSM coupling as can be read from [71].

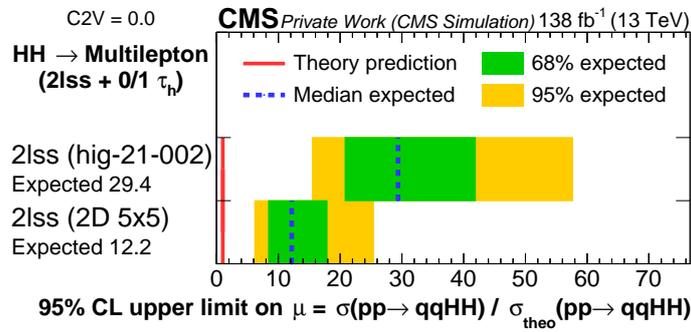


Figure 34. Limit on BSM like qqHH production in the $2lss + 0/1\tau_h$ channel, showing the 95% confidence level (CL) upper limit on the signal strength modifier $\mu = \sigma(pp \rightarrow qqHH) / \sigma_{theo}(pp \rightarrow qqHH)$. $HH \rightarrow$ multilepton is denoted as hig-21-002 and the qqHH focused search as 2D 5x5.

A blinded prefit BDT distribution in the $2lss + 0/1\tau_h$ channel can be seen in Figure 35 with the main backgrounds and the visualization of the separation of the qqHH and ggHH signals. The prefit BDT distribution is a 1-D transformation of the 2-D input. All results seen in this section are without optimized binning and are extracted using the SM output node, meaning that the results can be improved further.

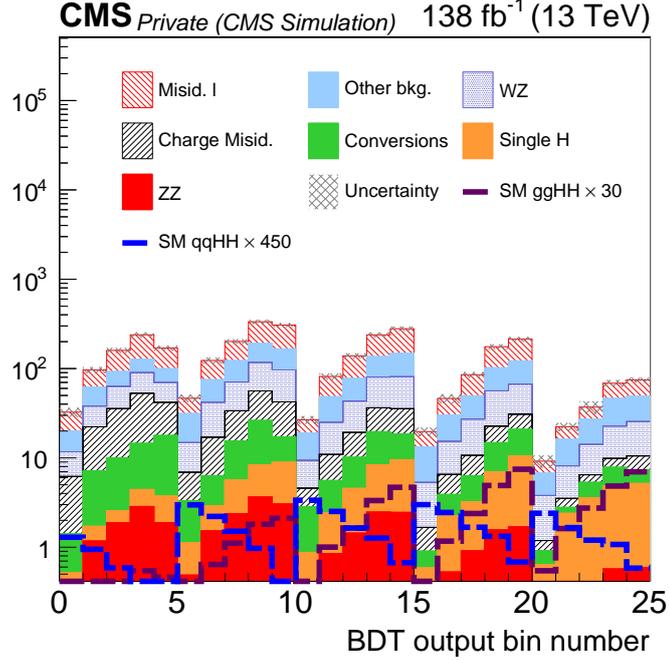


Figure 35. Prefit BDT distribution in the $2lss + 0/1\tau_h$ channel

6.4. Comparison to other analyses

The CMS $HH \rightarrow \text{multilepton}$ analysis is the first ever analysis searching for di-Higgs production in $HH \rightarrow \tau\tau\tau\tau$ and $HH \rightarrow WW\tau\tau$ decay modes, and the first analysis in CMS targeting the $HH \rightarrow WWWW$ decay mode. Figure 36 shows the comparison of the 95% upper limits on qqHH production cross section as a function of C_{2V} between the CMS $HH \rightarrow \text{multilepton}$ analysis, the extension of the analysis to the qqHH production mode, and other leading di-Higgs analyses ($HH \rightarrow bbbb$, $HH \rightarrow bb\tau\tau$, $HH \rightarrow bb\gamma\gamma$). The addition of the qqHH focused search improves the results of the CMS $HH \rightarrow \text{multilepton}$ analysis for both SM ($C_{2V} = 1$) and BSM ($C_{2V} = 0$) cases. The improvements around SM are greater than around BSM, allowing the CMS $HH \rightarrow \text{multilepton}$ analysis to compete around SM with the other leading di-Higgs analyses.

