This electronic thesis or dissertation has been downloaded from Explore Bristol Research, http://research-information.bristol.ac.uk

*Author:*
**Taylor, Joe**

*Title:*
**A search for light Higgs bosons in supersymmetric decay cascades with the CMS detector**

# A search for light Higgs bosons in supersymmetric decay cascades with the CMS detector

Joseph Ross Taylor

A dissertation submitted to the University of Bristol in accordance with the requirements for award of the degree of Doctor of Philosophy in the Faculty of Science, School of Physics.

June 2019

~ 49800 words

# Abstract

This thesis presents a search for pairs of light Higgs bosons produced in supersymmetric decay cascades. The final state targeted is that where both Higgs bosons decay into $b\bar{b}$ pairs. The analysis uses proton-proton collision data recorded with the CMS detector at a centre-of-mass energy of 13 TeV. The full data sets from 2016 and 2017 are used, corresponding to integrated luminosities of 35.9 fb$^{-1}$ and 41.5 fb$^{-1}$, respectively.

The signal model exists within the framework of the Next-to Minimal Supersymmetric Standard Model. It is of interest because, under certain mass configurations, the $E_{\mathrm{T}}^{\mathrm{miss}}$ in the events can be highly suppressed. This produces an all jet final state for which more conventional supersymmetry searches would have reduced sensitivity.

# Acknowledgements

Naturally, as a student of Bristol University and the Rutherford Appleton Laboratory, and being a member of a large, highly collaborative experiment, I am indebted to a myriad of people, whom it would not be possible to list coherently. I would, however, like to explicitly extend my gratitude to Jim Brooke, Henning Flaecher, and Claire Shepherd-Themistocleous; for their supervision over the last four years, and to Robin Aggleton; for being so generous with his time as I was getting started. Special acknowledgements should also be directed to all the different lunchtime and after-work-beer crews; for making me laugh and helping me maintain a certain degree of sanity.

On a personal note, I would like to thank my parents; for their unrelenting love and support, right from day zero. Finally, I would like to thank Sophia; for everything.

# Author's Declaration

I declare that the work in this dissertation was carried out in accordance with the requirements of the University's Regulations and Code of Practice for Research Degree Programmes and that it has not been submitted for any other academic award. Except where indicated by specific reference in the text, the work is the candidate's own work. Work done in collaboration with, or with the assistance of, others, is indicated as such. Any views expressed in the dissertation are those of the author.

Signed                                    Date

Chapters 2 and 3 are reviews of particle theory and the CMS experiment, respectively. They do not represent the Author's original work.

Chapter 4 begins with a review of the CMS Level-1 trigger. Consequently, within this chapter, only Sections 4.3.2 and 4.4 represent the Author's original work. To specify further:

- All the figures in Section 4.3.2 were made by the Author. They do, however, represent a much wider collection of work which was performed in collaboration with the calo-layer-2 subgroup of the CMS Level-1 trigger working group.

- The Level-1 jet energy corrections, described in Section 4.4, were performed by the Author using a framework already written by Robin Aggleton.

All the remaining chapters and appendices are the Author's original work, with the exception of:

- The cross sections in Table 5.6, which were calculated by Alex Titterton.

- The data/MC scale factors stated in Section 7.1, which were derived by various working groups within the CMS collaboration (explicit acknowledgements are provided in the main text).

It should also be noted that the statistical analysis described in Section 7.2 was conducted using the CMS 'combine' tool.

# Contents

# List of Figures

# List of Tables

# 1

---

# Introduction

---

The objective of particle physics is to describe all the different elementary particles that can exist in the Universe. Ninety years ago, there were only believed to be two fundamental particles; the proton and the electron. Since then, through an interplay between theory and experiment, a whole new set of fundamental particles, and forces, have been discovered. These particles, and their interactions, are elegantly described by the Standard Model (SM); a quantum field theory with local gauge symmetries. Despite its success, the Standard Model is not a complete representation of the Universe and a variety of Beyond the Standard Model (BSM) extensions have been proposed which would further expose the true reality of Nature. The primary goal of contemporary experimental particle physics is, therefore, to observe evidence of BSM physics. To date, however, this has not been achieved.

This thesis is a description of my search for BSM physics. It begins, in Chapter 2, with an overview of the Standard Model and the simplest consistent supersymmetric (SUSY) extensions. It is within the framework of the Next-to Minimal Supersymmetric Standard Model (NMSSM) that the signal model exists. The search for this signal model was conducted using data collected at the CMS detector. In Chapter 3, an overview of the detector, and the data reconstruction methods, is presented. This

is followed, in Chapter 4, by a more detailed description of the CMS Level-1 trigger, in order to facilitate a discussion about my contribution to the system. The remaining chapters provide a precise description of the analysis. In Chapter 5, the signal model is formally introduced, along with information about the data sets used. The signal model is rather unique and, hitherto, a search had not been conducted for it. As a consequence, the whole analysis strategy needed to be conceptualized; the details of which are provided in Chapter 6. These include descriptions of the triggering, event selection, and background estimation used in the analysis. In Chapter 7, the various systematic uncertainties are discussed, followed by a presentation of the analysis results. Finally, in Chapter 8, I provide my concluding remarks.

## 1.1   List of Definitions

The Cartesian coordinate system, used to represent positions within the CMS detector, is defined as follows: the origin is the nominal proton-proton collision point, the $x$-axis points towards the centre of the LHC ring, the $y$-axis points vertically upwards, and the $z$-axis points along the LHC beamline, in an anticlockwise direction when viewed from above. The radial distance from the beamline, $r$, is therefore given by $r = \sqrt{x^2 + y^2}$.

Directions, relative to the origin, are described by the azimuthal angle, $\phi$, defined by:

$$\phi = \arctan \frac{y}{x} \tag{1.1}$$

and the pseudorapidity, $\eta$, defined by:

$$\eta = -\ln\left(\tan\frac{\theta}{2}\right) \tag{1.2}$$

where the polar angle, $\theta$, is defined by:

$$\theta = \arctan\left(\frac{\sqrt{x^2 + y^2}}{z}\right) \tag{1.3}$$

The angular separation of two objects, $i$ and $j$, is described by the $\Delta R$ quantity, which is defined by:

$$\Delta R = \sqrt{(\eta_i - \eta_j)^2 + (\phi_i - \phi_j)^2} \tag{1.4}$$

Transverse momentum is defined as the component of an objects momentum perpendicular to the beamline:

$$\overrightarrow{p_{\mathrm{T}}} = \overrightarrow{p_x} + \overrightarrow{p_y} \tag{1.5}$$

Often, however, one is simply referring to the magnitude of this vector:

$$p_{\mathrm{T}} = \sqrt{p_x^2 + p_y^2} \tag{1.6}$$

Instantaneous luminosity, $\mathcal{L}$, is defined as:

$$\mathcal{L} = \frac{1}{\sigma} \frac{dN}{dt} \tag{1.7}$$

where $N$ is the number of events for a process with cross section $\sigma$. Instantaneous luminosity is measured in units of $\mathrm{cm}^{-2}\mathrm{s}^{-1}$. Cross sections are measured in units of picobarns (pb), where one barn equals $10^{-28}$ $\mathrm{m}^2$. Integrated luminosity (the instantaneous luminosity integrated over a given time period) is measured in inverse femtobarns ($\mathrm{fb}^{-1}$).

## 1.2 List of Conventions

Natural units, with $\hbar = c = 1$, are used. As a consequence, energy, momentum, and mass are measured in the same units (electronvolts).

Electrical charge is measured in units of $e$, the magnitude of the charge carried by an electron.

The symbol $\sqrt{s}$ is used to denote the centre-of-mass energy of a system.

# 2

---

# The Standard Model and Supersymmetry

---

This chapter begins with a brief theoretical overview of the Standard Model (SM); a representation of all the known elementary particles and their interactions. Despite the success of the Standard Model, it does have a number of shortcomings, which are discussed. Some of these problems can be solved by Supersymmetry (SUSY). An outline of Supersymmetry is provided, followed by a description of the two simplest supersymmetric extensions of the Standard Model.

## 2.1   The Standard Model

### 2.1.1   Theoretical Overview

The Standard Model of particle physics [7–9] is a quantum field theory [10] in which the excitations of the fields give rise to all the elementary particles.

There are twelve different spin-$\frac{1}{2}$ fermions, categorised as quarks and leptons. There are six quarks in total, divided into three different generations. Within each generation, there is a quark with electrical charge +2/3 (up, charm, top) and a quark

with electrical charge $-1/3$ (down, strange, bottom). The six leptons are also divided into three different generations. Within each generation, there is a lepton with electrical charge $-1$ (electron, muon, tau) and a corresponding neutrino with no electrical charge.

The SM lagrangian density is invariant under local $\mathrm{SU(3)_C} \otimes \mathrm{SU(2)_L} \otimes \mathrm{U(1)_Y}$ gauge transformations. Quantum chromodynamics (QCD) is represented by the $\mathrm{SU(3)_C}$ component. The symmetry acts on the quark fields, as they possess 'colour' quantum numbers, and has eight corresponding spin-1 gauge bosons called gluons. The electroweak interactions are represented by the $\mathrm{SU(2)_L} \otimes \mathrm{U(1)_Y}$ gauge symmetries. The $\mathrm{SU(2)_L}$ component acts on the left-handed chiral projections of all the fermion fields, organised into isospin doublets, and has three corresponding spin-1 gauge bosons ($W^+$, $W^0$, $W^-$). In the lepton sector, the isospin doublets are ordered by generation. In the quark sector, there is mixing between the generations, which is quantified by the Cabibbo-Kobayashi-Maskawa matrix. The $\mathrm{U(1)_Y}$ component acts on both the left and right-handed chiral projections of the fermion fields. These fields couple to the corresponding spin-1 gauge boson ($B^0$) with a strength proportional to their hypercharge, defined by $Y = 2\left(Q_{\mathrm{EM}} - T_3\right)$, where $T_3$ is the isospin projection. The photon is recovered as a mixture of the $W^0$ and $B^0$ gauge bosons. The orthogonal mixture provides the $Z^0$ boson.

The SM lagrangian density contains a Higgs field which, through the Higgs mechanism, enables gauge invariant mass terms. The Higgs field is an isospin doublet of complex scalars with hypercharge $Y = 1$. It has a non-zero vacuum expectation value which spontaneously breaks the electroweak gauge symmetry. The gauge bosons obtain their mass from the Higgs field kinetic term. The masses are dependent on the vacuum expectation value and the electroweak coupling parameters. The $W^\pm$ and $Z^0$ bosons are measured to have masses of 80.4 GeV and 91.2 GeV, respectively [7]. The photon does not acquire a mass and its interactions, with particles that are electrically charged, obey a residual U(1) symmetry. The fermions obtain their mass from Yukawa couplings to the Higgs field. Each mass term has a free

parameter that has to be determined experimentally. The two heaviest fermions are the bottom quark and the top quark, with masses of 4.2 GeV and 173 GeV, respectively [7].

Following electroweak symmetry breaking, the Higgs field has one remaining degree of freedom, the Higgs boson. This scalar field has an observed mass of 125 GeV [7], couples with all the massive particles, and is CP-even (i.e. does not change sign under the combination of a charge conjugation and parity transformation).

### 2.1.2 Problems with the Standard Model

To date, the Standard Model has successfully described all the phenomena observed in high-energy experiments. Despite this achievement, the Standard Model is not a complete representation of the Universe. Most notably, it does not contain a theory of gravity. This guarantees that new physics exists at a yet unexplored energy scale. In the case of quantum gravity, however, this corresponds to exceptionally high energies.

The Standard Model has two other notable issues. One is that the Standard Model does not contain a legitimate Dark Matter (DM) candidate [11]. The other is the hierarchy problem in relation to the Higgs boson mass [12]. Quantum loop corrections to the square of the Higgs boson mass are quadratically sensitive to the new physics energy scale ($\approx 10^{19}$ GeV for quantum gravity). Consequently, the Higgs boson, with an observed mass of 125 GeV (the electroweak energy scale), requires a finely tuned 'bare' mass to cancel the mass correction.

## 2.2 Supersymmetry

Supersymmetry is a symmetry relating bosons and fermions [12, 13], which can be used to extend the Standard Model. By incorporating supersymmetry, a new collection of particles are introduced, as every bosonic field obtains a fermionic superpartner and every fermionic field obtains a bosonic superpartner. Each superpartner pairing has the same number of bosonic and fermionic degrees of freedom. Moreover,

the fields have the same electric charge, weak isospin, and colour degrees of freedom. The SUSY particles do not, however, have the same masses as their SM counterparts, as supersymmetry has not been observed. As a consequence, supersymmetry is a broken symmetry.

Supersymmetry is an attractive BSM theory. It can provide both a legitimate DM candidate (see Section 2.2.1) and an elegant solution to the hierarchy problem. In the quantum loop corrections to the square of the Higgs boson mass, the contribution from each SUSY particle is such that it exactly cancels the quadratic divergence from its SM counterpart, due to its different spin. The correction instead scales with the square of the heaviest SUSY particle mass, constituting far less fine tuning. This solution does start to become problematic, however, if the SUSY particles have masses beyond the TeV scale. Another advantage of supersymmetry is that the $U(1)_Y$, $SU(2)_L$, and $SU(3)_C$ running gauge couplings can have equivalent values at an energy scale of around $10^{16}$ GeV, allowing for grand unified theories with a single force.

## 2.2.1   Particles of the MSSM

The Minimal Supersymmetric Standard Model (MSSM) is the simplest SUSY model consistent with the Standard Model. In this section, a brief summary of the new particles predicted by the MSSM is provided. In-depth reviews of the MSSM can be found in Ref. [12, 13].

The complex scalar superpartners of the quarks and leptons are called squarks ($\widetilde{q}$) and sleptons ($\widetilde{\ell}$), respectively. The spin-$\frac{1}{2}$ superpartners of the gauge bosons are called gauginos. Specifically, the gluino ($\widetilde{g}$) is the superpartner of the gluon, and the winos ($\widetilde{W}^+$, $\widetilde{W}^0$, $\widetilde{W}^-$) and bino ($\widetilde{B}$) are the superpartners of the electroweak gauge bosons.

Supersymmetric theories cannot be consistent with the Standard Model if superpartners are created from the SM Higgs field [12]. Instead, two Higgs fields, with hypercharges of $Y = \pm 1$, are required. The two isospin doublets employed are

$H_u = (H_u^+, H_u^0)$, with $Y = +1$, and $H_d = (H_d^0, H_d^-)$, with $Y = -1$. The spin-$\frac{1}{2}$ superpartners corresponding to the complex scalar fields are called higgsinos ($\widetilde{H}_u^+$, $\widetilde{H}_u^0$, $\widetilde{H}_d^0$, $\widetilde{H}_d^-$). Following electroweak symmetry breaking, there are five Higgs boson mass eigenstates; two CP-even neutral scalars, a single CP-odd neutral scalar, and two scalars which are charge conjugate partners. Note that the observed Standard Model-like Higgs boson can be recovered as the lightest CP-even neutral scalar.

Due to the effects of electroweak symmetry breaking, the charged winos ($\widetilde{W}^+$, $\widetilde{W}^-$) and higgsinos ($\widetilde{H}_u^+$, $\widetilde{H}_d^-$) mix to form two chargino mass eigenstates, $\widetilde{\chi}_i^\pm$, where $i = 1, 2$. Similarly, the neutrally charged gauginos ($\widetilde{B}$, $\widetilde{W}^0$) and higgsinos ($\widetilde{H}_u^0$, $\widetilde{H}_d^0$) mix to form four neutralino mass eigenstates, $\widetilde{\chi}_i^0$, where $i = 1, 2, 3, 4$ (labelled in ascending mass order). A good dark matter candidate is provided [11] when the lightest neutralino, $\widetilde{\chi}_1^0$, is the Lightest SUSY Particle (LSP). If this is the case, all the SUSY particle decay chains end up with an LSP in the final state.

The MSSM has a large parameter space in which a myriad of SUSY particle mass spectra are possible. In addition, there are numerous ways in which the SUSY particles can couple, both with themselves and with the SM particles. Some particle couplings are forbidden due to R-parity, an additional symmetry imposed on the fields in order to conserve baryon and lepton number. Note that it is R-parity that prevents the LSP from decaying into SM particles.

### 2.2.2 Particles of the NMSSM

The Next-to Minimal Supersymmetric Standard Model [14] (NMSSM) extends upon the MSSM by introducing an additional complex scalar field, which is a gauge singlet of $SU(3)_C \otimes SU(2)_L \otimes U(1)_Y$. The singlet field mixes with the neutrally charged Higgs fields to produce an additional CP-even neutral scalar and CP-odd neutral scalar, relative to the MSSM. Note that these new Higgs bosons can be configured so that they are lighter than the Standard Model-like Higgs boson. The spin-$\frac{1}{2}$ superpartner of the singlet field, the singlino, mixes with the higgsinos and the neutrally charged gauginos to produce an additional neutralino, relative to the MSSM.

The motivation for introducing the gauge singlet is that it solves the so-called $\mu$-problem of the MSSM [12, 14]. In the MSSM, the lagrangian must contain a SUSY mass term, $\mu$, which couples to the $H_u$ and $H_d$ fields, providing the higgsino mass terms and the squared-mass term in the Higgs scalar potential. To allow an appropriate Higgs vacuum expectation value, the value of $\mu$ has to be the order of the electroweak scale. The $\mu$-problem arises as there is no clear reason why this should be the case. The NMSSM solves this problem by dynamically generating the SUSY mass term. In the lagrangian, the $\mu$ term is replaced by the gauge singlet (multiplied by a dimensionless constant). The vacuum expectation value of this scalar field is of the desired mass scale and, thus, solves the $\mu$-problem. The introduction of the gauge singlet, $S$, and its solution to the $\mu$-problem, results in the following NMSSM specific soft SUSY breaking lagrangian:

$$ -\mathcal{L}_{\text{NMSSM}}^{\text{soft}} = m_S^2 |S|^2 + \left( \lambda A_\lambda H_u H_d S + \xi_S S + \frac{1}{2} m_S'^2 S^2 + \frac{1}{3} \kappa A_\kappa S^3 + \text{h.c.} \right) \quad (2.1) $$

where $m_S^2$, $\lambda$, $A_\lambda$, $\xi_S$, $m_S'^2$, $\kappa$, and $A_\kappa$ are parameters of the model.

The search for BSM physics conducted in this thesis is for a SUSY model that exists within the framework of the NMSSM. The signal model is formally introduced in Section 5.1. It utilizes the scenario where the LSP is a singlino-like neutralino, which leads to unique SUSY decay cascades that would have gone undetected by more conventional SUSY searches.

# 3

---

# The LHC and the CMS

# Experiment

---

This chapter begins by introducing the Larger Hadron Collider (LHC); a particle accelerator that provides proton-proton (pp) collisions at unprecedented energies. The Compact Muon Solenoid (CMS) detector measures the outgoing particles from these collisions. An overview of the CMS detector is presented, followed by an outline of the reconstruction methods. Finally, the jets reconstructed in CMS are discussed in detail, as they are of central importance in the remaining chapters.

## 3.1   The Large Hadron Collider

The LHC [15, 16] is a particle accelerator situated at the CERN laboratory in Geneva. It is approximately circular in shape with a circumference of 27 km. The LHC contains two beampipes in which proton beams circulate in opposite directions. Where the two beams intersect, proton-proton collisions are provided at an unprecedented centre-of-mass energy and luminosity. The LHC was designed to accelerate protons up to an energy of 7 TeV, thus producing pp interactions with $\sqrt{s} = 14$ TeV. The instantaneous luminosity was designed to be $1.0 \times 10^{34}$ cm$^{-2}$s$^{-1}$.

Figure 3.1: Schematic diagram of the CERN accelerator complex.

The protons are accelerated by a sequence of four machines before entering the LHC. Linac 2 first accelerates the protons to an energy of 50 MeV. Next, the Proton Synchrotron Booster (PSB) increases the proton energy to 1.4 GeV before the Proton Synchrotron (PS) further raises it to 25 GeV. The Super Proton Synchrotron (SPS) then accelerates the protons to an energy of 450 GeV. Following this, the protons are injected into the LHC. A schematic diagram of the particle accelerators at CERN is provided in Figure 3.1.

Radiofrequency cavities, driven by klystrons, are used to accelerate the protons within the LHC. After about 20 minutes, the protons reach their maximum energy of 7 TeV. The oscillating electric fields in radiofrequency cavities cause protons to be squeezed into bunches (the bunch structure is obtained in the PS). Consequently, the pp collisions occur within the crossing of two opposing proton bunches. The bunch crossings, also referred to as events, occur at a rate of 40 MHz.

The orbits of the protons are controlled by a series of electromagnets. The proton trajectories are bent around the LHC ring using superconducting dipoles, which generate magnetic fields of 8.3 T. The beams are repeatedly focussed in one spatial

**CMS Integrated Luminosity Delivered, pp**

Data included from 2010-03-30 11:22 to 2018-10-26 08:23 UTC



Figure 3.2: The cumulative integrated luminosity delivered to CMS by the LHC, separated into the different years of proton-proton operations.

dimension, and defocussed in the orthogonal dimension, by quadrupole magnets. The width of each beam is typically 200 µm, however, at the intersection points it is squeezed down to 16 µm to increase the instantaneous luminosity.

There are four collision points along the LHC ring. Each one has a particle detector built around it to study the outgoing particles. The ATLAS [17] and CMS [3] experiments are both general purpose particle detectors. The LHCb [18] experiment is designed to study the physics of b-quarks and c-quarks. The ALICE [19] detector is specialized to study heavy ion collisions, which the LHC provides at the end of each year of operations.

There have been two main eras of operations at the LHC. The first era, Run-I, occurred between 2010 and 2012. It provided pp collisions with $\sqrt{s} = 7$ TeV and $\sqrt{s} = 8$ TeV. The second era, Run-II, occurred between 2015 and 2018. It provided pp collisions with $\sqrt{s} = 13$ TeV, which is just below the design energy. Figure 3.2

shows the integrated luminosity delivered to CMS by the LHC for each operational year. During Run-II, the LHC was very successful in providing CMS with pp collision data. In 2018, the instantaneous luminosity was typically around $1.9 \times 10^{34}$ cm$^{-2}$s$^{-1}$; nearly double the value in the design specification. As a consequence, the mean number of interactions per bunch crossing was 37, significantly exceeding the average of 25 predicted in the design specification.

## 3.2    The CMS Detector

CMS is a general purpose particle detector that was designed to accommodate a broad physics programme [20]. The detector is capable of discovering a myriad of new physics processes, with diverse final states, and can also perform precision measurements of SM parameters [21]. It achieves this by identifying, and measuring the kinematic properties, of the stable outgoing particles from the LHC pp collisions.

As a general purpose detector, CMS has to be able to reconstruct all final-state particles. The design was focussed, however, on some key physics objectives. The primary goal was to search for the Higgs boson. Such a particle was discovered by CMS [22] and shown to be compatible with the SM Higgs boson [23]. The mass of the Higgs boson is measured to be 125 GeV, however, prior to its discovery, the mass was unknown. Therefore, the CMS detector required sensitivity to the Higgs boson across a large mass range spanning from 90 GeV (the expected discovery limit of the LEP [24] experiments) to the TeV scale. Figure 3.3 shows how the Higgs boson branching ratios vary in this mass range. For Higgs boson masses below 135 GeV, H $\rightarrow$ b$\bar{\text{b}}$ is the primary decay mode. This, however, is not a good discovery channel due to the very large QCD multi-jet background. Instead, the H $\rightarrow \gamma\gamma$ decay was identified as the most promising discovery mode because, despite having a small branching ratio, this channel has a clean signature in pp collisions. For Higgs boson masses above 135 GeV, H $\rightarrow$ WW is the primary decay mode. The best discovery channel, however, is H $\rightarrow$ ZZ where both Z bosons decay into di-electrons or di-muons. This channel has a clear signature and a sizeable branching ratio, especially

Figure 3.3: Branching ratios of the SM Higgs boson as a function of mass [1].

for Higgs boson masses above 200 GeV. In accordance with the anticipated Higgs boson discovery modes, CMS had to detect photons, electrons, and muons with large geometric acceptance across a broad $p_\mathrm{T}$ range. Very good energy and spatial resolution was required in order to clearly reconstruct the Higgs boson mass peak.

Another important objective during the design of CMS, was the search for Beyond the Standard Model (BSM) physics. One BSM theory targeted was that of heavy vector bosons, such as the Z' and W', with leptonic decay modes. Therefore, the ability to reconstruct multi-TeV electrons and muons was integral to the detector design. Another BSM theory targeted was SUSY, identified by the decay cascades following the production of squarks and gluinos. There are a multitude of final states possible, but most involve high $p_\mathrm{T}$ LSPs which traverse the detector without interacting, thus creating a sizeable momentum imbalance in the transverse plane of the detector. This property is quantified by the $E_\mathrm{T}^\mathrm{miss}$ in an event, the negative vector sum of the $p_\mathrm{T}$ of all the particles detected. Consequently, the measurement of $E_\mathrm{T}^\mathrm{miss}$ was an important consideration during the design process.

CMS DETECTOR

Total weight       : 14,000 tonnes
Overall diameter : 15.0 m
Overall length    : 28.7 m
Magnetic field     : 3.8 T

STEEL RETURN YOKE
12,500 tonnes

SILICON TRACKERS
Pixel (100x150 μm²) ~1 m² ~66M channels
Microstrips (80–180 μm) ~200 m² ~9.6M channels

SUPERCONDUCTING SOLENOID
Niobium titanium coil carrying ~18,000 A

MUON CHAMBERS
Barrel: 250 Drift Tube, 480 Resistive Plate Chambers
Endcaps: 468 Cathode Strip, 432 Resistive Plate Chambers

PRESHOWER
Silicon strips ~16 m² ~137,000 channels

FORWARD CALORIMETER
Steel + Quartz fibres ~2,000 Channels

CRYSTAL
ELECTROMAGNETIC
CALORIMETER (ECAL)
~76,000 scintillating PbWO₄ crystals

HADRON CALORIMETER (HCAL)
Brass + Plastic scintillator ~7,000 channels

Figure 3.4: Cutaway diagram of the CMS detector [2].

The CMS detector [3, 4] is the solution that was found to the full set of physics objectives, with the constraint of being affordable. A diagram of the CMS detector is provided in Figure 3.4. It is composed of multiple detector layers, cylindrical in shape, which surround the interaction point. The overall length is 28.7 m and the overall diameter is 15.0 m. The central feature of the CMS detector is the large superconducting solenoid magnet. Within the solenoid, travelling outwards from the interaction point, are: the inner tracker, which measures the trajectories of charged particles; the electromagnetic calorimeter, which absorbs electrons and photons; and the hadronic calorimeter, which absorbs hadrons. Outside the solenoid, the only (known) final-state particles remaining are muons and neutrinos. Muons lose only a small fraction of their energy in the calorimeters due to their relatively large mass ($m_\mu \approx 207\,m_{\mathrm{e}}$) and because they don't interact via the strong force. Instead, the muon trajectories are measured outside the solenoid by a set of muon detectors. Neutrinos only interact via the weak force, so exit CMS without interacting. They are detected indirectly through $E_{\mathrm{T}}^{\mathrm{miss}}$ measurements.

### 3.2.1 Solenoid Magnet

The CMS detector was designed around the choice of magnet system, in which a large superconducting solenoid was selected [25]. The solenoid has a length of 12.9 m and an inner diameter of 5.9 m. The superconducting coil contains four layers of niobium-titanium. Each conducting layer carries a current of 19.5 kA. The coil has 542 turns and, thus, there are 2168 conductor windings in total. This creates a uniform 3.8 T magnetic field, in the longitudinal direction, within the solenoid. There is an iron yoke around the outside, and at the extremities, of the solenoid, which is responsible for the return of the magnetic flux.

The magnetic field inside the solenoid bends the trajectories of all electrically charged particles in the transverse plane. The radius of curvature is proportional to the $p_T$ of the particles and inversely proportional to the magnetic field strength. Consequently, the inner tracker can be used to measure the $p_T$ of the charged particles. The calorimeters are also situated within the solenoid, which is possible due to its large diameter. This minimizes the amount of material in front of the calorimeters, allowing for better energy resolution. Outside the solenoid, the only charged particles remaining are muons. Their trajectories are bent by the magnetic flux in the iron yoke. Muon detectors, interspaced between the yoke layers, measure the trajectories and, thus, provide a $p_T$ measurement.

The necessity for such a strong magnetic field was driven by the requirements of the muon momentum resolution. The $p_T$ of the muons is measured solely by the trajectories that they take in the magnetic field, both inside and outside the solenoid. As the radius of curvature increases, the $p_T$ resolution degrades. Therefore, the $p_T$ resolution degrades for muons of greater $p_T$, as this increases the radius of curvature. This can be counteracted by increasing the magnetic field strength, which decreases the radius of curvature. In order to achieve the requirement that $\Delta p_T / p_T \approx 10\%$ for muons with $p_T = 1$ TeV, a 4 T magnetic field was needed.

## 3.2.2   Inner Tracker

The inner tracker is situated at the centre of the CMS detector, surrounding the interaction point [26, 27]. It is cylindrical in shape, with a length of 5.8 m and a diameter of 2.5 m. One role of the tracker is to determine the trajectories of the electrically charged particles created in the pp collisions. The trajectories form an important part in measuring the $p_T$ of the charged particles. The trajectories are also used, either directly or indirectly, in the identification of all reconstructed particles. The design specification was to reconstruct the trajectories of all charged particles with $p_T > 1$ GeV and $|\eta| < 2.5$. Another role of the tracker is to precisely identify the primary vertex, the location where the pp interaction of interest occurred in a bunch crossing. This is defined, in offline reconstruction, as the point where the quadratic sum of the $p_T$ of the emerging tracks is maximal. The tracker is also used to identify any potential secondary vertices, the displaced positions where hadrons containing c-quarks or b-quarks decay.

The high luminosity achieved by the LHC comes with the associated difficulties of a high collision rate and high pile-up (additional inelastic pp collisions); both of which placed requirements on the tracker design. With bunch crossings occurring every 25 ns, a fast detector response was necessary in order to associate trajectories to the correct event. The high pile-up conditions cause an average of around 1000 particles crossing the tracker every event. Consequently, the tracker was required to have high granularity to ensure the correct trajectories could be reconstructed. It also meant that the tracker had to be radiation-hard due to the intensity of the incident particles. Another significant consideration when designing the tracker, was to keep the amount of material to a minimum. This is important because it reduces the effects of multiple scattering, bremsstrahlung and photon conversion. The best solution found was to build the tracker solely out of silicon sensors [28].

The layout of the tracker is shown in Figure 3.5. It is composed of two subsystems; the pixel detector and the strip tracker. Both subsystems have a geometric acceptance of $|\eta| < 2.5$.

Figure 3.5: Diagram of the inner tracker in the $r$-$z$ plane [3]. Each line represents a detector layer. Note that the pixel detector corresponds to the legacy version.

The pixel detector lies closest to the beamline. It is composed of silicon pixel cells with size $100 \times 150$ µm$^2$ in the $r$-$\phi$ and $z$ dimensions, respectively. Initially, the pixel detector had three cylindrical layers at radii of 4.4, 7.3, and 10.2 cm together with two disks on either side. This was upgraded, for 2017 operations, to have four cylindrical layers at radii of 3.0, 6.8, 10.9, and 16.0 cm together with three disks on either side [29]. There are two main reasons why the area of the pixel cells are so small. One reason is that the pixel detector is used for the secondary vertex reconstruction. This is best achieved by pixel cells with high granularity in both the $r$-$\phi$ and $z$ dimensions. Another reason is to ensure a low occupancy in the pixel cells. The number of particles per unit area projected onto the pixel detector is very high. This is due to its close proximity to the interaction point, and the large number of particles emerging per event. In order to have a low occupancy, the pixel cell area must be small.

The strip tracker is made of silicon micro-strip sensors. It surrounds the pixel detector, with layers extending from 20 cm to 116 cm in the radial direction. As can be seen in Figure 3.5, the strip tracker is composed of three main parts: the inner barrel and disks (TIB/TID), the outer barrel (TOB), and the endcaps (TEC). The cell sizes vary throughout the strip tracker, but all roughly have a pitch of 100 µm

and a length of 100 mm. With regards to occupancy, they can afford to have a larger area than the pixel cells due to the reduced particle flux at larger radii. In the TIB and TOB, the strips lie parallel to the beamline and in the TID and TEC, the strips lie parallel to the radial. This configuration provides high granularity measurements in the transverse plane.

### 3.2.3 Electromagnetic Calorimeter

The electromagnetic calorimeter [30] (ECAL) stops, and measures the energies of, electrons and photons. It has a central barrel, which surrounds the inner tracker, and endcaps on either side. The ECAL is a homogeneous calorimeter [28] composed of lead tungstate ($PbWO_4$) crystals. When a high energy electron or photon enters a crystal, an electromagnetic shower is produced in which a cascade of bremsstrahlung and pair production processes occur. Using photodetectors, the resultant scintillation light is then used to calculate the energy of the initial particle.

There are a number of reasons why $PbWO_4$ crystals were chosen for the ECAL. First, they provide good energy resolution, a critical requirement for electrons and photons. The crystals are also radiation resistant, allowing them withstand the high intensity of incident particles arising from the LHC collisions. In addition, they have a fast response, ensuring that the energy deposits are attributed to the correct event. Finally, the crystals have a short radiation length [28] of 0.89 cm, allowing the ECAL to be compact, and a small Moliere radius [28] of 2.2 cm, allowing for good position resolution.

Figure 3.6 shows the layout of the ECAL. The barrel (EB) begins at $r = 1.29$ m and covers the range $|\eta| < 1.479$. The endcap (EE) begins at $|z| = 3.15$ m and covers the range $1.479 < |\eta| < 3.0$. In order avoid cracks, the crystals are aligned so that they point just beyond the interaction point. In the barrel, the crystals have a length of 230 mm, corresponding to 25.8 radiation lengths. This is sufficient to contain 98% of the energy of a 1 TeV electron or photon. The crystal front faces have an area of $22 \times 22$ mm$^2$, such that their width matches the Moliere radius. In the endcap,

Figure 3.6: Diagram of the ECAL in the $r$-$z$ plane [4]. Each blue rectangle represents a PbWO$_4$ crystal.

the crystals are slightly shorter (220 mm) and slightly wider ($29 \times 29$ mm$^2$). In front of the endcaps, in the $1.653 < |\eta| < 2.6$ region, there is a preshower detector (ES). Its purpose is to identify neutral pions decaying into closely separated photon pairs, ensuring that they are not reconstructed as single photons. The preshower detector is a sampling calorimeter [28] with only two layers, designed to start the electromagnetic showers. The absorbing layers are made of lead and the detection layers are made from silicon strip sensors. The silicon cells have a pitch of 1.9 mm and are aligned orthogonally in the two planes. This provides the preshower detector with a much finer granularity than the ECAL crystals.

The relative energy resolution of the ECAL, as a function of electron energy, was measured to be [31]:

$$\frac{\sigma}{E} = \frac{2.8\%}{\sqrt{E/\text{GeV}}} \oplus \frac{12\%}{E/\text{GeV}} \oplus 0.3\% \tag{3.1}$$

The first term is the stochastic term, driven by random fluctuations in the lateral shower containment and the photo-statistics. The second term is due to noise in the electronics and digitization. The final term is due to crystal non-uniformity and calibration uncertainty.

Figure 3.7: Diagram of CMS in the *r-z* plane with the endcap (HE), barrel (HB), outer (HO), and forward (HF) hadronic calorimeters labelled [3].

### 3.2.4   Hadronic Calorimeter

The hadronic calorimeter [32] (HCAL) is used to reconstruct hadron jets (see Section 3.4) and is critical in the calculation of $E_{\mathrm{T}}^{\mathrm{miss}}$. It works by stopping, and measuring the energy of, hadron particles. The HCAL is composed of four different subsystems, which can be seen in Figure 3.7. Within the solenoid magnet, there is the barrel hadronic calorimeter (HB), covering the range $|\eta| < 1.3$, and the endcap hadronic calorimeter (HE), covering the range $1.3 < |\eta| < 3.0$. There is also the outer hadronic calorimeter (HO), which surrounds the solenoid magnet, and the forward hadronic calorimeter (HF), which exists beyond the muon detectors at high pseudorapidity.

The HB and HE are sampling calorimeters. They are composed of several absorbing layers, made from brass plates (70% Cu and 30% Zn), alternated with tiles of plastic scintillator. In the HB, the layers lie parallel to the beamline and in the HE, they lie parallel to the radial direction. Brass plates were chosen as they are radiation tolerant, non-magnetic, and because they have a relatively low nuclear interaction length [28] of $\lambda_I = 16.42$ cm, allowing the detectors to be compact. The energy

deposited in an absorbing layer is measured by collecting the light in the following scintillating tile using a wavelength shifting fibre. The scintillating tiles are aligned in $(\eta, \phi)$ space, forming the concept of the energy tower; the energy deposited in a segment of the calorimeter about an axis pointing to the interaction point. In the HB, the angular size of each tile is such that the energy towers map onto a $5 \times 5$ array of ECAL crystals, corresponding to $(\Delta\eta, \Delta\phi) = (0.087, 0.087)$. This angular tile size is maintained in the HE for $|\eta| < 1.7$. Beyond this, the tiles have a larger $\Delta\phi$ size of 0.175 and a larger, irregular, $\Delta\eta$ size. The energy resolution of the HB, in combination with the ECAL barrel, was measured during a test beam [33] to be:

$$\frac{\sigma}{E} = \frac{84.7\%}{\sqrt{E/\text{GeV}}} \oplus 7.4\% \tag{3.2}$$

The depth of the HB is restricted to 1.18 m, due to the radial constraints of the solenoid magnet. There are 16 absorbing layers in total. The first eight brass plates are 50.5 mm thick and the last six are 56.5 mm thick (the two boundary layers are made of stainless steel for structural strength). At $\eta = 0$, the total absorber thickness is 5.4 $\lambda_I$. The effective thickness rises as $|\eta|$ increases and at $|\eta| = 1.3$, the total absorber thickness is 10.3 $\lambda_I$. In addition, there is around 1.1 $\lambda_I$ of material provided by the ECAL crystals. In the HE, which does not have the same depth constraints, all the brass plates are 79 mm thick and there is one additional layer. This provides a total absorber thickness of around 10 $\lambda_I$.

The reason for the HO is that, in the low $|\eta|$ regions, the HB does not contain enough absorbing material to contain all hadronic showers. The HO is a sampling calorimeter, using the solenoid as the absorbing material and plastic scintillator tiles as the active material. In addition, for the region $|\eta| < 0.25$, the total absorber thickness is enhanced by a 19.5 cm thick layer of iron. The HO scintillating tiles roughly map onto the HB energy towers.

The HF begins at $|z| = 11.2$ m and covers the range $2.9 < |\eta| < 5.2$. It enables the reconstruction of forward jets and provides greater angular coverage for the $E_\text{T}^\text{miss}$

calculation. Because it exists at such high pseudorapidity, the HF is exposed to an extremely high particle flux. In order to withstand this, radiation-hard quartz fibres are used as the active material. The fibres, which have a diameter of 0.6 mm, run parallel to the beamline, embedded in a 165 cm thick steel absorber. The signal arises from Cherenkov light generated in the fibres which is guided to photomultipliers at the back of the absorber. Two different fibres lengths are used in order to distinguish between electromagnetic and hadronic showers. The fibres are placed in a grid formation in the transverse plane, each separated by 5 mm. Their outputs are grouped such that they form towers with size $(\Delta\eta, \Delta\phi) = (0.175, 0.175)$.

### 3.2.5 Muon Detectors

The muon detectors [34] are used to identify muons and measure their $p_{\mathrm{T}}$. They exist outside the solenoid magnet, interspaced between the iron return yoke structures. The muon detectors work by reconstructing the muon trajectories, which are bent by the return flux of the solenoid magnet. Using these trajectories, the muon $p_{\mathrm{T}}$ can then be determined. Due to the high strength of the solenoid magnet, good $p_{\mathrm{T}}$ resolution is achieved.

The total muon detector system is comprised of three subsystems, which can be seen in Figure 3.8. The drift tube (DT) system covers the barrel region and the cathode strip chamber (CSC) system covers the endcap region. Both are complemented by the resistive plate chamber (RPC) system. All three subsystems are gaseous particle detectors [28].

The DT system is composed of drift chambers. They are inserted between the iron return yoke barrels, creating four cylindrical detector layers and covering the region $|\eta| < 1.2$. The drift chambers contain three 'superlayers', where each superlayer is formed from four staggered layers of rectangular drift cells. The two outermost superlayers are aligned with the beamline, providing information about the muon's $(r, \phi)$ coordinates. In order to attain the best angular resolution, these superlayers are maximally separated, by about 20 cm, within each drift chamber. The inner

Figure 3.8: Diagram of CMS in the *r-z* plane with the DT, RPC, and CSC muon detectors labelled [4].

superlayer is orthogonal to the beamline, providing information about the $(r, z)$ coordinate. The drift cells are 42 mm wide, giving a maximum drift path of 21 mm. They are small enough to ensure that the occupancy is negligible but also large enough to restrict the number of channels to a manageable level.

The CSC system is composed of cathode strip chambers. They are inserted between the iron return yoke endcaps, creating a set of detector layers which are perpendicular to the beamline covering the region $0.9 < |\eta| < 2.4$. Each chamber is trapezoidal in shape and contains seven cathode strip planes, separated by 9.5 mm gas gaps, alternated with six anode wire planes. The cathode strips run in the radial direction and have a pitch that varies between 6.7 mm and 16.0 mm. They are used to measure the $(\phi, z)$ coordinates of the muons. The anode wires run almost orthogonally to the strips, with an inter-wire separation of 3.2 mm. They are used to measure the $(r, z)$ coordinates.

The RPC system is composed of parallel-plate chambers. The chambers are formed of two pairs of anode and cathode plates, each with a 2 mm gas gap, either side of a readout strip. The RPC system has an excellent time resolution but only a modest spatial resolution. There are six RPC detector layers in the barrel and three in the endcaps with $|\eta| < 1.6$.

### 3.2.6   CMS Simulation

Monte Carlo (MC) simulations of the LHC bunch crossings and the subsequent CMS detector response are critical in most CMS physics analyses. They are formed from a sequence of different simulation stages. First, the events are generated by simulating a given interaction between two colliding protons. This is followed by parton showering and hadronization simulations. The output particles are then simulated as they propagate through, and interact with, the CMS detector. This is achieved using a detailed model, implemented in GEANT4 [35], of the detector geometry and materials. In the next stage, each event is merged with the simulation of a set of pile-up interactions which, on average, reflect the conditions provided by

the LHC. Finally, the readout electronics are emulated, allowing the events to then be reconstructed like real data.

## 3.3 Particle-flow Reconstruction in CMS

CMS uses a particle-flow (PF) technique to reconstruct events [36]. All the detector layers are collectively used to identify, and measure the kinematic properties of, every stable final-state particle (except neutrinos). The PF algorithm is composed of three main parts. First, the basic PF elements are reconstructed by each individual detector. Next, the PF elements are linked together into sets, called PF blocks. Finally, the PF blocks are used to identify and reconstruct the final-state particles. These stages are outlined in Sections 3.3.1, 3.3.2, and 3.3.3, respectively.

### 3.3.1 Basic PF Elements

In the first part of the PF algorithm, each individual detector is used to reconstruct the basic PF elements required for particle identification and reconstruction:

- Using the inner tracker, the trajectories of the charged particles are reconstructed. This is achieved using multiple iterations of a combinatorial track finding algorithm [37]. Due to the tracker's high granularity, the tracks are reconstructed with high efficiency and have a low fake rate, although these features do degrade as the particle $p_{\mathrm{T}}$ increases.

  There are a few complications that are accounted for due to incident particles interacting with the material in the tracker: electrons can lose a significant fraction of their energy emitting bremsstrahlung photons; photons can convert into di-electron pairs; and hadrons can undergo nuclear interactions that alter their trajectories or create secondary particles.

- The energy deposits in the calorimeters are grouped geometrically into energy clusters. This is done, separately, for the pre-shower, ECAL, and HCAL by specific clustering algorithms. The energy clusters are calibrated as a function

of energy and pseudorapidity. The ECAL clusters have distinct calibrations depending on what particle is being reconstructed. This is because the ECAL has a significantly different response to electrons/photons and hadrons. Due to the high granularity of the ECAL clusters, the constituent particles of a jet can be distinguished from each other.

The energy deposits in the ECAL are additionally grouped into 'superclusters', which are used to collect the energy of an electron and the bremsstrahlung photons it emitted. Consequently, superclusters are narrow in the $\eta$ dimension but long in the $\phi$ dimension.

- Using the DT, CSC and RPC muons detectors, the muon trajectories outside the solenoid are reconstructed. The tracks are reconstructed with high efficiency across the whole acceptance. The muon purity is very high because the calorimeters absorb the other final-state particles (except neutrinos).