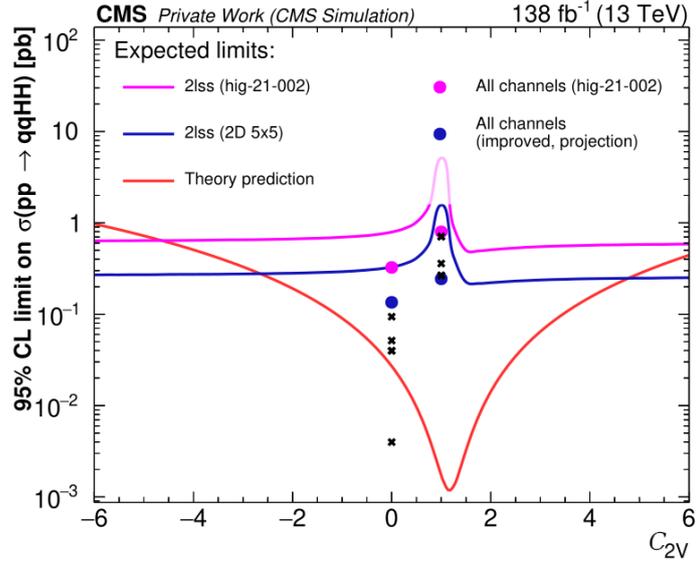


Figure 36. Comparison of upper limits at 95% CL for qqHH production cross section as a function of C_{2V} achieved for the CMS HH→multilepton analysis drawn in pink, the qqHH focused search drawn in blue, and HH→bbbb (CMS [72], ATLAS [73]), HH→bb $\tau\tau$ (ATLAS [74]), HH→bb $\gamma\gamma$ (CMS [75], ATLAS [76]) analyses denoted as black crosses.

Figure 37 shows the comparison of the 95% CL upper limit on the cross section for HH production as a function of the trilinear Higgs boson self coupling modifier κ_λ between the HH→multilepton analysis and the qqHH focused search in the $2lss + 0/1\tau_h$ channel. It can be seen that the addition of the qqHH focused search does not affect the previous results negatively, showing even a slight improvement.

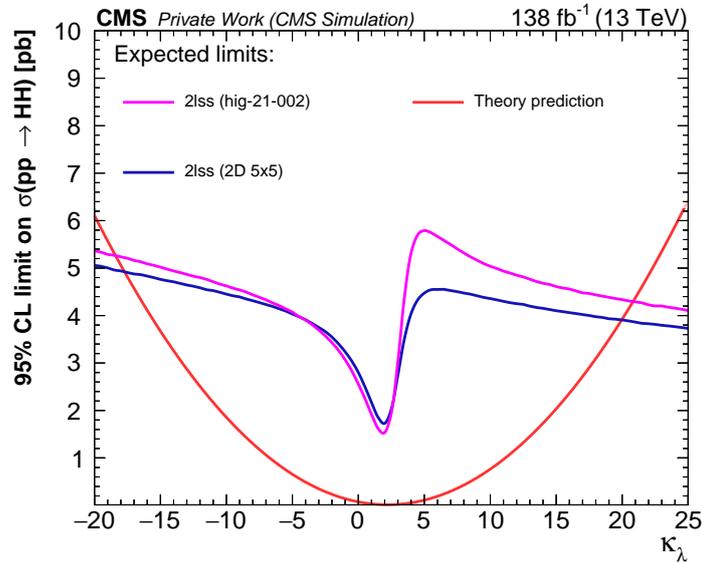


Figure 37. Comparison of upper limits at 95% CL for κ_λ achieved for the CMS HH→multilepton analysis (pink line) and the addition of the qqHH focused search (blue line) in the $2lss + 0/1\tau_h$ channel. The results from the CMS HH→multilepton analysis are not expected to diminish with the addition of the qqHH focused search.