### 3.3.2   Linking Algorithm

In the second part of the PF algorithm, the basic PF elements are connected into sets, called PF blocks, by a linking algorithm. The algorithm works by matching the reconstructed trajectories and energy clusters spatially. The different types of links are as follows:

- Links between inner tracks and energy clusters. These are established if the extrapolated trajectory of a track coincides with an energy cluster in the preshower, ECAL, or HCAL. In order to gather the bremsstrahlung photons emitted by an electron, links are formed if the tangents of a trajectory, at each tracker layer, coincide with an ECAL cluster.

- Links between two inner tracks. If two tracks are compatible with photon conversion, a link is formed. If such a photon is found to be compatible with the bremsstrahlung of an electron track, further links are formed. Moreover,

in order to account for nuclear interactions in the tracker, links are established if two tracks share a secondary vertex.

- Links between two energy clusters. These are established when an ECAL cluster is contained within a HCAL cluster and when a preshower cluster is contained within an ECAL cluster. Additionally, ECAL clusters are linked to ECAL superclusters when they share a crystal.

- Links between inner tracks and muon detector tracks. These are established if the extrapolated trajectory of an inner track is compatible with a track in the muon detectors.

### 3.3.3 Particle Identification and Reconstruction

Once the PF blocks have been established, the particle identification and reconstruction algorithms are executed on each of them. The algorithm first attempts to identify muons, followed by electrons and isolated photons, then, finally, hadrons and non-isolated photons. If a particle is identified in a PF block, the associated PF elements are removed from it before continuing.

Muons are identified by the connection of inner tracks with muon detector tracks. If they are not isolated from other inner tracks or energy clusters, the muons are required to pass further selection criteria. These criteria eliminate accidental track associations and fakes caused by hadron showers that penetrate the muon detectors. For low energy muons, the $p_\text{T}$ measurement in the muon detectors is limited by multiple scattering effects. Consequently, for values less than 200 GeV, the muon $p_\text{T}$ is determined using the inner tracker only. Above this threshold, the $p_\text{T}$ is determined after finding the best track fit using both the inner tracker and muon detector.

Isolated photons are identified from ECAL superclusters without links to a full inner track. They are required to be isolated from other energy clusters or inner tracks. If there is a connected HCAL cluster, it must have a small energy. The photon is assigned the energy, and direction, of the supercluster.

Electrons are identified from ECAL clusters with a connected inner track. The entire supercluster and all the track tangent elements are considered which, due to their complexity, must satisfy further identification criteria. Additionally, if there is a connected HCAL cluster, it must have a small energy. The electron is assigned a direction based on the primary track. The energy is determined using a combination of the ECAL supercluster energy and the momentum of the primary track. For high energy electrons, this is dominated by the ECAL supercluster energy measurement because it has superior resolution.

The remaining particles to be identified are those associated to jets. These include charged hadrons, neutral hadrons, and non-isolated photons (from $\pi^0$ decays). There is no attempt made to identify the hadron species.

- Non-isolated photons are identified from the ECAL clusters without a link to a track and neutral hadrons are identified from the HCAL clusters without a link to a track. Both particle types are assigned the energy, and direction, of their associated calorimeter cluster. This approach means that the ECAL energy deposits of neutral hadrons are incorrectly reconstructed as photons and, thus, receive the wrong kind of calibration. This is deemed acceptable due to the small fraction of the total jet energy deposited in the ECAL by neutral hadrons.

- Charged hadrons are then identified from inner tracks connected to a HCAL cluster. Connected ECAL clusters are used in the reconstruction but are not necessary for identification. Due to the strong magnetic field, the charged hadrons are separated from the neutral particles in a jet, often creating a distinct set of energy clusters. In order to test this, the calorimeter cluster energies are compared to the momenta of the associated tracks. If there is an excess of calorimetric energy, it is attributed to additional photons and neutral hadrons. After this, the charged hadrons are reconstructed, one for each track. If there was a calorimetric energy excess, their momenta are determined solely by their track information. If there was no calorimetric energy excess,

each charged hadron is assigned a momentum during a fit involving all the associated tracks and energy clusters.

- Outside the tracker acceptance, the same distinctions cannot be made. The ECAL clusters without a HCAL connection are identified as photons whilst connected ECAL and HCAL clusters are identified as hadrons, without distinguishing whether they are charged or neutral.

## 3.4 Jets

### 3.4.1 Introduction to Jets

When partons (quarks or gluons) are produced in high energy interactions, they promptly fragment and hadronize, creating a collimated spray of hadrons. A jet is the object formed when one attempts to group the final-state particles into a set originating from such a parton. The energy and direction of a jet is, therefore, an attempted reconstruction of an outgoing parton.

Jets are not fundamental objects. They depend on the particle grouping algorithm used, of which there are many [38]. In CMS, the anti-$k_T$ algorithm [38] is used to define jets. It works by sequentially clustering pairs of particles as follows:

1. For each pair of particles $i$, $j$ calculate the 'distance' $d_{ij}$ given by:

$$d_{ij} = \min(p_{Ti}^{-2}, p_{Tj}^{-2}) \cdot \frac{\Delta R_{ij}^2}{R^2} \qquad (3.3)$$

where $R$ is a free parameter called the distance parameter. Additionally, for each individual particle $i$ calculate the 'distance' $d_{iB}$ given by:

$$d_{iB} = p_{Ti}^{-2} \qquad (3.4)$$

2. Find the minimum $d_{ij}$ or $d_{iB}$ value and identify the corresponding particles. If the minimum is of type $d_{ij}$, go to step 3. If it is of type $d_{iB}$, go to step 4.

3. Merge particles $i$ and $j$ into a single new particle by summing their four-momenta. Then return to step 1.

4. Declare particle $i$ to be a jet and remove it from the collection. If particle $i$ was the final remaining particle, then stop. If not, return to step 1.

The resultant jets develop outwards around high $p_T$ particles. They remain unchanged if an event is modified by a soft emission or a collinear splitting (i.e. the jets are infrared and collinear safe). The jets have a circular shape in $(\eta, \phi)$ space, intuitively mapping onto the cone shape one imagines for a parton shower. The jet area, in $(\eta, \phi)$ space, is very close to $\pi R^2$. This area is extremely stable and has minimal $p_T$ dependence. Therefore, the free parameter $R$ sets the angular size of the jets.

### 3.4.2 Jets at CMS

The jets most commonly used in CMS are reconstructed using the anti-$k_T$ algorithm with a distance parameter of $R = 0.4$. They are called AK4 jets. Typically, detectors at hadron colliders reconstruct jets using calorimeter towers as the inputs. At CMS, this technique would be limited by the coarse segmentation and modest energy resolution of the HCAL. A significant improvement is achieved because CMS is able to perform PF reconstruction, allowing PF particles, rather than calorimeter towers, to be used as inputs to the jet algorithm. It has the following benefits:

- The charged hadrons, which carry around 65% of a jet's energy [36], are measured using the inner tracker. This provides far superior energy and spatial resolution compared to solely using the HCAL. Furthermore, nuclear interactions in the tracker are accounted for.

- The photons (from $\pi^0$ decays), which carry around 25% of a jet's energy [36], are distinguished from hadrons. Consequently, an excellent energy resolution is achieved because the corresponding ECAL energy deposits can be calibrated

specifically for photons. Additionally, photon conversion in the tracker is accounted for.

- Only around 10% of the jet energy, that carried by neutral hadrons [36], has to be measured using the HCAL directly.

- If one of the hadrons decays into an electron or muon, it will be incorporated into the jet.

Pile-up interactions produce additional quarks and gluons, which then undergo showering and hadronization, superimposed on the hard scatter (the high $p_\mathrm{T}$ parton-parton interaction). The corresponding tracks and energy deposits degrade the jet reconstruction, especially in the high pile-up environment provided by the LHC. In CMS, there are two main methods used to try and mitigate pile-up effects; charged hadron subtraction (CHS) and pile-up per particle identification (PUPPI). Both methods are made possible because of PF reconstruction.

The CHS technique [39, 40] removes all the reconstructed charged hadrons associated to pile-up vertices from the event before applying the jet algorithm. This eliminates a sizeable fraction of the pile-up energy from the resultant jets. The remaining pile-up energy in each jet, due to photons (from $\pi^0$ decays) and neutral hadrons, is then estimated and subtracted from the jet. This energy contamination is estimated by multiplying the jet area by the average $p_\mathrm{T}$ per unit area due to neutral pile-up.

The PUPPI technique [41] rescales the four-momentum of each PF particle before applying the jet algorithm. The rescaling weights vary between zero, for particles arising from pile-up, and unity, for particles arising from the hard scatter. For charged hadrons the weight assignment is trivial due to the vertex information. For neutral particles it is more complex. First, a local shape variable, which discriminates between pile-up and hard scatter particles, is calculated for each PF particle. Next, the distribution of this variable, in both the pile-up and hard scatter scenarios, is determined using the charged hadrons in the event. Finally, these pieces of information are combined to assign a non-integer weight to each neutral particle.

The PUPPI technique is especially useful when investigating jet substructure (see Section 3.4.3) because it removes, or suppresses, neutral pile-up particles in the jets, rather than just applying an energy correction.

Following the pile-up mitigation, jet energy corrections (JECs) are applied to the jets, in order to calibrate them to the correct energy scale [40]. The JECs are determined, as a function of jet $p_T$ and $|\eta|$, using QCD multi-jet MC. The reference jets are formed by clustering all the stable output particles (except neutrinos) at the generator level. After the primary JECs, residual corrections are applied to jets in data to account for the differences between data and simulation. Finally, the jets are required to pass a set of identification criteria designed to reject fake jets arising from instrumental effects.

### 3.4.3 Fat Jets and Substructure

Quark anti-quark pairs become collimated when they arise from the decay of a boosted particle. For a parent particle of mass $M$ and transverse momenta $p_T$, the q$\overline{\text{q}}$ angular separation follows a falling distribution with a minimum, and mode, of $\Delta R \approx 2M/p_T$ [42]. In scenarios where $\Delta R < 0.4$, the AK4 jet algorithm will not resolve the two partons. One solution, when working with boosted topologies, is to actively try and reconstruct both partons in a single 'fat-jet' by clustering with a larger distance parameter. These jets can then be distinguished from background due to their two prong substructure. Consequently, CMS also reconstructs anti-$k_T$ jets with a distance parameter of $R = 0.8$. They are called AK8 jets. Following their reconstruction, pile-up mitigation and JECs are applied to the AK8 jets as described above for AK4 jets.

In the remainder of this section, a signal jet is defined as an AK8 jet which reconstructs the q$\overline{\text{q}}$ pair from the decay of a boosted massive particle (e.g. a Higgs boson) and a background jet is defined as an AK8 jet in which a single parton is reconstructed. Jet mass, the invariant mass of a jet's constituent particles, is an important discriminating feature between the signal and background jets. At the parton level,

signal jets have a mass equal to the parent particle mass and background jets have a mass of zero. This distinction is obscured, however, by soft gluon emission. It significantly broadens the signal jet mass peak and causes the background jets to attain a mass which is, on average, proportional to their momentum [43]. In addition, the jet masses are increased by coincident pile-up particles, however, this is largely mitigated using the PUPPI technique.

Jet grooming techniques attempt to remove the soft and unassociated radiation whilst retaining the hard underlying substructure. Therefore, the jet mass, evaluated after a perfect grooming algorithm, would be zero for a background jet and equal to the parent particle mass for a signal jet (assuming perfect particle momentum resolution). There are multiple jet grooming techniques [44] but, currently, the most prevalent method in CMS is the soft-drop algorithm [45] with parameters $\beta = 0$ and $z_{\mathrm{cut}} = 0.1$.

The soft-drop algorithm works by recursively removing soft wide-angle radiation. The algorithm begins by re-clustering a jet's constituent particles using the Cambridge-Aachen algorithm [38]. This sequentially clusters pairs of particles with the smallest angular separation. The algorithm then de-clusters the jet as follows:

1. Undo the previous stage of the Cambridge-Aachen clustering, dividing jet $j$ into sub-jets $j_1$ and $j_2$.

2. Declare $j$ as the final jet (and then stop) if the sub-jets pass the condition:

$$\frac{\min(p_{T1}, p_{T2})}{p_{T1} + p_{T2}} > z_{\mathrm{cut}} \left( \frac{\Delta R_{1,2}}{R} \right)^{\beta} \tag{3.5}$$

3. Re-define $j$ as the sub-jet with the larger $p_{\mathrm{T}}$. If $j$ is now comprised of a single particle, declare $j$ as the final jet. Otherwise, return to step 1.

It is the failure to pass the condition in Equation 3.5 which results in the removal of the soft wide-angle radiation. In the CMS soft-drop configuration, where $\beta = 0$ and $z_{\mathrm{cut}} = 0.1$, a sub-jet is removed if it has less than 10% of the total (de-clustered)

jet $p_\mathrm{T}$. As will be seen in Section 6.2.3, this configuration is very successful at grooming background jets. As a consequence of this, however, it is slightly too aggressive with signal jets, and can incorrectly remove hard PF particles arising from the $\mathrm{q\bar{q}}$ decays.

### 3.4.4   B-Tagging

When b-quarks are produced in high energy interactions the resultant spray of hadrons will contain a b-hadron. Due to their relatively large mass ($\approx 5$ GeV [7]), b-hadrons carry away a significant fraction of the original b-quark momentum. The b-hadrons then decay at a point, called the secondary vertex, which is displaced with respect to the primary vertex. On average, the decays produce five charged particles [46]. With a lifetime of approximately $\tau \approx 1.5$ ps [7], the b-hadron decay length is:

$$
\begin{aligned}
L_\mathrm{b} &= c\tau_\mathrm{b}\beta_\mathrm{b}\gamma_\mathrm{b} \\
&= c\tau_\mathrm{b}\frac{|\vec{p_\mathrm{b}}|}{m_\mathrm{b}} \\
&\approx \frac{|\vec{p_\mathrm{b}}|}{\mathrm{GeV}}\ 0.1\ \mathrm{mm}
\end{aligned}
\tag{3.6}
$$

This is sufficiently large to enable the reconstruction of the secondary vertex using the charged tracks emerging from the b-hadron decay.

Using these b-hadron properties, jets that originate from b-quarks (b-jets) can be distinguished from those originating from lighter flavour quarks. This is done by b-tagging algorithms which assign to each jet a numerical discriminator representing the likelihood that it is a b-jet. There are multiple b-tagging algorithms employed by CMS [47]. The most powerful ones use multivariate techniques which combine information about both the displaced tracks and secondary vertices associated to a jet axis; the direction of a jet's momentum vector.

A special case of interest is when a boosted massive particle (e.g. a Higgs boson) decays into a $\mathrm{b\bar{b}}$ pair. As was discussed in Section 3.4.3, such topologies can be

reconstructed in a single fat-jet and then the (soft-drop) mass used to distinguish against background. In this scenario, however, further distinctions can be made due to the presence of the two b-hadrons. During Run-I, this feature was exploited using the standard b-tagging algorithms, by either applying them to the sub-jets or to the entire fat-jet. For Run-II, a specific double-b-tagging algorithm was developed [48]. In addition to using the standard b-tagging information, the algorithm exploits the two prong substructure of the jets.

The jet axis is an important property in the standard b-tagging algorithms. In accordance, the substructure variables used in the double-b-tagger are the sub-jet axes; where the two sub-jets are obtained by re-clustering a jet's constituent particles, using the $k_T$ algorithm [38], and then undoing the final step. The sub-jet axes are useful because they are each strongly correlated to the flight direction of one of the b-hadrons. No other substructure quantities are used in the double-b-tagger. The jet mass is not used in order to avoid a (strong) mass dependency.

There are 27 discriminating variables that are input into the multivariate discriminant, which is trained using a boosted decision tree (BDT). The variables rely on a jet's displaced tracks and secondary vertices, often in association with the sub-jet axes. In order to keep the double-b-tagger as general as possible, none of the variables used have a strong dependence on the jet $p_{\mathrm{T}}$ or mass.

### 3.4.5 Event Variables using Jets

There are two common variables that can be assigned to each event; $H_{\mathrm{T}}$ and $\vec{H}_{\mathrm{T}}^{\mathrm{miss}}$. These variables are calculated from all the jets in an event with $p_{\mathrm{T}}$ greater than a given threshold (typically around 40 GeV) and $|\eta|$ less than a given value (typically around 3.0). The $H_{\mathrm{T}}$ quantity is the scalar sum of the $p_{\mathrm{T}}$ of the selected jets and $\vec{H}_{\mathrm{T}}^{\mathrm{miss}}$ is the negative vector sum of the $\vec{p_{\mathrm{T}}}$ of the selected jets.

# 4

# Jets in the Level-1 Trigger

This chapter begins with an overview of the entire CMS trigger system. It then focuses on the Level-1 calorimeter trigger, which was upgraded for Run-II, in order to facilitate a discussion about jet reconstruction. This discussion includes the jet performance studies I conducted following the trigger upgrade and a description of the jet energy corrections I derived in 2016 and 2017. At the end of the chapter, two Level-1 trigger problems relevant to the main analysis are outlined.

## 4.1   Introduction to the CMS Trigger System

The CMS detector cannot record on disk all the events provided by the LHC. This is because the full CMS detector readout is of order 1 MB per event and because the LHC provides proton bunch collisions at a rate of 40 MHz. The maximum rate at which events can be archived is around 1000 Hz. Consequently, only one in every 40,000 events can be recorded. The CMS trigger is the system that selects these events. It uses a simplified detector readout to determine whether an event has the characteristics of a physics process of interest. The goal of the trigger is to select interesting events with the highest possible efficiency. The physics processes targeted by the CMS experiment, and hence the trigger system, are electroweak scale

SM physics and new physics at the TeV scale. The CMS trigger system reduces the rate, by selecting events, in two steps. The first step is performed by the Level-1 trigger and the second step is performed by the High-Level trigger (HLT).

## 4.1.1 Level-1 Trigger

The first step of the event selection is performed by the Level-1 trigger [49]. This system is made out of custom hardware. It analyses all the events, which occur at a rate of 40 MHz, and reduces the rate down to 100 kHz. While the Level-1 trigger operates, the full detector readout is held in pipeline memories in the front-end electronics. If the Level-1 trigger accepts an event, the full detector readout is sent to the HLT.

The latency of the Level-1 trigger is 3.2 µs. This is set by the maximum time for which the full detector readout can be held in pipeline memory. In this 3.2 µs time period, the Level-1 trigger must read in data from the event, execute its algorithms, and send the event selection decision to the front-end electronics of the detectors. The time constraint limits the sophistication of the physics objects that can be reconstructed in the Level-1 trigger. The inner tracker cannot be used due to the volume and complexity of the data it collects. Only coarse-grained information from the calorimeters and the muon detectors can be used. Furthermore, the algorithms executed on this data must be relatively simple compared to those used in the offline reconstruction.

There are three main parts to the Level-1 trigger system; the calorimeter trigger, the muon trigger, and the global trigger. The calorimeter trigger uses information from the ECAL, HCAL, and HF detectors to reconstruct e/$\gamma$, tau, and jet objects. It is also used to calculate energy sum quantities such as $E_T^{miss}$. The muon trigger, which is used to reconstruct muon objects, uses information from the RPC, DT, and CSC detectors. The global trigger receives the physics objects reconstructed in the calorimeter and muon triggers. It determines if these objects pass any of the event selection criteria and then sends the decision to the front-end electronics.

The event selection criteria implemented in the global trigger require the physics objects, both individually and in combinations, to be above certain $E_T$ thresholds. The angular information of the physics objects can also be used. For example, two objects can be required to be close, or far apart, from each other. Each event selection criterion represents at least one physics process that is of interest to CMS. Lower $E_T$ thresholds correspond to greater acceptance, but also to higher trigger rates. It is a balancing act to accommodate the broad CMS physics programme whilst maintaining a total event selection rate below 100 kHz.

## 4.1.2 High-Level Trigger

The second step of the event selection is performed by the HLT [50]. The system is implemented on a farm of commercial multi-processors. It reduces the rate down from 100 kHz to less than 1000 Hz, the maximum rate at which events can be written to disk. The HLT creates physics objects by partially reconstructing the full CMS detector readout. Like the Level-1 trigger, it selects events depending on whether the reconstructed objects meet certain criteria.

The physics objects that the HLT reconstructs are motivated by the candidate objects, and their locations, identified in the Level-1 trigger. This saves computational time, as the HLT knows which algorithms to execute and which regions of the detectors to use. The HLT saves further computational time by reconstructing the objects in stages and filtering events as it progresses. As the stages advance, the reconstruction complexity increases. The initial objects only use data from the calorimeter and muon detectors. The inner tracker is not used, due to the complexity of the data it contains. Only in the later stages of an objects reconstruction is the tracker information integrated. Once the final filter is applied, the quality of a physics object is comparable to that after full offline reconstruction.

## 4.1.3   Upgrade to the Level-1 Trigger

The Level-1 trigger was upgraded in CMS for Run-II of the LHC [5]. It was commissioned in 2015 and operational from 2016. It will also be used for Run-III of the LHC. The calorimeter trigger, muon trigger, and global trigger were all redesigned.

The motivation for the Level-1 trigger upgrade was the expectation that, during Run-II, the instantaneous luminosity delivered by the LHC would increase beyond the initial design specification. This was indeed the case. The LHC was designed to provide an instantaneous luminosity of $1.0 \times 10^{34}$ cm$^{-2}$s$^{-1}$, however, in 2018, the instantaneous luminosity delivered was typically around $1.9 \times 10^{34}$ cm$^{-2}$s$^{-1}$, with peak values of $2.1 \times 10^{34}$ cm$^{-2}$s$^{-1}$. This would have resulted in double the Level-1 trigger rates, for the same $E_{\mathrm{T}}$ thresholds, had the Level-1 trigger not been upgraded.

The expected increase in instantaneous luminosity would also intensify the pile-up conditions. The LHC design specification estimated an average of 25 pile-up interactions per bunch crossing, however, during 2018, the mean was 37 and some events had more than 60. Pile-up interactions deposit their energy in the calorimeters. If this additional energy is not accounted for when reconstructing physics objects, the energy resolution will degrade as the number of pile-up interactions increases. The legacy Level-1 trigger could not support algorithms that included pile-up subtraction methods. Consequently, the energy resolution of the Level-1 physics objects was expected to deteriorate in Run-II. This provided additional motivation to upgrade the Level-1 trigger, so that it could support algorithms that accounted for pile-up interactions.

It should also be noted that the pp centre-of-mass energy was increased from 8 TeV to 13 TeV for Run-II. This change was also going to result in higher trigger rates, because it increased the probability of high $E_{\mathrm{T}}$ events occurring. However, this was not the main motivation for upgrading the trigger because centre-of-mass energies of 14 TeV were expected in Run-I.