7. Summary and conclusions

This thesis presented results on the feasibility and prospects of extending the CMS $HH \rightarrow$ multilepton analysis focus to the sub-dominant $qqHH$ production mode, using data collected in LHC Run 2 during the years 2016-2018 with center of mass energy of $\sqrt{s} = 13$ TeV and a total luminosity of 138 fb^{-1} . The CMS $HH \rightarrow$ multilepton analysis was performed in seven different search categories in final states with multiple electrons, muons, or hadronically decaying τ leptons. From the seven categories, the $2lss + 0/1\tau_h$ search category was chosen for the extension of the analysis due to having the biggest HH signal. In order to target $qqHH$ production, two BDTs were trained, one focusing on separating $qqHH$ -like events from $ggHH$ -like events (VBF/ggF BDT), and one for separating signal from background (s/bkg BDT). Both BDTs include new variables revolving around the additional jets (VBF jets) expected in $qqHH$. The VBF jets were chosen from a selection of reconstructed high p_T jets, that was cleaned of jets used for W boson reconstruction, and had to satisfy the selection criteria $|\eta| < 4.7$, $p_T > 20$. The jet pair with the biggest invariant mass passing the VBF jet selection was chosen as the VBF jet pair. The signal extraction was carried out using the output of the two BDTs. Expected 95% CL upper limits on SM-like ($C_{2V} = 1$) non-resonant $qqHH$ production in the $2lss + 0/1\tau_h$ channel were calculated, and the expected limit on the production rate was found to be $908 \times \sigma_{qqHH}^{SM}$, resulting in an improvement of results by a factor of 3.3 compared to the CMS $HH \rightarrow$ multilepton analysis. The expected limit on the coupling strength modifier C_{2V} was set at the interval $-2.66 < C_{2V} < 4.78$. Regarding BMS-like ($C_{2V} = 0$) $qqHH$ production in the $2lss + 0/1\tau_h$ channel, 95% CL upper limits on μ were set at $12 \times \sigma_{theo}$, improving the result by a factor of 2.4. When comparing the $qqHH$ focused search results for κ_λ with the CMS $HH \rightarrow$ multilepton analysis, it can be said that the results do not diminish but rather improve slightly. The addition of the $qqHH$ focused search results to the CMS $HH \rightarrow$ multilepton analysis allows for rivaling other contemporary HH analyses by both the ATLAS and CMS collaborations for SM ($C_{2V} = 1$) $qqHH$ production. The results presented show very strong improvements in the extraction of the $qqHH$ signal component, without limiting the potential for already existing results, and can be further improved with the optimization of the BDT binning and output node choice.

8. Acknowledgements

First, I would like to thank my co-supervisor Christian Veelken for proposing the topic of this thesis, allowing me to work with the analysis group at NICPB. I thank my supervisor Torben Lange for guiding me and making the experience of writing this thesis a very educating and fun one. Thanks also go to Laurits Tani, who helped me settle into the new environment and assisted with the machine learning part of this thesis. Finally, I thank all of my colleagues at NICPB for giving me valuable input at group meetings and being fun people to be around.

Bibliography

- [1] Serguei Chatrchyan et al. “Observation of a New Boson at a Mass of 125 GeV with the CMS Experiment at the LHC”. In: *Phys. Lett. B* 716 (2012), pp. 30–61. DOI: 10.1016/j.physletb.2012.08.021. arXiv: 1207.7235 [hep-ex].
- [2] Georges Aad et al. “Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC”. In: *Phys. Lett. B* 716 (2012), pp. 1–29. DOI: 10.1016/j.physletb.2012.08.020. arXiv: 1207.7214 [hep-ex].
- [3] CMS Collaboration. *Search for Higgs boson pairs decaying to $WWWW$, $WW\tau\tau$, and $\tau\tau\tau\tau$ in proton-proton collisions at $\sqrt{s} = 13$ TeV*. 2022. DOI: 10.48550/ARXIV.2206.10268. URL: <https://arxiv.org/abs/2206.10268>.
- [4] Fady Bishara, Roberto Contino, and Juan Rojo. “Higgs pair production in vector-boson fusion at the LHC and beyond”. In: *The European Physical Journal C* 77.7 (July 2017). DOI: 10.1140/epjc/s10052-017-5037-9. URL: <https://arxiv.org/pdf/1611.03860.pdf>.
- [5] R. Mann. *An Introduction to Particle Physics and the Standard Model*. CRC Press, 2009. ISBN: 9781420083002. URL: <https://books.google.ee/books?id=wSrNBQAAQBAJ>.
- [6] M. Robinson. *Symmetry and the Standard Model: Mathematics and Particle Physics*. Springer New York, 2014. ISBN: 9781489997777. URL: <https://books.google.ee/books?id=RCjWoQEACAAJ>.
- [7] Donald H. Perkins. *Introduction to High Energy Physics*. 4th ed. Cambridge University Press, 2000. DOI: 10.1017/CBO9780511809040.
- [8] Mark Thomson. *Modern Particle Physics*. Cambridge University Press, 2013. DOI: 10.1017/CBO9781139525367.
- [9] Wikimedia Commons. *Standard Model of Elementary Particle Physics*. [Online; accessed December 20, 2022]. 2017. URL: https://commons.wikimedia.org/wiki/File:Standard_Model_of_Elementary_Particles_modified_version.svg.
- [10] S. Dawson. *Introduction to Electroweak Symmetry Breaking*. 1999. DOI: 10.48550/ARXIV.HEP-PH/9901280. URL: <https://arxiv.org/abs/hep-ph/9901280>.
- [11] John Ellis. *Higgs Physics*. 2013. DOI: 10.48550/ARXIV.1312.5672. URL: <https://arxiv.org/abs/1312.5672>.
- [12] Maxime Gouzevitch and Alexandra Carvalho. “A review of Higgs boson pair production”. In: *Reviews in Physics* 5 (2020), p. 100039. ISSN: 2405-4283. DOI: <https://doi.org/10.1016/j.revip.2020.100039>. URL: <https://www.sciencedirect.com/science/article/pii/S2405428320300022>.

- [13] Biagio Di Micco et al. “Higgs boson potential at colliders: Status and perspectives”. In: *Reviews in Physics* 5 (Nov. 2020), p. 100045. DOI: 10.1016/j.revip.2020.100045. URL: <https://www.sciencedirect.com/science/article/pii/S2405428320300083>.
- [14] R. Frederix et al. “Higgs pair production at the LHC with NLO and parton-shower effects”. In: *Physics Letters B* 732 (May 2014), pp. 142–149. DOI: 10.1016/j.physletb.2014.03.026. URL: <https://doi.org/10.1016%5C%2Fj.physletb.2014.03.026>.
- [15] M. Cepeda et al. *Higgs Physics at the HL-LHC and HE-LHC*. 2019. DOI: 10.48550/ARXIV.1902.00134. URL: <https://arxiv.org/abs/1902.00134>.
- [16] M. Grazzini et al. “Higgs boson pair production at NNLO with top quark mass effects”. In: *Journal of High Energy Physics* 2018.5 (May 2018). DOI: 10.1007/jhep05(2018)059. URL: <https://doi.org/10.48550/arXiv.1803.02463>.
- [17] Frédéric A. Dreyer and Alexander Karlberg. “Vector-boson fusion Higgs pair production at N³LO”. In: *Physical Review D* 98.11 (Dec. 2018). DOI: 10.1103/physrevd.98.114016. URL: <https://doi.org/10.48550/arXiv.1811.07906>.
- [18] Lyndon Evans and Philip Bryant. “LHC Machine”. In: *Journal of Instrumentation* 3.08 (Aug. 2008), S08001. DOI: 10.1088/1748-0221/3/08/S08001. URL: <https://dx.doi.org/10.1088/1748-0221/3/08/S08001>.
- [19] the ATLAS Collaboration. “The ATLAS Experiment at the CERN Large Hadron Collider”. In: *JINST* 3 (2008), S08003. DOI: 10.1088/1748-0221/3/08/S08003.
- [20] the CMS Collaboration. “The CMS Experiment at the CERN LHC”. In: *JINST* 3 (2008), S08004. DOI: 10.1088/1748-0221/3/08/S08004.
- [21] The ALICE Collaboration. “The ALICE experiment at the CERN LHC”. In: *Journal of Instrumentation* 3.08 (Aug. 2008), S08002. DOI: 10.1088/1748-0221/3/08/S08002. URL: <https://dx.doi.org/10.1088/1748-0221/3/08/S08002>.
- [22] The LHCb Collaboration. “The LHCb Detector at the LHC”. In: *Journal of Instrumentation* 3.08 (Aug. 2008), S08005. DOI: 10.1088/1748-0221/3/08/S08005. URL: <https://dx.doi.org/10.1088/1748-0221/3/08/S08005>.
- [23] Esma Anais Mobs. “The CERN accelerator complex. Complexe des accélérateurs du CERN”. In: (2016). General Photo. URL: <https://cds.cern.ch/record/2225847>.
- [24] J. T. Boyd. *LHC Run-2 and Future Prospects*. 2020. DOI: 10.48550/ARXIV.2001.04370. URL: <https://arxiv.org/abs/2001.04370>.
- [25] Jacques Gareyte. “LHC main parameters”. In: *Part. Accel.* 50 (1995), pp. 61–68. URL: <https://cds.cern.ch/record/304825>.
- [26] David Barney and Sergio Cittolin. “CMS Detector Drawings”. In: (2000). URL: <https://cds.cern.ch/record/2629816>.