## 4.2 Level-1 Calorimeter Trigger

This section explains how the Level-1 calorimeter trigger works in more detail. First, the data that is provided from the calorimeters is explained. This is followed by a description of the legacy Level-1 calorimeter trigger used by CMS. This system is representative of a traditional calorimeter trigger. In discussing it, one gains a clearer insight into the improvements achieved by the upgraded Level-1 calorimeter trigger, which is discussed in the final subsection.

### 4.2.1 Calorimeter Trigger Primitives

The inputs provided to the Level-1 calorimeter trigger are called trigger primitives [49]. They are supplied from both the ECAL and HCAL detectors. The trigger primitives correspond to different sections of the calorimeters in $(\eta, \phi)$ space. These sections are called trigger towers.

The trigger primitives represent the $E_\mathrm{T}$ deposited in each trigger tower. They are generated by a sub-system that interfaces with the front-end electronics of the ECAL and HCAL detectors. For the ECAL trigger primitives, the $E_\mathrm{T}$ of a trigger tower is encoded in an 8-bit quantity. In addition, there is a feature bit which flags whether the $E_\mathrm{T}$ deposits are compatible with those typically left by electrons or photons. For the HCAL trigger primitives, the $E_\mathrm{T}$ of a trigger tower is also encoded in an 8-bit quantity. Additionally, there is a feature bit which flags energy deposits caused by minimum ionizing particles.

In the barrel calorimeters, each trigger tower has an angular size of $(\Delta\eta, \Delta\phi) = (0.087, 0.087)$. This corresponds to a $5 \times 5$ array of ECAL crystals and one physical HCAL tower. As $|\eta|$ increases beyond 1.740, which is part way through the endcap calorimeters, the trigger towers become larger in $\Delta\eta$. This is so that the trigger towers continue to map onto the $\eta$ coordinates of the physical HCAL towers. Additionally, in the $|\eta| > 1.740$ region, the physical HCAL towers become twice as large in $\Delta\phi$. The azimuthal angular size of the trigger towers remains the same, however, with two trigger towers mapping on to one physical HCAL tower. The $E_\mathrm{T}$ in the

physical HCAL tower is divided equally between the two trigger towers. The sizes of all the trigger towers, and the detectors that they correspond to, are provided in Table 4.1.

In total, the Level-1 calorimeter trigger receives 9936 trigger primitives. The ECAL provides 4032 of the trigger primitives, corresponding to the 2448 and 1584 trigger towers in the EB and EE, respectively. The HCAL provides 5904 of the trigger primitives, corresponding to the 2304, 1728 and 1872 trigger towers in the HB, HE, and HF, respectively.

## 4.2.2   Legacy Calorimeter Trigger

The Level-1 calorimeter trigger used by CMS during Run-I [49] is representative of a traditional calorimeter trigger system. It is composed of two main parts; the Regional Calorimeter Trigger (RCT) and the Global Calorimeter Trigger (GCT). The RCT processes the trigger primitive data in parallel, with each processing node handling data from a different array of trigger towers. The GCT then receives all the objects created by the RCT. It applies some further algorithms, sorts the objects by category, and then sends the final set of objects to the global trigger.

A critical constraint when designing the legacy Level-1 trigger, was that each individual processing step had to occur within 25 ns, the time period between bunch crossings. Due to this constraint, it was not possible to transfer all the trigger primitive data onto a single processing node. This is why the trigger primitive data is partitioned into distinct regions. Each RCT node receives 64 trigger primitives, 32 from both the ECAL and HCAL, corresponding to a $4 \times 8$ array of trigger towers. These are divided into two $4 \times 4$ trigger tower arrays, on which two main operations are conducted. One operation sums the total $E_\mathrm{T}$. These energies are used by the GCT to cluster jets and calculate energy sums quantities. The other operation finds the four best isolated, and non-isolated, $\mathrm{e}/\gamma$ candidates. In order to achieve continuous coverage, this algorithm requires data sharing across nodes which process neighbouring trigger tower regions. After the RCT stage, the information

Table 4.1: The different trigger towers, their sizes, and the calorimeters that they correspond to. Note that $\Delta\phi = 5°$ for all trigger towers.

| Tower | $|\eta|$ range | | Size | Calorimeter | |
|---|---|---|---|---|---|
| $|\eta|$ index | Low | High | $\Delta\eta$ | ECAL | HCAL |
| 1 | 0.000 | 0.087 | 0.087 | EB | HB |
| 2 | 0.087 | 0.174 | 0.087 | EB | HB |
| 3 | 0.174 | 0.261 | 0.087 | EB | HB |
| 4 | 0.261 | 0.348 | 0.087 | EB | HB |
| 5 | 0.348 | 0.435 | 0.087 | EB | HB |
| 6 | 0.435 | 0.522 | 0.087 | EB | HB |
| 7 | 0.522 | 0.609 | 0.087 | EB | HB |
| 8 | 0.609 | 0.696 | 0.087 | EB | HB |
| 9 | 0.696 | 0.783 | 0.087 | EB | HB |
| 10 | 0.783 | 0.870 | 0.087 | EB | HB |
| 11 | 0.879 | 0.957 | 0.087 | EB | HB |
| 12 | 0.957 | 1.044 | 0.087 | EB | HB |
| 13 | 1.044 | 1.131 | 0.087 | EB | HB |
| 14 | 1.131 | 1.218 | 0.087 | EB | HB |
| 15 | 1.218 | 1.305 | 0.087 | EB | HB |
| 16 | 1.305 | 1.392 | 0.087 | EB | HB, HE |
| 17 | 1.392 | 1.479 | 0.087 | EB | HE |
| 18 | 1.479 | 1.566 | 0.087 | EE | HE |
| 19 | 1.566 | 1.653 | 0.087 | EE | HE |
| 20 | 1.653 | 1.740 | 0.087 | EE | HE |
| 21 | 1.740 | 1.830 | 0.090 | EE | HE |
| 22 | 1.830 | 1.930 | 0.100 | EE | HE |
| 23 | 1.930 | 2.043 | 0.113 | EE | HE |
| 24 | 2.043 | 2.172 | 0.129 | EE | HE |
| 25 | 2.172 | 2.322 | 0.150 | EE | HE |
| 26 | 2.322 | 2.500 | 0.178 | EE | HE |
| 27 | 2.500 | 2.650 | 0.150 | EE | HE |
| 28 | 2.650 | 3.000 | 0.350 | EE | HE |
| 29 | 2.853 | 2.964 | 0.111 | - | HE, HF |
| 30 | 2.964 | 3.139 | 0.175 | - | HF |
| 31 | 3.139 | 3.314 | 0.175 | - | HF |
| 32 | 3.314 | 3.489 | 0.175 | - | HF |
| 33 | 3.489 | 3.664 | 0.175 | - | HF |
| 34 | 3.664 | 3.839 | 0.175 | - | HF |
| 35 | 3.839 | 4.013 | 0.174 | - | HF |
| 36 | 4.013 | 4.191 | 0.178 | - | HF |
| 37 | 4.191 | 4.363 | 0.172 | - | HF |
| 38 | 4.363 | 4.538 | 0.175 | - | HF |
| 39 | 4.538 | 4.716 | 0.178 | - | HF |
| 40 | 4.716 | 4.889 | 0.173 | - | HF |
| 41 | 4.889 | 5.191 | 0.302 | - | HF |

content has been reduced sufficiently such that all the different regions of data can be brought together, and further processed, in the GCT.

There are a number of problems with the regional architecture of the legacy calorimeter trigger. One problem is the data sharing required between processing nodes in the RCT. It introduces a significant overhead and ultimately limits the algorithms that can be performed. Another problem is that the trigger primitive data used in the RCT is physically distinct from the physics objects in the GCT, limiting the flexibility of the algorithms. A final problem is that, due to the required reduction in information, the GCT does not cluster jets using the full trigger tower granularity. All these problems are addressed by the time-multiplexed architecture of the upgraded calorimeter trigger.

### 4.2.3   Time-Multiplexed Trigger

The Time-Multiplexed Trigger [5] (TMT) is the current Level-1 calorimeter trigger system. It was introduced, as part of the Level-1 trigger upgrade, during Run-II. The TMT enables the physics object algorithms to be executed, using all the trigger primitive data, on a single processing node. This is achieved by the time-multiplexed architecture of the system, which provides individual steps with processing times greater than a single bunch crossing period.

The TMT is composed of two main parts; Layer-1 and Layer-2. Layer-1 pre-processes the ECAL and HCAL trigger primitives. It consists of a set of cards, where each card maps onto a different strip of trigger towers. Layer-2 performs the physics object algorithms [51]. It also consists of multiple cards. For a given bunch crossing, $N$, all the Layer-1 cards transmit their trigger tower data to a single Layer-2 card. For the next bunch crossing, $N + 1$, the Layer-1 cards transmit their data to a different Layer-2 card. This continues until bunching crossing $N + 9$, where the Layer-1 data is transmitted to the original Layer-2 card. It is this time-multiplexed architecture that provides a sufficiently large time window to transfer all the trigger primitive data onto a single processing node.

An illustration of the TMT architecture is provided in Figure 4.1. The system utilizes two cards; the CTP7 and the MP7. Both cards have a Field-Programmable Gate Array (FPGA) in which the logic is programmed. This allows the algorithms to be reconfigured, which provides flexibility as one can respond to changing conditions or implement new techniques. Layer-1 uses 18 CTP7 cards. Each card receives the trigger primitives, from both the ECAL and HCAL, corresponding to an array of trigger towers spanning all of the $\eta$ dimension and four trigger towers in the $\phi$ dimension. Layer-2 uses 9 MP7 cards. They are connected to the Layer-1 cards via an optical patch panel, which facilitates the time-multiplexing.

The trigger tower energies are streamed into a Layer-2 card one ($\eta$-bin) row at a time. The card operates on the data as it is fed in, scanning across the trigger tower row. Consequently, all the physics object algorithms are reduced to 1D tasks and can be fully pipelined. An example is provided in Section 4.3.1, where the jet algorithm is described. Operating in 1D allows the FPGA to be laid out sequentially. This increases the speed and reduces the routing congestion.

The upgrade to the Level-1 calorimeter trigger has improved the quality of the physics objects reconstructed. Better spatial resolution is achieved because the algorithms now use the full trigger tower granularity. The energy resolution is improved due to more sophisticated algorithms that have access to all the trigger tower data. For example, all the objects now have pile-up subtraction, which none of them had in the legacy trigger. It should also be noted that the TMT has a dedicated tau algorithm [52], something that could not be implemented in the legacy calorimeter trigger.

The improvement in quality of the physics objects means that they can be selected, at a given $E_\mathrm{T}$ threshold, more efficiently. As a consequence, the upgraded Level-1 trigger accepts events at a lower rate than the legacy trigger whilst maintaining the same selection efficiency. This can be seen in Ref. [53] for e/$\gamma$ objects, and in Ref. [6] for jets and energy sums.

Figure 4.1: Schematic diagram of the time-multiplexed architecture in the upgraded Level-1 calorimeter trigger [5].

## 4.3 Jets and Energy Sums

This section reviews jets, and the energy sum quantities derived from them, in the upgraded Level-1 calorimeter trigger. First, the jet, $H_\mathrm{T}$, and $H_\mathrm{T}^\mathrm{miss}$ algorithms are described. The performance of these objects are then examined from data taken during the first few months of operations using the upgraded Level-1 trigger.

### 4.3.1 Algorithms

All the algorithms in the Level-1 calorimeter trigger are reduced to 1D tasks. The jet algorithm [6] works by looping through the individual trigger towers, one row at a time. Each trigger tower is considered as the centre of a candidate Level-1 jet, which is composed of the surrounding $9 \times 9$ array of trigger towers. Consequently, the size of a Level-1 jet, in the barrel, is $(\Delta\eta, \Delta\phi) = (0.783, 0.786)$. Despite being square shaped, this approximately matches an AK4 jet. In the endcap and HF, the Level-1 jet $\Delta\eta$ size increases because the trigger tower size increases for $|\eta| > 1.740$. The direction of a Level-1 jet is given by the $(\eta, \phi)$ coordinates of the central trigger tower.

A jet is reconstructed around a trigger tower if its total $E_\mathrm{T}$ meets the following two requirements:

- The $E_\mathrm{T}$ must be greater than a value called the jet seed threshold. This was set to 4 GeV for the duration of Run-II.

- The $E_\mathrm{T}$ must be a local maximum within the surrounding $9 \times 9$ array of trigger towers. This prevents double counting of jets as the algorithm loops through the trigger towers. As is illustrated in Figure 4.2, the local maxima conditions are asymmetric along the diagonal of the $9 \times 9$ array. This prevents trigger towers with the same $E_\mathrm{T}$ from vetoing each other. Note that, at high pseudorapidities, a full $9 \times 9$ array cannot be considered around a central trigger tower. In these cases, only the remaining subset of trigger towers are used.

Figure 4.2: The $E_T$ requirement of a trigger tower (green box), relative to its neighbouring trigger towers (purple and blue boxes), in order for a Level-1 jet to be reconstructed about it [6].

The reconstructed jets are initially assigned an $E_T$ equal to the sum of all the trigger tower energies in the $9 \times 9$ array. The total $E_T$ of each trigger tower is formed from the sum of its corresponding ECAL and HCAL trigger primitives. Both trigger primitives are 8-bit quantities which are linearly converted into physical energies using a least-significant-bit, set to equal 0.5 GeV. Thus, the $E_T$ of each trigger primitive ranges between 0 and 127.5 GeV in 0.5 GeV intervals. The $E_T$ of a Level-1 jet, the sum of 81 trigger tower energies, is represented by an 11-bit quantity with a least-significant-bit of 0.5 GeV. Thus, the available $E_T$ values range from 0 to 1023.5 GeV in 0.5 GeV steps. The Level-1 jets are given a maximum $E_T$ of 1023.5 GeV if the sum of the trigger tower energies exceed this value or if any of the trigger primitives have saturated energies. In the latter case, it means that the $E_T$ of the jets can be grossly overestimated, however, it ensures that the object will pass the single jet trigger.

Following the initial Level-1 jet energy assignment, the $E_T$ contribution from pile-up is estimated using the so-called chunky-donut algorithm [54]. The algorithm first calculates the total $E_T$ in the four $3 \times 9$ arrays of trigger towers sharing an edge with the jet, as illustrated in Figure 4.3. These quantities evaluate the $E_T$ density in the calorimeters, due to pile-up interactions, in the proximity of the jet. They can be overestimated if, close to the jet, there is another high $E_T$ object from the

Figure 4.3: The four $3 \times 9$ arrays of trigger towers (pink boxes) considered around a Level-1 jet (grey boxes) by the chunky-donut algorithm [6].

hard scatter. In order to avoid this, the highest of the four energy sums is omitted. The three remaining are then subtracted from the $E_T$ of the jet. The energies can be subtracted directly, without rescaling, as there are 81 trigger towers in the three $3 \times 9$ arrays, the same number of which a jet is comprised.

Following the application of the chunky-donut algorithm, jet energy corrections (JECs) are applied to the Level-1 jets, completing their reconstruction. The derivation of the JECs, and how they are implemented, is described in Section 4.4. The $H_T$ and $H_T^{\text{miss}}$ energy sum quantities are then calculated using the jets with $E_T > 30$ GeV and $|\eta| < 2.4$. If there are any saturated trigger primitives in the event then both $H_T$ and $H_T^{\text{miss}}$ are given a maximum energy of 2047.5 GeV, to ensure the event is selected. Following the calculation of the energy sum quantities, all the Level-1 trigger objects are sent to the global trigger.

Having described the jet reconstruction algorithm in the upgraded Level-1 trigger, it is worth highlighting the improvements relative to the legacy Level-1 trigger:

- **Position resolution:** Jets in the upgraded Level-1 trigger have a granularity equal to the size of a trigger tower. In the legacy Level-1 trigger, it was four times larger in both the $\eta$ and $\phi$ dimensions. This is because the jet algorithm was applied to trigger towers already clustered into $4 \times 4$ arrays.

- **Pile-up mitigation:** Jets in the upgraded Level-1 trigger have pile-up subtraction applied by the chunky-donut algorithm. In the legacy Level-1 trigger, there was no pile-up subtraction as such algorithms could not be implemented.

- **Energy quantization:** The $E_{\mathrm{T}}$ of jets in the upgraded Level-1 trigger is represented by an 11-bit quantity. Thus, for a maximum energy of 1023.5 GeV, the $E_{\mathrm{T}}$ is quantized in 0.5 GeV steps. In the legacy Level-1 trigger, the jet $E_{\mathrm{T}}$ was an 8-bit quantity. Consequently, for the same maximum energy, the $E_{\mathrm{T}}$ was quantized in 4.0 GeV steps.

## 4.3.2 Performance

The upgraded Level-1 trigger was operational from 2016. As the first collision data was recorded, it was imperative to check the trigger's performance and ensure it was operating successfully. One test was to compare the physics objects reconstructed in the Level-1 trigger with the equivalent objects reconstructed offline. This allows one to evaluate the position and energy resolution of the Level-1 objects. It also allows one to assess the efficiency at which the Level-1 objects are reconstructed. These studies were conducted regularly during the first few months of operations. They were summarised at the ICHEP conference in August 2016 [55]. In the remainder of this section, the results for jets and $H_{\mathrm{T}}$ are presented.

The data sample used for these studies was collected using a single muon trigger, providing a set of events unbiased by the jet trigger. The Level-1 jets were evaluated against offline AK4 jets reconstructed from PF particles. Only one offline jet was selected per event; the highest $E_{\mathrm{T}}$ jet passing additional lepton veto requirements. The Level-1 jet with the smallest $\Delta R$ separation from the offline jet was selected for the comparison. To ensure a good spatial matching, it was required to have $\Delta R < 0.3$ from the offline jet.

Level-1 jet turn-on curves, for jets with $|\eta| < 3.0$, are provided in Figure 4.4. The graphs quantify how efficiently single Level-1 jets, above a given $E_{\mathrm{T}}$ threshold, are reconstructed as a function of the offline jet $E_{\mathrm{T}}$. The perfect shape would be a step

Figure 4.4: Efficiency of reconstructing single jets in the Level-1 trigger, above a variety of $E_\mathrm{T}$ thresholds, as a function of the offline jet $E_\mathrm{T}$.

function, rising from zero to unity at the point where the offline jet $E_\mathrm{T}$ equals the Level-1 jet threshold. This would allow the selection of all desired jets, without wasting any readout bandwidth. However, due to the imperfect $E_\mathrm{T}$ resolution of the Level-1 jets, this step function is smeared into a shape resembling an error function. The worse the $E_\mathrm{T}$ resolution is, the larger the width of the error function. If the average $E_\mathrm{T}$ of the Level-1 jets is correct, the error function is centred on the corresponding $E_\mathrm{T}$ threshold. The centre is translated as the Level-1 jet energy scale becomes more inaccurate. In Figure 4.4, all the turn-on curves reach unity and are approximately centred on their corresponding $E_\mathrm{T}$ threshold. In addition, the widths of the error functions are not too large. For example, if the Level-1 jet threshold is 200 GeV, there is a 95% probability of accepting an offline jet with an $E_\mathrm{T}$ of 240 GeV. This study provided evidence that the Level-1 trigger was successfully reconstructing jets.

Figure 4.5 shows the $E_\mathrm{T}$ resolution of Level-1 jets with respect to the offline jets that they were matched to. The offline jets were required to have $E_\mathrm{T} > 30$ GeV and $|\eta| < 3.0$. The $E_\mathrm{T}$ resolutions are provided for three different pile-up scenarios, corresponding to the low, medium, and high pile-up conditions at the beginning of

Figure 4.5: The $E_T$ resolution of Level-1 jets with respect to the matched offline jet, shown for the low (blue), medium (green), and high (red) pile-up conditions at the beginning of 2016 operations.

the 2016 run. Comparing the three $E_T$ resolutions was a good test of the chunky-donut algorithm. In the higher pile-up conditions, the Level-1 jet energies were, on average, increased relative to the energies of the offline jets. The difference, however, was very small, indicating that the chunky-donut algorithm was performing well.

The $H_T$ turn-on curves, for a variety of different thresholds, are provided in Figure 4.6. The offline $H_T$ quantity was calculated using offline jets with $E_T > 30$ GeV and $|\eta| < 2.4$ in order to replicate the Level-1 $H_T$ quantity. The $H_T$ turn-on curves all reach efficiencies of unity. For an equivalent energy threshold, the widths of these error functions are larger than those in the single jet turn-on curves. This is because the $H_T$ quantity can be composed of multiple low $E_T$ jets which, in the the Level-1 trigger, have a worse relative energy resolution than high $E_T$ jets. In addition, the $H_T$ turn-on curves are not centred on their corresponding threshold. The 50% efficiency points occur at offline $H_T$ values greater than the threshold values, meaning that the Level-1 $H_T$ was being underestimated. This shift was attributed to the energy scale of the sub-leading jets, because the leading jet energy scales agreed well in Figure 4.4. Ultimately, the shift in energy scale had only a small impact on the

Figure 4.6: Efficiency of reconstructing $H_T$ in the Level-1 trigger, above a variety of thresholds, as a function of the offline $H_T$ quantity.

performance. Furthermore, the study has a strong dependence on the type of offline jet used as a reference. Jets composed of calorimeter energy towers, rather than PF particles, were used when the study was repeated on 2017 data and resulted in the Level-1 $H_T$ being overestimated [56].

## 4.4 Jet Energy Corrections

The relative simplicity of the Level-1 jet algorithm gives rise to jets that are susceptible to losses. Due to the magnetic field, charged hadrons are bent away from neutral particles arising from the same outgoing parton. If they are bent outside the corresponding $9 \times 9$ trigger tower array, the reconstructed jet will be missing some energy. To account for this, corrections are applied to the $E_T$ of the Level-1 jets. It should be noted that the trigger primitives are calibrated upstream in Layer-1. Therefore, to leading order, the JECs should not account for incorrect trigger tower energies.

The degree to which a charged hadron is separated in the calorimeters depends on its $p_T$ and the distance it propagates in the transverse plane (i.e. its $|\eta|$ coordinate).

As a consequence, the JECs are applied as a function of the Level-1 jet energy, $E_\mathrm{T}^\mathrm{L1}$, and pseudorapidity, $|\eta^\mathrm{L1}|$. In this section, the derivation and implementation of the JECs are described followed by an example of the closure tests employed.