- [27] Izaak Neutelings. *CMS coordinate system with the a cylindrical detector*. [Online; accessed September 21, 2022]. 2021. URL: https://tikz.net/wp-content/uploads/2021/09/axis3D_CMS-004.png.
- [28] Izaak Neutelings. *CMS coordinate system with the a cylindrical detector*. [Online; accessed September 21, 2022]. 2021. URL: https://tikz.net/wp-content/uploads/2022/12/axis2D_pseudorapidity.png.
- [29] The Tracker Group Of The CMS Collaboration. *The CMS Phase-1 Pixel Detector Upgrade*. 2020. DOI: 10.48550/ARXIV.2012.14304. URL: <https://arxiv.org/abs/2012.14304>.
- [30] Hendrik Jansen and Thomas Hebbeker. “Study of Unparticle plus Lepton Signatures at CMS”. In: (Dec. 2022).
- [31] Paolo Azzurri. “The CMS Silicon Strip Tracker”. In: *J. Phys.: Conf. Ser.* 41 (2006). 8 pages, 8 figures, talk given at XIX EPS NPDC Conference on New Trends in Nuclear Physics Applications and Technology, September 5-9, 2005 Pavia, Italy Subj-class: Instrumentation and Detectors, pp. 127–134. DOI: 10.1088/1742-6596/41/1/011. URL: <https://cds.cern.ch/record/914891>.
- [32] Christian W. Fabjan and Fabiola Gianotti. “Calorimetry for particle physics”. In: *Rev. Mod. Phys.* 75 (4 Oct. 2003), pp. 1243–1286. DOI: 10.1103/RevModPhys.75.1243. URL: <https://link.aps.org/doi/10.1103/RevModPhys.75.1243>.
- [33] Carsten Hof. “Implementation of a model-independent search for new physics with the CMS detector exploiting the world-wide LHC Computing Grid”. PhD thesis. RWTH Aachen U., 2009.
- [34] The CMS collaboration. “The CMS trigger system”. In: *Journal of Instrumentation* 12.01 (Jan. 2017), P01020–P01020. DOI: 10.1088/1748-0221/12/01/p01020. URL: <https://arxiv.org/pdf/1609.02366.pdf>.
- [35] *Particle-Flow Event Reconstruction in CMS and Performance for Jets, Taus, and MET*. Tech. rep. Geneva: CERN, 2009. URL: <https://cds.cern.ch/record/1194487>.
- [36] The CMS collaboration. “Particle-flow reconstruction and global event description with the CMS detector”. In: *Journal of Instrumentation* 12.10 (Oct. 2017), P10003–P10003. DOI: 10.1088/1748-0221/12/10/p10003. URL: <https://doi.org/10.1088/1748-0221/12/10/p10003>.
- [37] Kim Albertsson et al. “Machine Learning in High Energy Physics Community White Paper”. In: *Journal of Physics: Conference Series* 1085.2 (Sept. 2018), p. 022008. DOI: 10.1088/1742-6596/1085/2/022008. URL: <https://dx.doi.org/10.1088/1742-6596/1085/2/022008>.
- [38] Alan S. Cornell et al. “Boosted decision trees in the era of new physics: a smuon analysis case study”. In: *Journal of High Energy Physics* 2022.4 (Apr. 2022). DOI: 10.1007/jhep04(2022)015. URL: <https://doi.org/10.1007/jhep04/2022/29015>.

- [39] Byron P. Roe et al. “Boosted decision trees as an alternative to artificial neural networks for particle identification”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 543.2-3 (May 2005), pp. 577–584. DOI: 10.1016/j.nima.2004.12.018. URL: <https://www.sciencedirect.com/science/article/abs/pii/S0168900205000355>.
- [40] Tianqi Chen and Carlos Guestrin. “XGBoost”. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, Aug. 2016. DOI: 10.1145/2939672.2939785. URL: <https://doi.org/10.1145%5C%2F2939672.2939785>.
- [41] Yann Coadou. “Boosted decision trees”. In: (Mar. 2022). DOI: 10.1142/9789811234033_0002. arXiv: 2206.09645 [physics.data-an].
- [42] Andreas Döpp et al. “Data-driven Science and Machine Learning Methods in Laser-Plasma Physics”. In: (Nov. 2022). arXiv: 2212.00026 [cs.LG].
- [43] Torben Lange. “Search for rare Higgs boson decays at a center of mass energy of $\sqrt{s} = 13$ TeV with the CMS Experiment at the LHC”. PhD thesis. Universität Hamburg, U. Hamburg, Dept. Phys., 2022.
- [44] Laurits Tani et al. “Evolutionary algorithms for hyperparameter optimization in machine learning for application in high energy physics”. In: *The European Physical Journal C* 81.2 (Feb. 2021). DOI: 10.1140/epjc/s10052-021-08950-y. URL: <https://doi.org/10.1140%5C%2Fepjc%5C%2Fs10052-021-08950-y>.
- [45] Alexandra Carvalho et al. “Higgs pair production: choosing benchmarks with cluster analysis”. In: *Journal of High Energy Physics* 2016.4 (Apr. 2016), pp. 1–28. DOI: 10.1007/jhep04(2016)126. URL: <https://arxiv.org/abs/1507.02245>.
- [46] Thomas Junk. “Confidence level computation for combining searches with small statistics”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 434.2-3 (Sept. 1999), pp. 435–443. DOI: 10.1016/S0168-9002(99)00498-2. URL: <https://doi.org/10.48550/arXiv.hep-ex/9902006>.
- [47] A L Read. “Presentation of search results: the CLs technique”. In: *Journal of Physics G: Nuclear and Particle Physics* 28.10 (Sept. 2002), p. 2693. DOI: 10.1088/0954-3899/28/10/313. URL: <https://dx.doi.org/10.1088/0954-3899/28/10/313>.
- [48] 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”. In: *Journal of High Energy Physics* 2014.7 (July 2014). DOI: 10.1007/jhep07(2014)079. URL: <https://doi.org/10.1007%5C%2Fjhep07%5C%282014%5C%29079>.
- [49] Paolo Nason. “A new method for combining NLO QCD with shower Monte Carlo algorithms”. In: *Journal of High Energy Physics* 2004.11 (Dec. 2004), p. 040. DOI: 10.1088/1126-6708/2004/11/040. URL: <https://dx.doi.org/10.1088/1126-6708/2004/11/040>.

- [50] Stefano Frixione, Paolo Nason, and Carlo Oleari. “Matching NLO QCD computations with parton shower simulations: the POWHEG method”. In: *Journal of High Energy Physics* 2007.11 (Nov. 2007), pp. 070–070. DOI: 10.1088/1126-6708/2007/11/070. URL: <https://doi.org/10.1088%5C%2F1126-6708%5C%2F2007%5C%2F11%5C%2F070>.
- [51] Simone Alioli et al. “A general framework for implementing NLO calculations in shower Monte Carlo programs: the POWHEG BOX”. In: *Journal of High Energy Physics* 2010.6 (June 2010). DOI: 10.1007/jhep06(2010)043. URL: <https://doi.org/10.1007%5C%2Fjhep06%5C%282010%5C%29043>.
- [52] Torbjörn Sjöstrand et al. “An introduction to PYTHIA 8.2”. In: *Computer Physics Communications* 191 (June 2015), pp. 159–177. DOI: 10.1016/j.cpc.2015.01.024. URL: <https://doi.org/10.1016%5C%2Fj.cpc.2015.01.024>.
- [53] The CMS collaboration. “Event generator tunes obtained from underlying event and multiparton scattering measurements”. In: *The European Physical Journal C* 76.3 (Mar. 2016). DOI: 10.1140/epjc/s10052-016-3988-x. URL: <https://doi.org/10.1140%5C%2Fepjc%5C%2Fs10052-016-3988-x>.
- [54] The CMS collaboration. “Extraction and validation of a new set of CMS pythia8 tunes from underlying-event measurements”. In: *The European Physical Journal C* 80.1 (Jan. 2020). DOI: 10.1140/epjc/s10052-019-7499-4. URL: <https://doi.org/10.1140%5C%2Fepjc%5C%2Fs10052-019-7499-4>.
- [55] P. Skands, S. Carrazza, and J. Rojo. “Tuning PYTHIA 8.1: the Monash 2013 tune”. In: *The European Physical Journal C* 74.8 (Aug. 2014). DOI: 10.1140/epjc/s10052-014-3024-y. URL: <https://doi.org/10.1140%5C%2Fepjc%5C%2Fs10052-014-3024-y>.
- [56] Richard D. Ball et al. “Parton distributions from high-precision collider data”. In: *The European Physical Journal C* 77.10 (Oct. 2017). DOI: 10.1140/epjc/s10052-017-5199-5. URL: <https://doi.org/10.1140%5C%2Fepjc%5C%2Fs10052-017-5199-5>.
- [57] The CMS collaboration. “Electron and photon reconstruction and identification with the CMS experiment at the CERN LHC”. In: *Journal of Instrumentation* 16.05 (May 2021), P05014. DOI: 10.1088/1748-0221/16/05/P05014. URL: <https://dx.doi.org/10.1088/1748-0221/16/05/P05014>.
- [58] *Search for ttH production in multilepton final states at sqrt(s) = 13 TeV*. Tech. rep. Geneva: CERN, 2016. URL: <https://cds.cern.ch/record/2141078>.
- [59] *Search for associated production of Higgs bosons and top quarks in multilepton final states at sqrt(s) = 13 TeV*. Tech. rep. Geneva: CERN, 2016. URL: <https://cds.cern.ch/record/2205282>.