## 4.4.1  Derivation

In order to derive the JECs, the Level-1 jets must be compared against reference jets which best represent the 'true' jet momenta. The reference jets are formed from generator level particles. All the stable output particles (except neutrinos) are clustered using the anti-$k_T$ algorithm with distance parameter $R = 0.4$. The Level-1 jet and reference jet collections are produced using QCD multi-jet MC samples. The simulation of the detector and trigger electronics are configured to the state of CMS in the upcoming collisions. The only difference is that JECs are not applied to the Level-1 jets. The QCD sample used has a uniform pile-up distribution. The calibration procedure only uses events with $n_\mathrm{PU}$ (the number of pile-up interactions) in a range representing the expected pile-up conditions.

The Level-1 and reference jets are matched spatially, to avoid any $E_\mathrm{T}$-dependant bias, using the following method:

- Loop through all the Level-1 jets, for a given event, in descending $E_\mathrm{T}^\mathrm{L1}$ order. Jets with saturated trigger towers are omitted from the procedure because they are given a maximum energy of 1023.5 GeV.

- Find the closest reference jet in $(\eta,\,\phi)$ space using $\Delta R$. In order to avoid the reference jet being arbitrarily soft, it must have $E_\mathrm{T}^\mathrm{ref} > 10$ GeV.

- If $\Delta R < 0.25$, the jets are paired. The reference jet is then removed from the collection to avoid multiple pairings.

The matching condition of $\Delta R < 0.25$ is relatively tight. This was chosen to ensure that the reference jets are not paired with an unrelated Level-1 jet, which would degrade the quality of the calibrations. Nevertheless, a matching efficiency of 95%

is achieved for Level-1 jets with $30 < E_{\mathrm{T}}^{\mathrm{L1}} < 40$ GeV and of at least 99% for jets with $E_{\mathrm{T}}^{\mathrm{L1}} > 75$ GeV.

The calibration metric used is the response, defined as $E_{\mathrm{T}}^{\mathrm{L1}}/E_{\mathrm{T}}^{\mathrm{ref}}$. As was discussed above, such a quantity varies as a function of $E_{\mathrm{T}}^{\mathrm{L1}}$ and $|\eta^{\mathrm{L1}}|$, due to the loss of charged hadrons in Level-1 jets. The JECs are obtained by a procedure that, for a given $E_{\mathrm{T}}^{\mathrm{L1}}$ and $|\eta^{\mathrm{L1}}|$, inverts the average response. First, the jet pairs are divided into 16 $|\eta^{\mathrm{L1}}|$ bins. For each $|\eta^{\mathrm{L1}}|$ bin, a smooth correction factor is found as a function of $E_{\mathrm{T}}^{\mathrm{L1}}$. This is done by further categorizing the jet pairs into $E_{\mathrm{T}}^{\mathrm{ref}}$ bins. Within each $E_{\mathrm{T}}^{\mathrm{ref}}$ bin, $\langle E_{\mathrm{T}}^{\mathrm{L1}} \rangle$ is determined using the arithmetic mean and $\langle E_{\mathrm{T}}^{\mathrm{L1}}/E_{\mathrm{T}}^{\mathrm{ref}} \rangle$ is determined by fitting a Gaussian to the response distribution. Using all the $E_{\mathrm{T}}^{\mathrm{ref}}$ bins, a graph of $\langle E_{\mathrm{T}}^{\mathrm{L1}}/E_{\mathrm{T}}^{\mathrm{ref}} \rangle^{-1}$ vs. $\langle E_{\mathrm{T}}^{\mathrm{L1}} \rangle$ is then constructed. The fit to this graph then gives the correction factor as a function of $E_{\mathrm{T}}^{\mathrm{L1}}$.

For $10 < E_{\mathrm{T}}^{\mathrm{ref}} < 320$ GeV, the $E_{\mathrm{T}}^{\mathrm{ref}}$ bins are 4 GeV wide, providing sufficient statistics for the calculations of $\langle E_{\mathrm{T}}^{\mathrm{L1}} \rangle$ and $\langle E_{\mathrm{T}}^{\mathrm{L1}}/E_{\mathrm{T}}^{\mathrm{ref}} \rangle$ whilst accurately capturing the functional shape of the correction factor. For $E_{\mathrm{T}}^{\mathrm{ref}} > 320$ GeV, the bin size is increased to 20 GeV to handle the lower statistics. This coarser binning is adequate as the correction factor has much less variation at high energies. The number of $|\eta^{\mathrm{L1}}|$ bins is dictated by the implementation of the JECs in hardware. How these bins are selected is described at the end of this subsection.

Figure 4.7 shows an example graph, in the first $|\eta^{\mathrm{L1}}|$ bin, of $\langle E_{\mathrm{T}}^{\mathrm{L1}}/E_{\mathrm{T}}^{\mathrm{ref}} \rangle^{-1}$ vs. $\langle E_{\mathrm{T}}^{\mathrm{L1}} \rangle$. At large $E_{\mathrm{T}}^{\mathrm{L1}}$ values, only small losses are incurred and the correction factors are just above unity. As $E_{\mathrm{T}}^{\mathrm{L1}}$ decreases, the charged hadrons are, on average, bent further from the neutral particles. This causes larger losses and, therefore, an increase in the correction factor, especially for $E_{\mathrm{T}}^{\mathrm{L1}} < 150$ GeV. The correction factors typically reach a maximum at $E_{\mathrm{T}}^{\mathrm{L1}} \approx 20$ GeV. As the jet energy is reduced further, the correction factors rapidly decrease. This was discovered to be an artefact of the jet seed threshold used in the Level-1 jet reconstruction. A reference jet can only be (correctly) matched if the corresponding Level-1 jet passes the jet seed threshold condition. At low $E_{\mathrm{T}}^{\mathrm{ref}}$, this creates a bias (relative to the full set of Level-1 jets

Figure 4.7: Graph of the jet energy correction factor, as a function of $E_T^{L1}$, in the first $|\eta^{L1}|$ bin $(0 < |\eta| < 0.435)$. The fit represents the corrections that will be implemented.

when no jet seed threshold is applied) towards Level-1 jets with higher energies, reducing the correction factor.

The fit to the graphs is such that the corrected Level-1 jet energy is:

$$E_T^{L1,corr} = E_T^{L1} \times \left( p_0 + p_1 \cdot \text{erf}\left( p_2 \log\left(E_T^{L1} - p_3\right) + p_4 \cdot \exp\left(p_5 (\log\left(E_T^{L1} - p_6\right))^2\right)\right)\right) \quad (4.1)$$

where $p_i$ parameterize the functional form. This function was selected after trialling a number of different functions used in the offline JECs. As can be seen in Figure 4.7, the fit range typically has a minimum corresponding to the greatest correction factor and a maximum at a point where the correction factor has plateaued (or the statistics have become too low). Outside the fit range, the fit function values at the extrema of the fit range are used for the JECs. A future JEC procedure could improve upon this by attempting to capture the behaviour of the response at low $E_T^{L1}$. Such low energy jets are not, however, particularly relevant to the Level-1 trigger event selection criteria.

Due to their trigger tower structure, there are 40 $|\eta^{\mathrm{L1}}|$ values available to the Level-1 jets (see Table 4.1 and note that trigger towers 28 & 29 overlap). For the JECs, these are arranged into 16 $|\eta^{\mathrm{L1}}|$ bins, the maximum number that can be implemented in hardware due to memory limitations. This is achieved by performing the calibrations on all 40 Level-1 jet $|\eta^{\mathrm{L1}}|$ categories and grouping them together depending on their correction curves. In 2016, this was a new addition to the JEC procedure because prior to this, in the legacy calorimeter trigger, the Level-1 jets were significantly coarser objects and only had 11 different $|\eta^{\mathrm{L1}}|$ values available.

Figure 4.8 shows the correction curves for all the Level-1 jet $|\eta^{\mathrm{L1}}|$ categories. Each curve has a similar shape, following the above description of charged hadron loss. The overall correction values, however, have a large dependence on $|\eta^{\mathrm{L1}}|$. This is because $|\eta^{\mathrm{L1}}|$ maps onto the transverse distance between the interaction point and the calorimeters. The greater this distance, the larger the charged hadron angular separation and, thus, the greater the energy loss. Consequently, in the endcap calorimeters, the correction factors significantly decrease as $|\eta^{\mathrm{L1}}|$ increases. In order to capture this behaviour, the $|\eta^{\mathrm{L1}}|$ bins are chosen to be narrower in the endcaps at the expense of being wider in the barrel. The HF correction curves represent small energy losses because the HF is at such high pseudorapidities. In fact, the correction factors drop below unity. This is due to the calibration of the HF trigger primitives and the degradation of the chunky-donut algorithm at high $|\eta^{\mathrm{L1}}|$. The HF fit ranges are not as long because there are less high $E_{\mathrm{T}}$ statistics at larger pseudorapidities. Furthermore, low statistics meant that correction curves could not be derived for Level-1 jets centred on trigger towers with an $|\eta|$ index between 38 and 41. As a consequence, having the necessary statistics to perform the JECs was also considered when determining the $|\eta^{\mathrm{L1}}|$ bins. The arrangement of the 16 $|\eta^{\mathrm{L1}}|$ bins can be seen in Table 4.2.

### 4.4.2 Lookup Tables

Following their derivation, the functional forms of the JECs are implemented in hardware. This is done using a set of lookup tables (LUTs).

Figure 4.8: Corrections curves for all the Level-1 jet $|\eta^{\mathrm{L1}}|$ categories, labelled by their central trigger tower $|\eta|$ index. Jets with indices 38-41 are omitted due to low statistics when deriving the corrections.

Table 4.2: The arrangement of the 16 $|\eta^{\mathrm{L1}}|$ bins in terms of the Level-1 jet central trigger tower $|\eta|$ index ($i_{\mathrm{CTT}}$) and the $|\eta|$ range.

| $|\eta^{\mathrm{L1}}|$ bin | $i_{\mathrm{CTT}}$ | $|\eta|$ range |
|---|---|---|
| 1 | 1 - 5 | 0.000 - 0.435 |
| 2 | 6 - 9 | 0.435 - 0.783 |
| 3 | 10 - 13 | 0.783 - 1.131 |
| 4 | 14 - 15 | 1.131 - 1.305 |
| 5 | 16 - 17 | 1.305 - 1.479 |
| 6 | 18 - 19 | 1.479 - 1.653 |
| 7 | 20 - 21 | 1.653 - 1.830 |
| 8 | 22 | 1.830 - 1.930 |
| 9 | 23 | 1.930 - 2.043 |
| 10 | 24 | 2.043 - 2.172 |
| 11 | 25 | 2.172 - 2.322 |
| 12 | 26 | 2.322 - 2.500 |
| 13 | 27 - 28 | 2.500 - 2.964 |
| 14 | 30 - 32 | 2.964 - 3.489 |
| 15 | 33 - 36 | 3.489 - 4.191 |
| 16 | 37 - 41 | 4.191 - 5.191 |

Each correction curve is approximated by a set of linear relationships between the Level-1 jet energy before ($E_{\mathrm{T}}^{\mathrm{L1,pre}}$) and after the corrections are applied ($E_{\mathrm{T}}^{\mathrm{L1,corr}}$):

$$E_{\mathrm{T}}^{\mathrm{L1,corr}} = \left( (m \times E_{\mathrm{T}}^{\mathrm{L1,pre}}) >> 9 \right) + c \tag{4.2}$$

where $m$ is an unsigned 10-bit quantity and $c$ is a signed 8-bit quantity. This is done in 16 different $E_{\mathrm{T}}^{\mathrm{L1,pre}}$ bins. These bins are tuned to give the smoothest representation of the correction curves and, therefore, are finer at low $E_{\mathrm{T}}^{\mathrm{L1,pre}}$ values.

The $E_{\mathrm{T}}^{\mathrm{L1,pre}}$ and $|\eta^{\mathrm{L1}}|$ binning conventions are each compressed into LUTs. Using these LUTs, every Level-1 jet receives an address pointing to the relevant correction information, from which Equation 4.2 can be applied. The correction information is encoded in another LUT as a set of 18-bit quantities, given by $(c << 10) + m$. There are 256 entries, corresponding to the 16 $|\eta^{\mathrm{L1}}|$ bins and the 16 $E_{\mathrm{T}}^{\mathrm{L1,pre}}$ bins.

### 4.4.3 Closure Test

It is imperative to check the JECs before they are put online in the Level-1 trigger. This is done by recreating the simulation used in the derivation procedure, but with the new JECs applied to the Level-1 jets. The Level-1 and reference jets are then matched spatially, using the same method explained above. Their transverse energies can then be compared, providing a test of both the correction functions and their implementation in the LUTs. A variety of tests are conducted:

- The JECs are re-derived. They should yield a correction factor close to unity across all $E_{\mathrm{T}}^{\mathrm{L1}}$ and $|\eta^{\mathrm{L1}}|$ space.

- Turn-on curves, like in Figure 4.4, are created. These ensure that the Level-1 jets can be triggered on efficiently and assess the overall jet $E_{\mathrm{T}}$ scale.

- Scatter plots of $E_{\mathrm{T}}^{\mathrm{L1}}$ vs. $E_{\mathrm{T}}^{\mathrm{ref}}$ are produced in each $|\eta^{\mathrm{L1}}|$ bin for a direct comparison of the transverse energies. An example, in the first $|\eta^{\mathrm{L1}}|$ bin, is provided in Figure 4.9.

Figure 4.9: Scatter plot of $E_{\mathrm{T}}^{\mathrm{L1}}$ vs. $E_{\mathrm{T}}^{\mathrm{ref}}$ in the first $|\eta^{\mathrm{L1}}|$ bin ($0 < |\eta| < 0.435$). The diagonal line represents $E_{\mathrm{T}}^{\mathrm{L1}} = E_{\mathrm{T}}^{\mathrm{ref}}$. Top: without JECs applied. Bottom: with the JECs applied.

## 4.5  Problems with the Level-1 Trigger

This section discusses two problems that occurred in the Level-1 trigger which are relevant to the main analysis described later in this thesis.

### 4.5.1  $H_{\mathrm{T}}$ Saturation Issue

Before the last month of pp collisions in 2016, the $H_{\mathrm{T}}$ algorithm was updated so that the maximum energy would be obtained if there were any saturated trigger towers in the event. It was, however, not implemented correctly, creating an overflow error which caused an inefficiency in accepting high $H_{\mathrm{T}}$ events. All the events lost by the $H_{\mathrm{T}}$ trigger are recovered with the inclusion of a single jet trigger, because the jet algorithm was correctly handling any trigger tower saturation. The $H_{\mathrm{T}}$ saturation issue was fixed for 2017 operations.

### 4.5.2  Prefiring

The prefire issue was a problem that occurred during 2016 and 2017 operations. During these years, a gradual timing shift in the ECAL was not correctly propagated to the trigger primitives. As a result, a sizeable fraction of high pseudorapidity ECAL trigger primitives were incorrectly associated with the previous bunch crossing. This caused a serious problem because when an event is accepted by the Level-1 trigger, the next two bunch crossings cannot be accepted due to trigger rules (required to prevent readout buffer overflows). Consequently, events could self veto (prefire) if they deposited a significant amount of energy in the $2.0 < |\eta| < 3.0$ region of the ECAL, as this could cause a fake $e/\gamma$ trigger to occur in the previous bunch crossing. The solution to this problem is discussed in Section 7.1.4.

<center>5</center>

# Introduction to the Analysis

In the remaining chapters of this thesis, I present my search for pairs of light Higgs bosons produced in supersymmetric decay cascades. The analysis was performed using proton-proton collision data recorded with the CMS detector at a centre-of-mass energy of 13 TeV. The full data sets from 2016 and 2017 are used, corresponding to integrated luminosities of 35.9 fb$^{-1}$ and 41.5 fb$^{-1}$, respectively.

In this chapter, the signal model is described and the reason it is of interest is explained. This is followed by information about the different data sets and MC samples used in the analysis.

## 5.1 Signal Model and Search Motivation

### 5.1.1 Introduction to the Signal Model

This analysis searches for light Higgs bosons in supersymmetric decay cascades following the production of $\widetilde{q}\widetilde{q}$, $\widetilde{q}\widetilde{g}$, and $\widetilde{g}\widetilde{g}$ states in pp collisions at a centre-of-mass energy of 13 TeV. The signal model exists within the framework of the Next-to Minimal Supersymmetric Standard Model (NMSSM), which was outlined in Section 2.2.2. An in-depth description of the signal model can be found in Ref. [57].

Figure 5.1: Feynman diagram of the signal model following squark pair production. The particle $\tilde{\chi}_2^0$ is the NLSP, the particle $\tilde{\chi}_1^0$ is the LSP, and the particle $H$ is the (light) CP-even neutral Higgs boson.

Figure 5.1 shows the signal model Feynman diagram in the case of $\tilde{q}\tilde{q}$ production. The decays of $\tilde{q}_i \rightarrow q_i + \tilde{\chi}_2^0$ and $\tilde{\chi}_2^0 \rightarrow H + \tilde{\chi}_1^0$ occur with branching ratios of unity. The Higgs boson in the decay cascade is not necessarily the SM-like Higgs boson. It can instead be the CP-even Higgs boson in the NMSSM with a lower mass than the SM Higgs boson. This unknown mass is denoted by the parameter $M_H$. This analysis targets the final state where both Higgs bosons decay into $b\bar{b}$ pairs. The light Higgs boson has the same (relative) couplings as the SM Higgs. The $H \rightarrow b\bar{b}$ branching ratio is enhanced, however, due to the lower Higgs boson mass.

In this signal model, the light-flavour squarks are degenerate in mass such that $m_{\tilde{q}} = \{m_{\tilde{u}}, m_{\tilde{d}}, m_{\tilde{s}}, m_{\tilde{c}}\}$. The mass of these squarks is denoted by $M_{\mathrm{SUSY}}$, the supersymmetric (SUSY) production mass scale. In the $\tilde{q}\tilde{g}$ and $\tilde{g}\tilde{g}$ production modes, the gluinos decay as $\tilde{g} \rightarrow \tilde{q}_i + q_i$. The squarks then decay as they do in Figure 5.1. There are two limit cases considered for the gluino mass relative to the squark masses:

1. The gluino mass is set to be 1% higher than that of the squarks. This nominal mass gap means that, for a given squark mass, the total production cross section is almost maximal whilst still allowing the $\tilde{g} \rightarrow \tilde{q}_i + q_i$ decay to occur. Because the mass gap is so small, only a small amount of momentum

is transferred to the quark in the gluino decay. Consequently, the kinematics of the final state particles are very similar for the $\widetilde{q}\widetilde{q}$, $\widetilde{q}\widetilde{g}$, and $\widetilde{g}\widetilde{g}$ production modes.

2. The gluino is considered too massive to be produced. Because the gluino is decoupled, the signal process can only occur through $\widetilde{q}\widetilde{q}$ production. This means that, for a given squark mass, the total production cross section is minimal.

In this thesis, the default case used is that where the gluino mass is nominally higher than that of the squarks. It will be stated explicitly if the other case is used. It should also be noted that the additional SUSY particles in the NMSSM are assumed to be too massive to be produced directly at a significant rate.

In addition to $M_{\mathrm{H}}$ and $M_{\mathrm{SUSY}}$, there are two other unknown masses in the signal model; those of the $\widetilde{\chi}_1^0$ and $\widetilde{\chi}_2^0$ neutralinos. These are the Lightest SUSY Particle (LSP) and the Next-to Lightest SUSY Particle (NLSP), respectively. The degrees of freedom due to the LSP and NLSP are parameterized by $\Delta_M \equiv M_{\widetilde{\chi}_2^0} - M_{\mathrm{H}} - M_{\widetilde{\chi}_1^0}$ and $R_M \equiv M_{\mathrm{H}} / M_{\widetilde{\chi}_2^0}$. Using these parameters, the signal model is configured so that it has similar kinematic characteristics to the original benchmark model.

## 5.1.2 Mass Configuration of Interest

Under initial inspection, it would appear that the key features required to identify the signal model are multiple jets and $E_{\mathrm{T}}^{\mathrm{miss}}$, where the $E_{\mathrm{T}}^{\mathrm{miss}}$ arises from the LSPs which leave the detector without interacting. This is not a unique final state, it is common to many SUSY/DM models and there have been extensive searches for such final states at CMS [58–60] and ATLAS [61].

The signal model becomes of interest when one considers the dynamics of the decay cascades in further detail. The important part is the $\widetilde{\chi}_2^0 \rightarrow \mathrm{H} + \widetilde{\chi}_1^0$ process which occurs within the decay cascade. As the parameter $R_M$ tends towards unity, the momentum transferred to the LSP tends towards zero (see Appendix A). Consequently, the $E_{\mathrm{T}}^{\mathrm{miss}}$ in the event due to the LSPs is highly suppressed. This provides

an all hadronic final state with low $E_\mathrm{T}^\mathrm{miss}$; one which the more conventional SUSY searches performed at CMS and ATLAS would not be sensitive to. Without the presence of $E_\mathrm{T}^\mathrm{miss}$, reconstructing the pair of light Higgs bosons becomes the key way to identify the signal model.

The other degree of freedom in the signal model, $\Delta_M$, controls the amount of phase space that the products of the $\widetilde{\chi}_2^0 \to \mathrm{H} + \widetilde{\chi}_1^0$ decay have (see Appendix A). The larger $\Delta_M$ is, the more phase space available. When $\Delta_M = 0$ GeV, the Higgs boson and the LSP are produced at rest in the rest frame of the NLSP. The $\Delta_M$ parameter is constrained to be small as $R_M$ tends towards unity.

Signal samples were produced for a variety of $M_\mathrm{H}$ and $M_\mathrm{SUSY}$ values. The other two mass parameters were fixed with values of $R_M = 0.99$ and $\Delta_M = 0.1$ GeV in order to maximally utilize the behaviour of interest. A study of the effect that reducing $R_M$ from unity has on this analysis is presented in Appendix B. Note that in the rest of this thesis, the values of $R_M$ and $\Delta_M$ are always held constant. More information on the production of the signal samples is available in Section 5.2.

In the signal model, the Higgs boson mass is an unknown parameter that satisfies $2M_\mathrm{b} \leq M_\mathrm{H} \leq M_{H_\mathrm{SM}}$. It should be stressed, however, that it is the region $M_\mathrm{H} > \frac{1}{2}M_{H_\mathrm{SM}}$ which is of primary interest. When this is not the case, the SM Higgs boson can decay into a pair of the light Higgs bosons; a process that could be discovered in other searches. Nevertheless, to preserve generality in this analysis, the search attempts to probe as low a Higgs boson mass as possible.

## 5.1.3 NMSSM Formulation

The signal model cannot exist within the framework of the MSSM. The mass configuration of interest has a light LSP which, in the MSSM, must be a bino (it cannot be a higgsino or wino as they would have electrically charged partners with comparable masses, that would have been discovered at LEP). Consequently, the squarks could decay directly into the LSP, as they couple to binos. This would produce large $E_\mathrm{T}^\mathrm{miss}$ final states that would be observed in other SUSY searches.