- [60] *Measurement of the associated production of a Higgs boson with a top quark pair in final states with electrons, muons and hadronically decaying τ leptons in data recorded in 2017 at $\sqrt{s} = 13$ TeV*. Tech. rep. Geneva: CERN, 2018. URL: <https://cds.cern.ch/record/2649199>.
- [61] “Performance of electron reconstruction and selection with the CMS detector in proton-proton collisions at $\sqrt{s} = 8$ TeV”. In: *Journal of Instrumentation* 10.06 (June 2015), P06005–P06005. DOI: 10.1088/1748-0221/10/06/p06005. URL: <https://doi.org/10.48550/arXiv.1502.02701>.
- [62] Matteo Cacciari, Gavin P Salam, and Gregory Soyez. “The anti- k_r subjet clustering algorithm”. In: *Journal of High Energy Physics* 2008.04 (Apr. 2008), pp. 063–063. DOI: 10.1088/1126-6708/2008/04/063. URL: <https://doi.org/10.1088/1126-6708/2008/04/063>.
- [63] Jesse Thaler and Ken Van Tilburg. “Identifying boosted objects with N-subjettiness”. In: *Journal of High Energy Physics* 2011.3 (Mar. 2011). DOI: 10.1007/jhep03(2011)015. URL: [https://doi.org/10.1007/jhep03\(2011\)015](https://doi.org/10.1007/jhep03(2011)015).
- [64] E. Bols et al. “Jet flavour classification using DeepJet”. In: *Journal of Instrumentation* 15.12 (Dec. 2020), P12012–P12012. DOI: 10.1088/1748-0221/15/12/p12012. URL: <https://doi.org/10.1088/1748-0221/15/12/p12012>.
- [65] The CMS collaboration. “Reconstruction and identification of lepton decays to hadrons and at CMS”. In: *Journal of Instrumentation* 11.01 (Jan. 2016), P01019. DOI: 10.1088/1748-0221/11/01/P01019. URL: <https://dx.doi.org/10.1088/1748-0221/11/01/P01019>.
- [66] The CMS collaboration. “Identification of hadronic tau lepton decays using a deep neural network”. In: *Journal of Instrumentation* 17.07 (July 2022), P07023. DOI: 10.1088/1748-0221/17/07/p07023. URL: <https://doi.org/10.1088/1748-0221/17/07/p07023>.
- [67] The CMS collaboration. “Evidence for associated production of a Higgs boson with a top quark pair in final states with electrons, muons, and hadronically decaying leptons at $\sqrt{s} = 13$ TeV”. In: *Journal of High Energy Physics* 2018.8 (Aug. 2018). DOI: 10.1007/jhep08(2018)066. URL: <https://arxiv.org/abs/1803.05485>.
- [68] Glen Cowan et al. “Asymptotic formulae for likelihood-based tests of new physics”. In: *The European Physical Journal C* 71.2 (Feb. 2011). DOI: 10.1140/epjc/s10052-011-1554-0. URL: <https://doi.org/10.1140/epjc/s10052-011-1554-0>.
- [69] *Procedure for the LHC Higgs boson search combination in Summer 2011*. Tech. rep. Geneva: CERN, 2011. URL: <https://cds.cern.ch/record/1379837>.
- [70] The CMS collaboration. . *Combine - Statistical analysis software tools for CMS*. URL: <https://cms-analysis.github.io/HiggsAnalysis-CombinedLimit/>.