Table 5.1: Parameters (left and middle column) and particle masses (right column) of the NMSSM benchmark point in Ref. [57].

| Parameter | Value | Parameter | Value | Particle(s) | Mass |
|---|---|---|---|---|---|
| $\lambda$ | $6.5 \times 10^{-3}$ | $M_1$ | 90 GeV | $M_{H_1}$ | 83 GeV |
| $\kappa$ | $1.9 \times 10^{-5}$ | $M_2$ | 950 GeV | $M_{H_2}$ | 123.2 GeV |
| $\tan\beta$ | 20 | $M_3$ | 830 GeV | $M_{H_3, A_2, H^\pm}$ | $\sim 950$ GeV |
| $\mu_{\text{eff}}$ | 900 GeV | $A_t$ | $-1500$ GeV | $M_{A_1}$ | 12.9 GeV |
| $\xi_S$ | $-1.02 \times 10^9$ GeV$^3$ | $A_b$ | $-1000$ GeV | $M_{\text{squarks}}$ | $\sim 860$ GeV |
| $m_S'^2$ | $3.6 \times 10^3$ GeV$^2$ | $m_{\text{sleptons}}$ | 600 GeV | $M_{\text{stop1}}$ | 810 GeV |
| $A_\kappa$ | 0 GeV | $m_{\text{squarks}}(\widetilde{u},\widetilde{d},\widetilde{s},\widetilde{c})$ | 830 GeV | $M_{\text{stop2}}$ | 1060 GeV |
| $A_\lambda$ | 50 GeV | $m_{\text{squarks}}(\widetilde{t},\widetilde{b})$ | 900 GeV | $M_{\text{gluino}}$ | 893 GeV |
| | | | | $M_{\widetilde{\chi}_1^0}$ | 5.26 GeV |
| | | | | $M_{\widetilde{\chi}_2^0}$ | 89 GeV |

In the NMSSM, however, the light LSP can be the singlino. This enables the signal model decay cascades, as the squarks cannot decay directly into a singlino LSP. The squarks must instead decay into a (bino) NLSP, which then decays into the LSP and a CP-even neutral Higgs boson. Furthermore, it is the NMSSM framework that allows for a CP-even neutral Higgs boson with a mass lower than that of the SM Higgs boson.

Table 5.1 provides an example set of NMSSM parameters, taken from Ref. [57], which give the desired particle masses and decay branching fractions of the signal model. The parameters $\lambda$, $\kappa$, $\xi_S$, $m_S'^2$, $\kappa$, and $A_\kappa$ were presented Equation 2.1; $\mu_{\text{eff}} = \lambda\langle S\rangle$; $\tan\beta = \langle H_u^0\rangle\big/\langle H_d^0\rangle$; $M_1$, $M_2$, and $M_3$ denote the bino, wino, and gluino mass terms, respectively; $A_t$ and $A_b$ denote the Higgs-stop and Higgs-sbottom trilinear couplings. The particle masses and decay branching fractions were calculated using NMSSMTOOLS version 4.2.1 [62]. The decays are not shown in Table 5.1, but $\widetilde{q}_i \to q_i + \widetilde{\chi}_2^0$ and $\widetilde{\chi}_2^0 \to H + \widetilde{\chi}_1^0$ occur with branching fractions of 100%.

## 5.1.4 Key Properties of the Signal Model

The extent to which the mass parameterization described above suppresses $E_T^{\text{miss}}$ can be seen in Figure 5.2. Signal samples with $M_{\text{SUSY}} = 2000$ GeV have only a negligible fraction of events with $E_T^{\text{miss}} > 30$ GeV. It should be noted that variation in $M_H$ has almost no impact on these distributions.

Figure 5.2: Normalised distribution of the magnitude of the sum of the two LSP transverse momenta vectors (i.e. the $E_T^{miss}$ due to the LSPs). This is shown for a variety of $M_{SUSY}$ values as indicated by the legend. For all the distributions the value of $M_H$ is 70 GeV.

In suppressing the momentum transferred from the NLSP to the LSP, the mass parameterization instead causes nearly all the NLSP momentum to be transferred to the Higgs boson. Consequently, the Higgs bosons are highly boosted objects and this results in a small $\Delta R$ separation of the $b\bar{b}$ pairs that they decay into. This can be seen in Figure 5.3 for signal samples with different $M_{SUSY}$ values and $M_H$ = 70 GeV. In all the distributions, the majority of $b\bar{b}$ pairs have $\Delta R < 0.8$. This is an important threshold because when $\Delta R < 0.8$, the Higgs boson can be reconstructed as a single AK8 jet. Doing so avoids problems resolving the Higgs boson as two separate AK4 jets when $\Delta R < 0.4$, as was discussed in Section 3.4.3.

The fraction of events where both $b\bar{b}$ pairs have an angular separation of $\Delta R < 0.8$ is presented in Table 5.2. One can see that both $M_H$ and $M_{SUSY}$ have an impact on the $\Delta R$ separation. In fact, this is the only aspect of the final state kinematics on which $M_H$ has a noticeable effect. In accordance with $\Delta R_{b\bar{b}} \approx 2M_H/(p_T)_H$, the more massive Higgs bosons have fewer events where both $b\bar{b}$ pairs have a separation of $\Delta R < 0.8$. The larger $M_{SUSY}$ is, the greater the energy provided to the decay cascades. This results, on average, in Higgs bosons with higher $p_T$ and hence smaller $\Delta R$ separations. For most of the mass configurations, the dominant scenario is that

Figure 5.3: Normalised distribution of the $\Delta R(b\bar{b})$ separation resulting from $H \to b\bar{b}$ in the decay cascades. This is shown for a variety of $M_{\mathrm{SUSY}}$ values as indicated by the legend. For all the distributions the value of $M_{\mathrm{H}}$ is 70 GeV.

Table 5.2: Fraction of signal events where both $H \to b\bar{b}$ pairs have an angular separation of $\Delta R < 0.8$.

| $M_{\mathrm{H}}$ | $M_{\mathrm{SUSY}}$ | | | | | |
|---|---|---|---|---|---|---|
| | 800 GeV | 1200 GeV | 1600 GeV | 2000 GeV | 2400 GeV | 2800 GeV |
| 30 GeV | 0.94 | 0.98 | 0.99 | 0.99 | 0.99 | 1.00 |
| 50 GeV | 0.78 | 0.90 | 0.95 | 0.96 | 0.98 | 0.98 |
| 70 GeV | 0.57 | 0.78 | 0.88 | 0.92 | 0.95 | 0.96 |
| 90 GeV | 0.39 | 0.65 | 0.78 | 0.86 | 0.90 | 0.93 |
| 110 GeV | 0.27 | 0.54 | 0.69 | 0.79 | 0.84 | 0.88 |
| 125 GeV | 0.20 | 0.45 | 0.62 | 0.73 | 0.80 | 0.85 |

where both $b\bar{b}$ pairs have $\Delta R < 0.8$ and, thus, the best method for reconstructing the Higgs bosons are in AK8 jets. Although this is not the case for the high $M_{\mathrm{H}}$ and low $M_{\mathrm{SUSY}}$ combinations, it is compensated for by the production cross section, which increases rapidly as the value of $M_{\mathrm{SUSY}}$ is reduced.

Another important property of the signal model is the $H_{\mathrm{T}}$ distribution (defined henceforth as; $H_{\mathrm{T}} \equiv \sum_{\mathrm{AK4\ jets}} p_{\mathrm{T}}$, using AK4 jets with $p_{\mathrm{T}} > 40$ GeV and $|\eta| < 3.0$). This is shown in Figure 5.4 for signal samples with different $M_{\mathrm{SUSY}}$ values and $M_{\mathrm{H}} = 70$ GeV. The $H_{\mathrm{T}}$ distribution has very large values. This is because essentially all of the energy of the initial $\tilde{q}\tilde{q}$, $\tilde{q}\tilde{g}$, or $\tilde{g}\tilde{g}$ state is converted into jets. It also means that the $H_{\mathrm{T}}$ distribution has a significant dependence on $M_{\mathrm{SUSY}}$. The larger

Figure 5.4: Normalised $H_{\mathrm{T}}$ distribution. This is shown for a variety of $M_{\mathrm{SUSY}}$ values as indicated by the legend. For all the distributions the value of $M_{\mathrm{H}}$ is 70 GeV.

$M_{\mathrm{SUSY}}$ is, the greater the $H_{\mathrm{T}}$ in the events. This feature is exploited in the event selection (see Section 6.2). By binning in $H_{\mathrm{T}}$, one can gain sensitivity to signal models with high $M_{\mathrm{SUSY}}$ values, despite the rapidly falling production cross sections. Additionally, the high valued $H_{\mathrm{T}}$ distribution of the signal model enables efficient triggering using $H_{\mathrm{T}}$ (see Section 6.1).

## 5.2 Data Sets and Simulation

### 5.2.1 Data Sets

This analysis uses the pp collision data, at a centre-of-mass energy of 13 TeV, collected by the CMS detector in 2016 and 2017. The full 2016 data set, which is certified for analyses, corresponds to an integrated luminosity of $35.92 \pm 0.90$ fb$^{-1}$. The full certified 2017 data set corresponds to an integrated luminosity of $41.53 \pm 0.96$ fb$^{-1}$. The JetHT primary data sets, collected using jet and $H_{\mathrm{T}}$ triggers, are used for the main analysis. The SingleMuon primary data sets, collected using muons triggers, are used to study the trigger strategy of the main analysis.

### 5.2.2 Standard Model MC Samples

The 2016 and 2017 Standard Model MC samples are listed in Table 5.3 and Table 5.4, respectively. The production cross section and the effective integrated luminosity of each sample are also provided. The QCD samples are used to help validate the data driven background estimation of QCD multi-jet processes and the other samples are used to directly estimate the background yields of the processes they represent.

The Standard Model MC samples used in this analysis were produced centrally by the CMS collaboration. For a given year, the simulation of the CMS detector corresponds to its operational state. In addition, the pile-up distributions were set to represent those observed during the collisions.

Table 5.3: The 2016 Standard Model MC samples used in the analysis. For each sample the production cross section (to four significant figures) and the integrated luminosity (to one decimal place) are provided.

| Data set | $\sigma$ (pb) | $\int \mathcal{L} dt$ (fb$^{-1}$) |
|---|---|---|
| QCD_HT1000to1500_madgraph-pythia8 | 1206 | 12.5 |
| QCD_HT1500to2000_madgraph-pythia8 | 120.4 | 98.2 |
| QCD_HT2000toInf_madgraph-pythia8 | 25.25 | 239.2 |
| ZJetsToQQ_HT600toInf_madgraph-pythia8 | 52.79 | 18.9 |
| WJetsToQQ_HT600ToInf_madgraph-pythia8 | 95.14 | 10.8 |
| TT_powheg-pythia8 | 831.8 | 186.3 |

Table 5.4: The 2017 Standard Model MC samples used in the analysis. For each sample the production cross section (to four significant figures) and the integrated luminosity (to one decimal place) are provided.

| Data set | $\sigma$ (pb) | $\int \mathcal{L}dt$ (fb$^{-1}$) |
|---|---|---|
| QCD_HT1000to1500_madgraph-pythia8 | 1005 | 16.2 |
| QCD_HT1500to2000_madgraph-pythia8 | 101.8 | 109.1 |
| QCD_HT2000toInf_madgraph-pythia8 | 20.54 | 261.3 |
| ZJetsToQQ_HT-800toInf_madgraph-pythia8 | 18.69 | 418.3 |
| WJetsToQQ_HT800toInf_madgraph-pythia8 | 34.00 | 237.7 |
| TTtoHadronic_powheg-pythia8 | 378.0 | 446.7 |
| TTtoSemiLeptonic_powheg-pythia8 | 365.3 | 407.1 |
| TTto2L2Nu_powheg-pythia8 | 88.29 | 835.9 |

The QCD, Z→q$\overline{\text{q}}$, and W→q$\overline{\text{q}}$ events were generated, in association with up to four additional partons, at leading order using MadGraph5_aMC@NLO version 2.3.3 and 2.4.2 [63] for 2016 and 2017, respectively. The $H_{\text{T}}$ cuts applied to these MC samples were conducted using the partons at generator level. The t$\overline{\text{t}}$ events were generated at next-to leading order precision using powheg v2 [64]. For the 2016 MC samples, the parton distribution function (PDF) used was NNPDF30 [65] and for the 2017 MC samples, the PDF used was NNPDF31 [66]. The showering and hadronization of the partons was conducted using the pythia version 8.2 program [67], with tune CUETP8M1 [68] (CUETP8M2T4 [69] for t$\overline{\text{t}}$) and tune CP5 [70] for the 2016 and 2017 MC samples, respectively. The cross sections of the MC samples were attained by multiplying the cross sections calculated by the event generators with the jet matching rejection factors. In the case of the t$\overline{\text{t}}$ samples, the cross sections were rescaled to those calculated at next-to-next-to leading order precision using the Top++ v2.0 program [71].

It should be noted that the 2016 Z+jets MC sample was discovered to have not had jet matching correctly applied. As the sample was not regenerated in time, the 2017 Z+jets MC sample is being used in its place. This is a satisfactory approximation for two reasons. First, the analysis has only a very small dependence on the Z+jets process, as will be shown in Section 7.1.5. Second, there is very little difference, between 2016 and 2017 MC, in the distributions of the variables relevant to this analysis. This was tested by comparing the 2016 and 2017 W+jets MC samples.

Table 5.5: The Higgs boson masses used for the signal MC samples. For each value of $M_H$, the corresponding $H \to b\bar{b}$ branching ratio (to three decimal points) is provided.

| $M_H$ | BR ($H \to b\bar{b}$) |
|---|---|
| 30 GeV | 0.868 |
| 35 GeV | 0.867 |
| 40 GeV | 0.865 |
| 50 GeV | 0.858 |
| 70 GeV | 0.840 |
| 90 GeV | 0.816 |
| 110 GeV | 0.749 |
| 125 GeV | 0.581 |

### 5.2.3 Signal Model MC Samples

The signal model MC samples were produced, privately, for different combinations of $M_H$ and $M_{SUSY}$ values. As was discussed in Section 5.1.2, the two other degrees of freedom in the model were held constant with values of $R_M = 0.99$ and $\Delta_M = 0.1$ GeV. The different values of $M_H$ used, along with the corresponding branching ratio of $H \to b\bar{b}$, are shown in Table 5.5. The branching ratios were calculated using the HDECAY package [72]. As the value of $M_H$ increases, the $H \to b\bar{b}$ branching ratio decreases. For $30 < M_H < 90$ GeV, the rate of change in the branching ratio is small. For $M_H > 90$ GeV, the rate of change increases as the $WW^*$ and $ZZ^*$ decay channels start to become more accessible. Each value of $M_H$ is combined with the set of $M_{SUSY}$ values shown in Table 5.6. For each $M_{SUSY}$ value, there are two production cross sections provided. The first cross section is for the case where the gluino mass is decoupled and, thus, the production modes involving gluinos are kinematically inaccessible. The second cross section is for the case where the gluino mass is 1% higher than that of the squarks. The cross sections were calculated using PROSPINO [73] at next-to leading order. They quickly decrease as the value of $M_{SUSY}$ increases.

The simulations of the initial state squarks and gluinos, in association with up to one additional parton, were generated at leading order using MADGRAPH5_aMC@NLO version 2.3.3. The PDF used was NNPDF23 [74]. The PYTHIA version 8.2 program, with tune CUETP8M1, was used to describe the decay cascades of the squarks and

Table 5.6: The different SUSY production mass scales used for the signal MC samples. For each value of $M_{\text{SUSY}}$, the corresponding production cross sections (to three significant figures) are provided for the case where the gluino mass is decoupled and the case where it is 1% higher than that of the squarks. Note that the relative error on all the cross section calculations is less than 0.1%.

| | $\sigma$ (pb) | |
| $M_{\text{SUSY}}$ | pp $\to \widetilde{q}\widetilde{q}$ only | pp $\to \widetilde{q}\widetilde{q},\widetilde{q}\widetilde{g},\widetilde{g}\widetilde{g}$ |
|---|---|---|
| 800 GeV | $2.05 \times 10^0$ | $6.47 \times 10^0$ |
| 1200 GeV | $2.04 \times 10^{-1}$ | $4.95 \times 10^{-1}$ |
| 1600 GeV | $2.97 \times 10^{-2}$ | $6.04 \times 10^{-2}$ |
| 2000 GeV | $5.05 \times 10^{-3}$ | $9.11 \times 10^{-3}$ |
| 2200 GeV | $2.13 \times 10^{-3}$ | $3.68 \times 10^{-3}$ |
| 2400 GeV | $9.08 \times 10^{-4}$ | $1.51 \times 10^{-3}$ |
| 2600 GeV | $3.84 \times 10^{-4}$ | $6.17 \times 10^{-4}$ |
| 2800 GeV | $1.85 \times 10^{-4}$ | $2.75 \times 10^{-4}$ |

gluinos, along with the parton showering and hadronization. The reconstruction of these events in CMS was performed using the full detector simulation (see Section 3.2.6). It was done twice, to replicate 2016 and 2017 data taking, respectively. Around 200,000 events were generated for each signal sample. In order to validate this private MC production, some signal samples were reproduced centrally by the CMS collaboration. They had parameters $M_{\text{H}} = 70$ GeV and $M_{\text{SUSY}} = 1200, 2000$, and 2600 GeV; chosen so that the private signal samples could be tested at low, medium, and high energy scales. The comparisons showed very good agreement.

## 5.2.4  Multiple Years of Data and Simulation

The data and simulation used in this analysis are handled separately for 2016 and 2017. They are only brought together, at the end, when extracting limits. For a given physics process, the distributions of the quantities of interest are very similar between the two years. Consequently, when a point just needs to be illustrated in this thesis, only the figure/table for one of the years is provided. When the information is critical to the analysis, both the 2016 and 2017 versions are provided.

# 6

# Analysis Strategy

This chapter presents three key aspects of the analysis. First, information about the trigger strategy is provided. This is followed by a description of the event selection. Finally, a data driven technique to estimate the QCD multi-jet background is presented.

## 6.1 Trigger Strategy

In this analysis, the minimum $H_T$ required in the kinematic event selection is 1500 GeV (see Section 6.2.1). Due to this large $H_T$ requirement, highly efficient triggering can be achieved for both 2016 and 2017.

### 6.1.1 2016 Triggers

The events considered from the 2016 data are those collected by the logical OR of the `HLT_PFHT900` and `HLT_AK8PFJet450` trigger paths. These triggers are for $H_T$ and single jet $p_T$, respectively. They have the lowest thresholds that were unprescaled throughout 2016 data taking. The inclusion of the single jet trigger is to mitigate an inefficiency in $H_T$ triggering that occurred at Level-1 towards the end of 2016 data taking, as was explained in Section 4.5.1.

Figure 6.1: The 2016 trigger efficiency as a function of offline $H_{\mathrm{T}}$. The common criterion for the numerator and denominator is the pre-selection requirement. In addition, events containing offline muons with $p_{\mathrm{T}} > 200$ GeV are vetoed.

The trigger efficiency is evaluated using the 2016 SingleMuon data set. This data has no dependency on jet triggers, so provides an unbiased sample of events. Figure 6.1 shows the trigger efficiency as a function of offline $H_{\mathrm{T}}$. The common criterion for the numerator and denominator is the pre-selection, which requires at least two AK8 jets with $p_{\mathrm{T}} > 170$ GeV. In addition, events containing an offline muon with $p_{\mathrm{T}} > 200$ GeV are vetoed. For offline $H_{\mathrm{T}} > 1500$ GeV (the $H_{\mathrm{T}}$ requirement in the kinematic event selection) the trigger is 100% efficient.

## 6.1.2   2017 Triggers

The events considered from the 2017 data are those collected by the logical OR of the `HLT_PFHT1050` and `HLT_AK8PFJet500` trigger paths. These are the same triggers that are used for 2016, but with different thresholds. The unprescaled thresholds are higher in 2017 because of the greater the instantaneous luminosity provided by the LHC.

The trigger efficiency is evaluated in the same way as described in Section 6.1.1. The results are shown in Figure 6.2. Due to the higher trigger thresholds, 100% efficiency is achieved at a higher value of offline $H_{\mathrm{T}}$ when compared to 2016. Despite this, for offline $H_{\mathrm{T}} > 1500$ GeV, the trigger is still 100% efficient.

Figure 6.2: The 2017 trigger efficiency as a function of offline $H_T$. The common criterion for the numerator and denominator is the pre-selection requirement. In addition, events containing offline muons with $p_T > 200$ GeV are vetoed.

## 6.2 Event Selection

The event selection used in this analysis focusses on the properties of two AK8 jets. These jets are the candidates for the boosted topologies of the two $H \rightarrow b\overline{b}$ decays in the signal model. They are used to discriminate against background processes whilst retaining a significant fraction of signal events. The key properties of the AK8 jets are the double-b-tag discriminator and the soft-drop mass; both of which were introduced in Section 3.4. Due to the existence of other high energy outgoing quarks in the signal model, requirements on additional jets and the $H_T$ in the event are used to further discriminate against the background.

### 6.2.1 Kinematic Event Selection

The kinematic event selection uses both AK4 jets and AK8 jets. The $p_T$ and the angular coordinates of these jets are determined using the default jets stored in the data/simulation. This means that, to mitigate the effects of pile-up, the 2016 AK4 jets, 2017 AK4 jets, and the 2016 AK8 jets have CHS applied, whereas the 2017 AK8 jets have PUPPI applied.

The event pre-selection requires two AK8 jets with $p_{\mathrm{T}} > 170$ GeV (AK8 jets in 2016 data/simulation only exist for this $p_{\mathrm{T}}$ range) and $|\eta| < 2.4$ (so that they are within the acceptance of the tracker). If there are more than two candidate AK8 jets, the two with the highest double-b-tag discriminator values are selected, as these jets are most likely to have come from the boosted $H \rightarrow b\overline{b}$ decays. The two selected AK8 jets are then randomly allocated as 'fatJetA' and 'fatJetB'.

The following requirements are applied after pre-selection. They define the kinematic event selection:

- Both of the selected AK8 jets must have $p_{\mathrm{T}} > 300$ GeV.

- There must be at least one AK4 jet with $p_{\mathrm{T}} > 300$ GeV and $|\eta| < 3.0$. It must also satisfy $\Delta R > 1.4$ with respect to both of the selected AK8 jets, to avoid being composed of the same PF particles.

- $H_{\mathrm{T}}$ is binned in the following categories: 1500-2500, 2500-3500, and 3500+ GeV. The $H_{\mathrm{T}}$ is calculated using all AK4 jets with $p_{\mathrm{T}} > 40$ GeV and $|\eta| < 3.0$.

The $p_{\mathrm{T}}$ cuts are placed on the two selected AK8 jets because, given they are the boosted $H \rightarrow b\overline{b}$ candidates, they should be high $p_{\mathrm{T}}$ objects. The requirement of an additional high $p_{\mathrm{T}}$ AK4 jet is because such a jet often arises from one of the two high energy quarks produced by $\widetilde{q}_i \rightarrow q_i + \widetilde{\chi}_2^0$ in the signal model decay cascades. This requirement is very good at suppressing dijet background events. The AK4 jet $p_{\mathrm{T}}$ threshold is the same as that of the two AK8 jets because the $p_{\mathrm{T}}$ distributions of the quarks and the Higgs bosons are very similar (see Appendix C). This is because, in the rest frame of the squark, the energy is almost equally divided between the quark and the NLSP in the $\widetilde{q}_i \rightarrow q_i + \widetilde{\chi}_2^0$ decay.

The jet $p_{\mathrm{T}}$ distributions, after pre-selection is applied, for signal samples with different $M_{\mathrm{SUSY}}$ values, are provided in Figure 6.3. The $t\overline{t}$ distributions are also provided as a reference. In both the AK4 and AK8 categories, there are many signal model jets with transverse momenta that far exceed the 300 GeV threshold, especially as

the $M_{\text{SUSY}}$ parameter increases. Additional sensitivity to the events with high $p_{\text{T}}$ jets is achieved through the $H_{\text{T}}$ binning. The $H_{\text{T}}$ is a good kinematic quantity by which to categorize events because, for the signal model, it captures the energy scale of the whole event. As was discussed in Section 5.1.4, the $H_{\text{T}}$ distribution of the signal model has a large dependence on the $M_{\text{SUSY}}$ parameter. This can also be seen in Figure 6.3. As $M_{\text{SUSY}}$ increases, a greater fraction of the signal events populate the higher $H_{\text{T}}$ bins. This allows sensitivity to high $M_{\text{SUSY}}$ values to be retained, despite the rapid decrease in production cross section, as very few background events enter the highest $H_{\text{T}}$ bin. A minimum $H_{\text{T}}$ requirement of 1500 GeV allows for a trigger efficiency of 100%. Although reducing this $H_{\text{T}}$ threshold would increase the yields of the lowest $M_{\text{SUSY}}$ signal models, this is not necessary because, owing to their much larger production cross sections, these samples already have a significant number of events passing the kinematic cuts.

For the signal samples in Figure 6.3, the $p_{\text{T}}$ distribution of the leading separated AK4 jet does not have a trivial shape, due to the angular requirements imposed on the jet. As expected, the distributions have broad peaks arising from the reconstruction of one of the two high energy quarks produced by $\widetilde{q}_i \rightarrow q_i + \widetilde{\chi}_2^0$ in the decay cascades. However, as the $p_{\text{T}}$ tends from around 200 GeV towards zero, the distribution rises again. In these events, the AK4 jets corresponding to the two high energy quarks are not selected because they are not sufficiently separated from the two selected AK8 jets (they can also be outside the $|\eta| < 3.0$ acceptance). It is very rare that an overlap occurs between a Higgs boson and quark belonging to the same decay cascade arm. This is because, in conserving $p_{\text{T}}$, they are most likely to have $\phi$ coordinates separated by 180 degrees. Instead, the overlaps typically occur between the Higgs boson from one decay cascade arm and the quark from the other.

Table 6.3 shows how the cuts (including those discussed in Sections 6.2.2 and 6.2.3, thus, the table is presented in Section 6.2.4) reduce the event yields for the signal model with different values of $M_{\text{H}}$ whilst $M_{\text{SUSY}}$ is held constant. The value of $M_{\text{H}}$ has very little effect on the fraction of events passing the kinematic cuts. The cut

Figure 6.3: Normalised distributions of the quantities used in the kinematic event selection, after pre-selection has been applied, for 2017 signal samples with different $M_{\mathrm{SUSY}}$ values and $t\bar{t}$. Top left: The $p_{\mathrm{T}}$ distribution of one of the selected AK8 jets. Top right: The $p_{\mathrm{T}}$ distribution of the leading separated AK4 jet. Bottom: The $H_{\mathrm{T}}$ distribution. All signal samples have $M_{\mathrm{H}} = 70$ GeV. Red lines: $M_{\mathrm{SUSY}} = 1200$ GeV. Blue lines: $M_{\mathrm{SUSY}} = 2000$ GeV. Green lines: $M_{\mathrm{SUSY}} = 2600$ GeV. Orange lines: $t\bar{t}$.

flow table for the other case, where $M_{\mathrm{H}}$ is held constant and $M_{\mathrm{SUSY}}$ is varied, can be found in Table 6.4. As $M_{\mathrm{SUSY}}$ increases, a greater fraction of events pass the $p_{\mathrm{T}}$ cuts on the jets and, as was mentioned above, there is an increase in the populations of the higher $H_{\mathrm{T}}$ bins. For $M_{\mathrm{SUSY}} = 2400$ GeV, over 80% of events pass the kinematic event selection. Despite the kinematic cuts being more aggressive on the signal samples with lower $M_{\mathrm{SUSY}}$ values, the losses incurred are more than compensated for by the much larger production cross sections.

Table 6.2 shows how the cuts (including those discussed in Sections 6.2.2 and 6.2.3, thus, the table is presented in Section 6.2.4) reduce the event yields for the different background processes considered. After the kinematic event selection, the dominant background is from QCD multi-jet processes, with $7.6 \times 10^5$ events passing the cuts. This is due to the very large QCD cross sections and because the selection criteria are based only on jet objects. The full extent of the background suppression can only be seen for the $t\bar{t}$ process, as it is the only background MC sample without a $H_{\mathrm{T}}$ cut applied during generation. Every step in the event selection causes a significant reduction in the yield. The fraction of events passing the kinematic event selection in the three $H_{\mathrm{T}}$ bins are $2.9 \times 10^{-4}$, $2.6 \times 10^{-5}$, and $2.4 \times 10^{-6}$, respectively. For the Z+jets and W+jets samples, the fraction of events passing each step of the kinematic cuts is very similar. This is to be expected, as the only difference between the samples is the type of vector boson decaying into the quark pairs. After the kinematic cuts are applied, the $t\bar{t}$, Z+jets, and W+jets processes have total event yields of $1.1 \times 10^4$, $8.2 \times 10^3$, and $1.5 \times 10^4$, respectively.

## 6.2.2   Double-b-tag Event Selection

After pre-selection, there are two AK8 jets that have been identified as the boosted $\mathrm{H} \rightarrow \mathrm{b\bar{b}}$ candidates. They are chosen because they have the highest double-b-tag discriminator values. The discriminator, which varies between -1 and +1, represents how likely it is that an AK8 jet has originated from two b-quarks.

Figure 6.4: Scatter plot of signal events in the double-b-tag plane. The signal sample is from 2016 MC with the parameters $M_H = 70$ GeV and $M_{SUSY} = 2000$ GeV. Full kinematic cuts are applied with $H_T \in 3500+$ GeV. In addition, both AK8 jets are required to have a soft-drop mass greater than 15 GeV. The red triangle represents the tag double-b-tag region.

As was stated above, the two AK8 jets are randomly allocated as 'fatJetA' and 'fatJetB'. They create a 2D plane in double-b-tag space in which events exist. Figure 6.4 shows how this 2D plane is occupied by signal events with parameters $M_H = 70$ GeV and $M_{SUSY} = 2000$ GeV. As the signal process contains two $H \rightarrow b\bar{b}$ decays, the distribution primarily occupies the corner where both AK8 jets have double-b-tag discriminators of around unity. To capture this shape, a 'tag' double-b-tag region is defined as the region within the red triangle in Figure 6.4. This is the space in the 2D plane that satisfies Equation 6.1, where $A$ and $B$ represent the double-b-tag discriminator values of fatJetA and fatJetB, respectively.

$$B - 1.0 > -1 \cdot (A - 0.3) \tag{6.1}$$

The tag double-b-tag region is constructed so that the lowest double-b-tag discriminator value it contains is 0.3. This is because there is no scale factor (weights applied to MC events so that they better describe a certain aspect of the data) information provided for the double-b-tag discriminator below this value.

Table 6.1: Fraction of signal events in the tag double-b-tag region (after pre-selection is applied). This was calculated using the 2016 signal MC.

| $M_\mathrm{H}$ | $M_\mathrm{SUSY}$ | | | | | |
| | 800 GeV | 1200 GeV | 1600 GeV | 2000 GeV | 2400 GeV | 2800 GeV |
|---|---|---|---|---|---|---|
| 30 GeV | 0.56 | 0.61 | 0.61 | 0.60 | 0.57 | 0.55 |
| 50 GeV | 0.51 | 0.56 | 0.56 | 0.54 | 0.52 | 0.50 |
| 70 GeV | 0.46 | 0.51 | 0.51 | 0.50 | 0.48 | 0.47 |
| 90 GeV | 0.42 | 0.47 | 0.47 | 0.46 | 0.44 | 0.43 |
| 110 GeV | 0.38 | 0.44 | 0.44 | 0.43 | 0.42 | 0.40 |
| 125 GeV | 0.36 | 0.41 | 0.42 | 0.41 | 0.40 | 0.38 |

The fraction of signal events that are within the tag double-b-tag region varies with the $M_\mathrm{H}$ and $M_\mathrm{SUSY}$ parameters. This can be seen in Table 6.1, which shows this fraction for different signal models after pre-selection is applied. The variation occurs because different combinations of $M_\mathrm{H}$ and $M_\mathrm{SUSY}$ give rise to different kinematics of the $b\bar{b}$ pairs. The key kinematic properties of a $b\bar{b}$ pair, with regards to double-b-tagging, are the momenta and angular separation of the b-quarks. The average $p_\mathrm{T}$ of the b-quarks increases as $M_\mathrm{SUSY}$ increases, because the $M_\mathrm{SUSY}$ parameter sets the energy scale of the events. The angular separation of the b-quarks follows $\Delta R_{b\bar{b}} \approx 2M_\mathrm{H}/(p_\mathrm{T})_H$. Consequently, the average angular separation decreases as $M_\mathrm{H}$ gets smaller and $M_\mathrm{SUSY}$ gets larger (the average $p_\mathrm{T}$ of the Higgs bosons increases as $M_\mathrm{SUSY}$ increases). From Table 6.1 it can be seen that the $M_\mathrm{H}$ parameter has the largest influence on the double-b-tagging performance and that signal models with lower $M_\mathrm{H}$ values are favoured. Information of how the tag double-b-tag requirement affects the signal yields after the full kinematic cuts can be found in Table 6.3 and Table 6.4 (these tables are presented in Section 6.2.4).

Although the tag double-b-tag region contains a large fraction of signal events, it only covers 6.125% of the total double-b-tag plane, which results in it being a powerful discriminator against background processes. This is especially so for QCD, the primary background in this analysis. Figure 6.5 shows how the double-b-tag plane is occupied by simulated QCD events. The distribution covers all of the double-b-tag plane. The highest density of events are in the opposite corner to the tag double-b-tag region. This is because the majority of the AK8 jets coming from

Figure 6.5: Scatter plot of QCD MC events in the double-b-tag plane. The 2016 MC samples are used. Full kinematic cuts are applied with $H_T \in$ 1500-2500 GeV. In addition, both AK8 jets are required to have soft-drop masses greater than 15 GeV.

QCD do not resemble AK8 jets arising from two b-quarks. Table 6.2 shows how the cuts in the analysis reduce the event yields for the different background processes. Only around 3% of the QCD events meet the tag double-b-tag requirement after the kinematic event selection is applied. The W+jets process also has a 3% selection efficiency, however, it is 6% for the Z+jets process. The increase in efficiency is because the Z boson can decay into true $b\bar{b}$ pairs. The highest tag double-b-tag selection efficiency, for a background process, comes from $t\bar{t}$. This is because, with two W bosons and two b-quarks in each event, there are multiple candidates that could fake high double-b-tag discriminator values. In the $H_T \in$ 1500-2500 GeV bin, the $t\bar{t}$ double-b-tag selection efficiency is 13%. This reduces, however, to 9% in the $H_T \in$ 3500+ GeV bin.

For a given physics process, there is a difference in the double-b-tag discriminator distribution between 2016 and 2017 (and between the data/MC scale factors). This occurs because the CMS pixel detector was upgraded for 2017 operations and because the double-b-tag BDT was re-trained. The variation between the two years is not large. There is no significant difference in the results presented above de-

pending on whether they were obtained using 2016 MC or 2017 MC. Despite only being small, the change in the double-b-tag discriminator distribution is the largest difference, in a quantity of interest, between samples from the two years. A cut flow table comparing equivalent 2016 and 2017 signal samples is provided in Table 6.5.

### 6.2.3 Soft-drop Mass Event Selection

Another important property of the two selected AK8 jets is their invariant mass. This is because these jets are the boosted H $\rightarrow$ b$\overline{\text{b}}$ candidates and, for signal events, they should ideally have masses that reflect the $M_{\text{H}}$ parameter. In this analysis, the mass of the AK8 jets is evaluated following the soft-drop grooming algorithm with PUPPI pile-up mitigation. It is referred to as the soft-drop mass.

The soft-drop masses of the two selected AK8 jets create a 2D plane in which events exist. Figure 6.6 shows the 30 mass regions within this space that are used in this analysis. The equations of the lines that form the mass regions can be found in Appendix D. The regions labelled $S_n$ are the signal mass regions. They are designed to contain the events from signal models with different $M_{\text{H}}$ values. The $S_n$ mass regions are centred on the diagonal line $A = B$ (where $A$ and $B$ represent the soft-drop masses of fatJetA and fatJetB, respectively) because the two selected AK8 jets should have equal but unknown masses corresponding to the $M_{\text{H}}$ parameter. The $S_n$ mass regions broaden the further they are from the origin, because as the soft-drop mass increases, so too does its absolute resolution.

For each $S_n$ mass region there are two corresponding mass regions, $U_n$ and $D_n$. These mass regions are used as sidebands for the data driven QCD estimation (see Section 6.3). The regions labelled $U_n$ are the 'up-sideband' mass regions and those labelled $D_n$ are the 'down-sideband' mass regions. The regions $U_n$ and $D_n$ are reflections of each other in the diagonal line $A = B$. This means that the two mass regions are essentially equivalent because the event distribution in the 2D soft-drop mass plane is also symmetric under the transformation $A \leftrightarrow B$ (this is due to the random allocation of fatJetA and fatJetB). For $n > 1$, the area within the

Figure 6.6: The different mass regions used in the 2D soft-drop mass plane. The regions labelled $S_n$ are the signal mass regions, those labelled $U_n$ are the 'up' side-band mass regions, and those labelled $D_n$ are the 'down' sideband mass regions. The equations of the lines that form the mass regions can be found in Appendix D.

sideband regions $U_n + D_n$ is the same as that in the corresponding signal region, $S_n$. The sideband regions $U_1$ and $D_1$ take on a unique shape. The lowest mass node is removed to create a triangular shape rather than a quadrilateral one. This prevents the $U_1$ and $D_1$ mass regions from containing jets with very low soft-drop masses; the reason for which will be explained at the end of this subsection.

The event distributions in the 2D soft-drop mass plane, for a set of signal models with different $M_{\mathrm{H}}$ values, and with the 30 mass regions overlaid, are shown in Figure 6.7. For each of the scatter plots, there are four distinct distributions that smear into each other. The most prominent distribution is the circular distribution centred on the diagonal line $A = B$. The centre corresponds to where both soft-drop masses roughly equal $M_{\mathrm{H}}$ (it is always slightly below $M_{\mathrm{H}}$). It is this primary distribution that the analysis targets, hence why it is contained within the $S_n$ mass regions. The lengths of the $S_n$ mass regions, along the diagonal line $A = B$, are set so that this distribution populates at least two of the $S_n$ mass regions. The $S_n$ mass regions are designed to contain this distribution for signal models with values of $M_{\mathrm{H}}$ up to 125 GeV (the upper limit of the Higgs boson mass in this search). Typically, around

50% of signal events fall within the collection of $S_n$ mass regions. As the value of $M_H$ decreases from 40 GeV, the main distribution starts to fall outside the lowest signal mass region, $S_1$, leading to a rapid decrease in signal sensitivity. The reason why the mass regions do not probe lower soft-drop masses is explained at the end of this subsection.

The smallest distinct distribution, of signal events in the 2D mass plane, occurs when both of the AK8 jets have a soft-drop mass just above 0 GeV. A low soft-drop mass AK8 jet suggests it is without substructure. Although the selected AK8 jets are the best candidates for the $H \rightarrow b\bar{b}$ decays, there are a few reasons why this can occur. The primary reason is that the soft-drop grooming algorithm can incorrectly remove hard PF particles arising from the $b\bar{b}$ decays (see Appendix C). Another is that the PF particles resulting from one of the b-quarks are not contained in the jet. A final reason is that the jet simply does not originate from a $H \rightarrow b\bar{b}$ decay.

The two remaining distinct distributions occur when one of the AK8 jets has a soft-drop mass around $M_H$, but the other AK8 jet has a soft-drop mass just above 0 GeV. These distributions primarily lie outside of the 2D mass regions, however, where they merge into the main circular distribution there is signal contamination in the $U_n$ and $D_n$ sideband mass regions (a similar effect can be seen, but to a lesser extent, when one of the the soft-drop masses is significantly greater than $M_H$). Ideally the $U_n$ and $D_n$ mass regions would not contain any signal events, as they are used to estimate the QCD background. This is not a problem, however, for a few reasons. Firstly, the amount of contamination in the sideband mass regions is far exceeded by the yield in the signal mass regions. Secondly, the fit which estimates the QCD yield operates simultaneously on a signal mass region and the corresponding sidebands, allowing the signal contamination to be taken into account. Finally, because the mass region indices corresponding to the sideband contamination are different to those where the yields in the signal mass regions are highest, and because the QCD prediction is done independently for each mass region index, the sideband contamination actually becomes a discriminating feature of the signal processes.

Figure 6.7: Distribution of signal events in the 2D soft-drop mass plane, with the mass regions overlaid. The sub-figures correspond to $M_\mathrm{H}$ values of 30, 40, 50, 70, 90, and 125 GeV, respectively. The value of $M_\mathrm{SUSY}$ is 2000 GeV for all the plots. The events have passed the tag double-b-tag requirement and the kinematic cuts have been applied with $H_\mathrm{T} \in 3500+$ GeV.

The mass regions provide a powerful discrimination against the QCD multi-jet background. The distribution of the QCD MC events in the 2D soft-drop mass plane, with the 30 mass regions overlaid, is shown in Figure 6.8. For the majority of events, at least one of the AK8 jets is evaluated to have a low soft-drop mass because jets arising from QCD do not have substructure. Consequently, the distribution lies near the origin and along the graph axes. Only around 5% of the QCD events fall within the set of $S_n$ mass regions, with a higher density of events in the lower mass regions. For the W+jets and Z+jets processes, this fraction increases slightly because one of the AK8 jets can originate from the vector boson, if it is boosted and decays hadronically. The background process with the highest fraction of events entering the set of $S_n$ mass regions is $t\bar{t}$. This is because both AK8 jets can reconstruct a W boson, or the entire top quark. The fraction is 19% for the lowest $H_{\mathrm{T}}$ bin and 15% for the other $H_{\mathrm{T}}$ bins. It should be noted that the background processes populate all of the $S_n$ mass regions whereas the signal processes primary populate only two or three.

The reason why the mass regions do not probe to lower soft-drop masses, and why the sideband regions $U_1$ and $D_1$ have a unique shape (to avoid spanning lower soft-drop masses), is due to the data driven estimation of the QCD background. The method is described in full in Section 6.3. One of the key requirements is that the density of QCD events in a given $S_n$ mass region is similar to the density in its corresponding sidebands, $U_n$ and $D_n$. From Figure 6.8, it can be seen that the QCD event density is roughly proportional to the inverse of the product of the two fatJet soft-drop masses. This gives a higher average QCD event density in the $U_n$ and $D_n$ mass regions compared to the corresponding $S_n$ mass region. The disparity increases the closer the $S_n$ region is to the origin because the rate of change of the event density increases as the soft-drop masses tend towards zero. This effect is accounted for in the QCD estimation method, but to ensure that the method remains robust the mass regions do not probe to lower soft-drop masses.

Figure 6.8: Scatter plot of the QCD MC events in the 2D soft-drop mass plane, with the mass regions overlaid. The 2016 MC samples are used. The events have passed the tag double-b-tag requirement and the kinematic cuts have been applied with $H_T \in$ 1500-2500 GeV.

## 6.2.4  Simulated Results

The cut flow tables, for the event selection requirements described above, are provided for MC processes corresponding to 2017 data taking. The 2017 MC is used because the background processes have smaller event weightings. Table 6.2 contains the cut flows for the background processes, Table 6.3 contains the cut flows of signal models with $M_{\mathrm{SUSY}} = 2000$ GeV and different $M_{\mathrm{H}}$ values, whereas Table 6.4 is for signal models with $M_{\mathrm{H}} = 70$ GeV and different $M_{\mathrm{SUSY}}$ values. The cut flow tables branch into three sections, one for each of the $H_{\mathrm{T}}$ bins. They do not branch into separate sections to represent the 10 different $S_n$ regions. Instead, they just show the number of events contained within all of the $S_n$ bins. Despite the fact that this analysis searches for an all hadronic final state, the tables show that the combination of the kinematic cuts, the tag double-b-tag region, and the mass regions, suppress the background to manageable levels whilst retaining a large fraction of signal events. In order to examine the differences between the 2016 and 2017 MC, a cut flow table comparing equivalent 2016 and 2017 signal samples is provided in

Table 6.5. The energy scale of the jets is, on average, slightly lower in the 2017 MC. This is evident as, in the 2017 samples, a smaller fraction of events pass the kinematic cuts. In addition, the 2016 signal MC meets the tag double-b-tag requirement more efficiently than the 2017 signal MC.

The final yields are visualized by mapping the three different $H_\mathrm{T}$ bins and the 10 different mass region indices onto a 30 bin histogram. The bins 1-10 represent the ten different mass region indices, in ascending order, for $H_\mathrm{T} \in$ 1500-2500 GeV. The bins 11-20 represent the ten different mass region indices for $H_\mathrm{T} \in$ 2500-3500 GeV, and bins 21-30 do so for $H_\mathrm{T} \in$ 3500+ GeV. All the figures using this format in this thesis will explicitly state whether the mass index represents mass regions of type $S_n$ (signal mass regions) or of type $U_n + D_n$ (sideband mass regions). The figures will also explicitly state which type of double-b-tag region is being used (new double-b-tag regions are introduced in the QCD estimation methodology).

An example of this 1D representation can be seen in Figure 6.9, which shows the simulated yields for the 2017 signal and background. The primary background process is QCD, although it is not as dominant as it was after only applying the kinematic event selection. In the higher $H_\mathrm{T}$ regions, the QCD yields are greatly reduced. The QCD populations also decrease as the mass index increases. This is because, as the mass index increases, the average QCD event density in the signal mass regions decreases faster than the area of the mass regions increases. The $t\bar{t}$ process is also a significant background, especially in the bins corresponding to the higher mass indices. This is because the selected AK8 jets can be from the decay products of the entire top quarks or their daughter W bosons, and hence can have high soft-drop masses. The yields from the Z+jets and W+jets processes are small compared to those from QCD and $t\bar{t}$. Consequently, the estimation of the vector boson backgrounds is not a critical part of the analysis. Although the acceptance of the di-boson backgrounds is greater than that of the Z+jets and W+jets processes, they are neglected in this analysis due to their low production cross sections.

Figure 6.9: Simulated yields for the 2017 signal and background. The signal samples are parameterized by $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 1200$, 1600, and 2000 GeV. The tag double-b-tag region and the $S_n$ mass regions are used.

The three signal samples in Figure 6.9 are parameterized by $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 1200$, 1600 and 2000 GeV. Each sample has peaks in the bins corresponding to the $S_3$ and $S_4$ mass regions. Due to its large production cross section, the $M_\mathrm{SUSY} = 1200$ GeV sample has peaks that stand out clearly above the background in all three $H_\mathrm{T}$ regions. The $M_\mathrm{SUSY} = 2000$ GeV sample illustrates how the $H_\mathrm{T}$ binning provides sensitivity to signal processes with larger $M_\mathrm{SUSY}$ values. As the $H_\mathrm{T}$ region index increases, the background yields diminish significantly whilst the signal yield increases, and in the final $H_\mathrm{T}$ region the signal has a clear peak over the background.

Figure 6.10 shows the simulated yields for different sets of signal samples. It can be seen how the signal peaks shift and broaden for larger values of $M_\mathrm{H}$, and how the sizes of the peaks get smaller as the value of $M_\mathrm{SUSY}$ increases beyond 2000 GeV.

Figure 6.10: Simulated yields for the 2017 signal and background. The signal samples used for sub-figures (a)-(d) each have $M_{SUSY}$ = 1200, 1600, and 2000 GeV but with $M_H$ = 30, 50, 90, and 125 GeV, respectively. The signal sample used for sub-figure (e) has $M_H$ = 70 GeV and $M_{SUSY}$ = 2200, 2400, and 2600 GeV. The tag double-b-tag region and the $S_n$ mass regions are used for all the sub-figures.

Table 6.2: Cut flow table for each of the 2017 MC background samples. The yields are scaled to an integrated luminosity of 41.53 fb$^{-1}$, that of the 2017 data set.

| | QCD ($H_\mathrm{T} > 1000\mathrm{GeV}$) | $t\bar{t}$ | $Z \to q\bar{q}$ ($H_\mathrm{T} > 800\mathrm{GeV}$) | $W \to q\bar{q}$ ($H_\mathrm{T} > 800\mathrm{GeV}$) |
|---|---|---|---|---|
| Before Cuts | 46,817,302.86 | 34,535,101.11 | 776,177.01 | 1,411,986.00 |
| Pre-selection | 45,992,514.74 | 2,913,434.04 | 737,117.80 | 1,354,392.10 |
| 2*AK8Jet $p_\mathrm{T} > 300\mathrm{GeV}$ | 34,394,415.50 | 463,625.93 | 435,134.13 | 859,488.91 |
| 1*AK4Jet $p_\mathrm{T} > 300\mathrm{GeV}$ | 1,834,684.83 | 17,571.50 | 15,807.14 | 29,846.64 |
| $H_\mathrm{T} \in 1500\text{-}2500$ GeV | 716,668.20 | 9,943.47 | 7,441.15 | 13,656.53 |
| tag double-b-tag | 24,257.78 | 1,281.73 | 476.02 | 476.14 |
| within $S_n$ mass regions | 1,157.34 | 245.67 | 36.24 | 32.50 |
| $H_\mathrm{T} \in 2500\text{-}3500$ GeV | 40,561.19 | 898.99 | 671.89 | 1,158.28 |
| tag double-b-tag | 1,221.62 | 92.34 | 37.03 | 39.31 |
| within $S_n$ mass regions | 47.89 | 13.55 | 1.89 | 1.75 |
| $H_\mathrm{T} \in 3500+$ GeV | 3,345.60 | 82.15 | 67.71 | 126.33 |
| tag double-b-tag | 91.18 | 7.10 | 3.77 | 3.67 |
| within $S_n$ mass regions | 5.56 | 0.96 | 0.20 | 0.35 |

Table 6.3: Cut flow table for 2017 signal samples with different values of the Higgs boson mass. The $M_\mathrm{SUSY}$ parameter is fixed at 2000 GeV. The yields are scaled to an integrated luminosity of 41.53 fb$^{-1}$, that of the 2017 data set.

| | $M_\mathrm{H}$ | | | | | |
|---|---|---|---|---|---|---|
| | 30 GeV | 50 GeV | 70 GeV | 90 GeV | 110 GeV | 125 GeV |
| Before Cuts | 284.89 | 278.36 | 266.80 | 251.77 | 212.13 | 127.64 |
| Pre-selection | 284.81 | 278.32 | 266.75 | 251.70 | 212.08 | 127.61 |
| 2*AK8Jet $p_\mathrm{T} > 300$ GeV | 254.84 | 248.36 | 237.12 | 222.77 | 187.20 | 112.15 |
| 1*AK4Jet $p_\mathrm{T} > 300$ GeV | 223.81 | 218.56 | 208.75 | 195.79 | 164.20 | 98.19 |
| $H_\mathrm{T} \in 1500\text{-}2500$ GeV | 11.37 | 10.63 | 10.14 | 9.43 | 8.15 | 4.75 |
| tag double-b-tag | 6.60 | 5.86 | 5.19 | 4.52 | 3.72 | 2.00 |
| within $S_n$ mass regions | 0.18 | 2.90 | 2.65 | 2.19 | 1.80 | 0.90 |
| $H_\mathrm{T} \in 2500\text{-}3500$ GeV | 79.69 | 77.84 | 74.21 | 69.76 | 58.53 | 35.11 |
| tag double-b-tag | 49.88 | 45.24 | 40.04 | 35.46 | 27.53 | 15.92 |
| within $S_n$ mass regions | 1.55 | 23.49 | 20.99 | 17.95 | 13.60 | 7.49 |
| $H_\mathrm{T} \in 3500+$ GeV | 132.59 | 129.94 | 124.30 | 116.50 | 97.43 | 58.29 |
| tag double-b-tag | 78.60 | 70.96 | 62.89 | 54.89 | 42.92 | 24.35 |
| within $S_n$ mass regions | 4.21 | 36.31 | 31.58 | 26.97 | 20.58 | 11.18 |

Table 6.4: Cut flow table for 2017 signal samples with different values of the SUSY mass scale. The $M_{\mathrm{H}}$ parameter is fixed at 70 GeV. The yields are scaled to an integrated luminosity of 41.53 fb$^{-1}$, that of the 2017 data set.

| | $M_{\mathrm{SUSY}}$ | | | | |
| | 800 GeV | 1200 GeV | 1600 GeV | 2000 GeV | 2400 GeV |
|---|---|---|---|---|---|
| Before Cuts | 189,648.13 | 14,507.85 | 1,769.60 | 266.80 | 44.10 |
| Pre-selection | 188,449.12 | 14,490.79 | 1,768.88 | 266.75 | 44.09 |
| 2*AK8Jet $p_{\mathrm{T}} > 300$ GeV | 100,469.73 | 10,828.11 | 1,489.42 | 237.12 | 40.30 |
| 1*AK4Jet $p_{\mathrm{T}} > 300$ GeV | 58,642.81 | 8,795.73 | 1,288.84 | 208.75 | 35.46 |
| $H_{\mathrm{T}} \in$ 1500-2500 GeV | 43,250.87 | 5,021.70 | 213.84 | 10.14 | 0.67 |
| tag double-b-tag | 21,729.20 | 2,734.27 | 114.28 | 5.19 | 0.33 |
| within $S_n$ mass regions | 10,024.44 | 1,390.93 | 58.29 | 2.65 | 0.17 |
| $H_{\mathrm{T}} \in$ 2500-3500 GeV | 4,539.91 | 3,269.09 | 774.86 | 74.21 | 5.75 |
| tag double-b-tag | 2,011.49 | 1,707.73 | 421.48 | 40.04 | 3.07 |
| within $S_n$ mass regions | 762.37 | 828.83 | 219.75 | 20.99 | 1.59 |
| $H_{\mathrm{T}} \in$ 3500+ GeV | 466.28 | 406.55 | 297.82 | 124.30 | 29.03 |
| tag double-b-tag | 183.02 | 193.31 | 149.46 | 62.89 | 14.51 |
| within $S_n$ mass regions | 61.43 | 84.25 | 71.37 | 31.58 | 7.38 |

Table 6.5: Cut flow table comparing equivalent 2016 and 2017 signal MC samples. Three different SUSY mass scales are used, as indicated in the table. The $M_{\mathrm{H}}$ parameter is always 70 GeV. All the yields are scaled to an integrated luminosity of 41.53 fb$^{-1}$ to allow a direct comparison between the samples.

| $M_{\mathrm{SUSY}}$ | 800 GeV | | 1600 GeV | | 2400 GeV | |
| Year | 2016 | 2017 | 2016 | 2017 | 2016 | 2017 |
|---|---|---|---|---|---|---|
| Before Cuts | 189,648.13 | 189,648.13 | 1,769.60 | 1,769.60 | 44.10 | 44.10 |
| Pre-selection | 188,600.91 | 188,449.12 | 1,768.27 | 1,768.88 | 44.08 | 44.09 |
| 2*AK8Jet $p_{\mathrm{T}} > 300$ GeV | 104,986.48 | 100,469.73 | 1,521.49 | 1,489.42 | 40.90 | 40.30 |
| 1*AK4Jet $p_{\mathrm{T}} > 300$ GeV | 62,554.07 | 58,642.81 | 1,316.94 | 1,288.84 | 35.86 | 35.46 |
| $H_{\mathrm{T}} \in$ 1500-2500 GeV | 45,772.48 | 43,250.87 | 218.66 | 213.84 | 0.65 | 0.67 |
| tag double-b-tag | 24,340.76 | 21,729.20 | 121.08 | 114.28 | 0.33 | 0.33 |
| within $S_n$ mass regions | 11,467.85 | 10,024.44 | 63.59 | 58.29 | 0.17 | 0.17 |
| $H_{\mathrm{T}} \in$ 2500-3500 GeV | 4,799.47 | 4,539.91 | 784.93 | 774.86 | 5.70 | 5.75 |
| tag double-b-tag | 2,176.53 | 2,011.49 | 453.07 | 421.48 | 3.18 | 3.07 |
| within $S_n$ mass regions | 866.74 | 762.37 | 234.67 | 219.75 | 1.66 | 1.59 |
| $H_{\mathrm{T}} \in$ 3500+ GeV | 479.36 | 466.28 | 310.50 | 297.82 | 29.51 | 29.03 |
| tag double-b-tag | 187.37 | 183.02 | 164.01 | 149.46 | 15.45 | 14.51 |
| within $S_n$ mass regions | 59.94 | 61.43 | 77.95 | 71.37 | 7.74 | 7.38 |

# 6.3 Data Driven QCD Estimation

QCD multi-jet processes are the dominant background in this analysis due to the fully hadronic final state of the signal model. In this section, a data driven approach to estimate the QCD background is introduced and then tested using a variety of techniques. Due to the importance of correctly estimating this background, all the tests of the methodology are presented for both the 2016 and 2017 datasets.

## 6.3.1 QCD Estimation Method

In this analysis there are 30 search regions. They are indexed by the symbol $i$. The search regions arise from the combination of the three $H_T$ bins in the kinematic event selection and the 10 different mass region indices. Within the 30 search regions there are further categorisations based on mass region type ($U_n$, $S_n$, and $D_n$) and the double-b-tag requirement (tag and anti-tag; described below). Signal model events primarily exist in the $S_n$ mass regions with tag double-b-tag. It is this categorisation of events for which the QCD yield needs to be estimated. This is done, independently for each search region, using the data yields with different mass region types and/or different double-b-tag requirements. These yields are largely free of signal contamination.

The QCD estimation calculation is provided in Equation 6.2. Here is an example symbol for clarity of notation: the symbol $\hat{U}_i^{\text{tag}}$ represents the number of events meeting the tag double-b-tag requirement in the $i^{\text{th}}$ search region with $U_n$ mass region type. Note that a mass region symbol with a hat corresponds to an event yield, whereas a mass region symbol without a hat represents the actual space in the 2D soft-drop mass plane.

$$\hat{S}_i^{\text{tag PRED}} = \frac{\hat{S}_i^{\text{anti-tag}}}{\hat{U}_i^{\text{anti-tag}} + \hat{D}_i^{\text{anti-tag}}} \cdot (\hat{U}_i^{\text{tag}} + \hat{D}_i^{\text{tag}}) \tag{6.2}$$

In Equation 6.2, a new type of double-b-tag region is introduced; the anti-tag double-b-tag region. This region corresponds to where both AK8 jets have a double-b-tag discriminator less than 0.3. In Figure 6.11, the anti-tag double-b-tag region

Figure 6.11: Scatter plot of signal events in the double-b-tag plane. The signal sample is from 2016 MC with the parameters $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. Full kinematic cuts are applied with $H_\mathrm{T} \in 3500+$ GeV. In addition, both AK8 jets are required to have soft-drop masses greater than 15 GeV. The red triangle represents the tag double-b-tag region, the red square represents the anti-tag double-b-tag region, and the combination of the grey rectangles represent the control double-b-tag region.

is marked by the red square. The distribution in the figure is an example of how signal events occupy the 2D double-b-tag plane. It can be seen that only a very small fraction of signal events enter the anti-tag double-b-tag region.

To help communicate how Equation 6.2 works, first consider the simplified scenario where the QCD event density is the same within each signal mass region and its corresponding sidebands. In this case, the yield $\hat{S}_i^\mathrm{tag}$ would be equal to the yield $\hat{U}_i^\mathrm{tag} + \hat{D}_i^\mathrm{tag}$, because the area of $U_n + D_n$ equals that of $S_n$ (except for the search regions $i = 1, 11,$ and $21$, which use the first mass region where the sidebands have a unique shape). However, in reality, the average QCD event density is different between a signal mass region and its corresponding sideband regions. The difference increases the closer these mass regions are to the origin, as was discussed in Section 6.2.3. The factor $F_i$, defined in Equation 6.3, is introduced to the QCD estimation calculation to account for this disparity. Using events that meet the anti-tag

double-b-tag requirement, it is the ratio of the yield in the signal mass regions to the yield in the corresponding sidebands.

$$F_i \equiv \frac{\hat{S}_i^{\text{anti-tag}}}{\hat{U}_i^{\text{anti-tag}} + \hat{D}_i^{\text{anti-tag}}} \tag{6.3}$$

If there was zero correlation between the soft-drop mass and the double-b-tag discriminator of an AK8 jet, the factor $F_i$ would work perfectly. This is because the (normalised) QCD distribution in the 2D soft-drop mass plane would be independent of the double-b-tag requirements imposed on the AK8 jets. The relationship between the soft-drop mass and the double-b-tag discriminator, and the impact this has on the QCD estimation method, is explored in Section 6.3.2.

The anti-tag double-b-tag region used to calculate the $F_i$ factors is dominated by QCD events and, thus, there is no need to compensate for the other backgrounds, or signal contamination, in the data. However, when evaluating $\hat{U}_i^{\text{tag}} + \hat{D}_i^{\text{tag}}$, the other backgrounds, and signal contamination, cannot be ignored. In the fits used to obtain results, the values used for $\hat{U}_i^{\text{tag}} + \hat{D}_i^{\text{tag}}$ correspond to the number events remaining in data after the other backgrounds, and signal contamination, are subtracted. This is explained in greater detail in Section 7.2.1. In Appendix E.1, the background and signal yields are compared for the different search region categories that are used in the QCD estimation calculation.

## 6.3.2 Double-b-tag Variation with Soft-drop Mass

The QCD estimation method is theoretically correct when there is no correlation between the soft-drop mass and the double-b-tag discriminator of the AK8 jets. The double-b-tag BDT is trained without explicitly using any jet mass information. Some of the discriminating variables do, however, have a small dependence on the soft-drop mass. In addition, the training is only conducted on jets with soft-drop masses greater than 40 GeV. Consequently, there is a correlation between the soft-drop mass and the double-b-tag discriminator. This correlation, and the impact it has on the QCD estimation method, is examined in this subsection.

All the tests in this subsection are performed on data in an event selection space designed to be dominated by QCD events. The requirements for this space, applied after pre-selection, are:

- Both of the selected AK8 jets must have $p_\mathrm{T} > 300$ GeV.

- $1500 < H_\mathrm{T} < 2500$ GeV.

- For the selected AK8 jet labelled 'fatJetB', the double-b-tag discriminator must be less than 0.3.

The QCD dominated event selection space corresponds to the first $H_\mathrm{T}$ bin of the ordinary kinematic event selection, but with the AK4 jet requirement removed to increase the number of QCD events. The double-b-tag cut on fatJetB is applied to suppress any signal contamination. Following these cuts, the soft-drop mass and double-b-tag distributions of the other selected AK8 jet, fatJetA, can be studied. This is because these distributions are free of any significant contributions from the other backgrounds or potential signal (see Appendix E.2) and because these distributions have minimal dependence on the double-b-tag cut (and lack of soft-drop mass cut) applied to fatJetB (see Appendix E.3).

The graphs at the top of Figure 6.12 show the normalised double-b-tag discriminator distributions of fatJetA, in different soft-drop mass bins, for 2016 and 2017 data. The lowest soft-drop mass bin corresponds to the lowest masses used in the 2D mass regions and the highest bin corresponds to soft-drop masses just large enough to be used in the training of the double-b-tagger. There is a distinction between the double-b-tag discriminator shapes for 2016 and 2017 data, due to the different trainings of the BDT, but the general trends remain the same. The shape of the discriminator has a clear dependence on the soft-drop mass. The biggest difference occurs at the lowest double-b-tag discriminator values. There is also a bump in the distribution for discriminator values of around zero, a feature that grows in size for lower soft-drop masses. This occurs because the lower the soft-drop mass, the more

Figure 6.12: Top: Normalised double-b-tag discriminator distribution of fatJetA, in different soft-drop mass bins, in the QCD dominated event selection space for 2016 data (left) and 2017 data (right). Bottom: As above but with only two bins.

likely it is that the BDT will not recognise the jet and return a central discriminator value. For the QCD estimation method, what is critical is not how the shape of the double-b-tag discriminator varies with soft-drop mass, but simply if the ratio of events falling into the tag and anti-tag double-b-tag regions is changing (ideally $\hat{S}_i^{\text{tag}}\big/\hat{S}_i^{\text{anti-tag}} = (\hat{U}_i^{\text{tag}} + \hat{D}_i^{\text{tag}})\big/(\hat{U}_i^{\text{anti-tag}} + \hat{D}_i^{\text{anti-tag}}) \; \forall\, i)$. The 1D equivalent of the anti-tag double-b-tag region is an AK8 jet with a double-b-tag discriminator less than 0.3. The tag double-b-tag region cannot be represented by a single AK8 jet, due to its triangular shape in the 2D double-b-tag plane. It can, however, be roughly thought of as the case where the double-b-tag discriminator is greater than 0.3. The plots at the bottom of Figure 6.12 are the same as those above, but binned only in these two categories. The ratio of the two categories does vary with soft-drop mass, but it does so in a relatively steady manner. This is because the largest differences in the double-b-tag discriminator shape all occur for values less than 0.3.

Figure 6.13: Normalised soft-drop mass distribution of fatJetA in the cases where the double-b-tag discriminator of fatJetA is greater than, and less than, 0.3. This is done within the QCD dominated event selection space for 2016 data (left) and 2017 data (right).

In Figure 6.13, the normalised soft-drop mass distributions of fatJetA are shown in both the case where the double-b-tag discriminator of fatJetA is greater than 0.3, and where it is less than 0.3 (the same two categories that are described above). In the ideal case, where the double-b-tag discriminator and the soft-drop mass of the AK8 jets are uncorrelated, the ratio between these distributions would be constant. However, in reality, it varies as a function of the soft-drop mass of fatJetA. As the soft-drop mass increases from 15 GeV to 200 GeV (the masses spanned by the 2D mass regions), the ratio continually falls, smoothly, at an ever decreasing rate. This initially suggests that the $F_i$ factors will not give the correct values, as they require the (normalised) soft-drop mass distributions of the tag and anti-tag double-b-tag regions to be the same. However, because a given 2D signal mass region has sidebands that are larger in one mass dimension and smaller in the other, the variation is roughly balanced out when evaluating $F_i$. A more thorough check of this is now presented.

The correlation between the soft-drop mass and the double-b-tag discriminator cause the values obtained for $F_i$ to be incorrect by the factor $C_i$, given in Equation 6.4. The term $\rho_h^{\text{(anti-)tag}}(m_a, m_b)$ is the event density in the 2D soft-drop mass plane for QCD events after full kinematic cuts, in the $h^{\text{th}}$ $H_{\text{T}}$ bin ($h = 1, 2, 3$), that meet the (anti-)tag double-b-tag requirements.

$$C_i = \frac{\underset{S_n}{\iint} \mathrm{d}m_a \mathrm{d}m_b \cdot \rho_h^{\text{tag}}(m_a, m_b)}{\underset{U_n+D_n}{\iint} \mathrm{d}m_a' \mathrm{d}m_b' \cdot \rho_h^{\text{tag}}(m_a', m_b')} \Bigg/ \frac{\underset{S_n}{\iint} \mathrm{d}m_a'' \mathrm{d}m_b'' \cdot \rho_h^{\text{anti-tag}}(m_a'', m_b'')}{\underset{U_n+D_n}{\iint} \mathrm{d}m_a''' \mathrm{d}m_b''' \cdot \rho_h^{\text{anti-tag}}(m_a''', m_b''')} \tag{6.4}$$

where:  $i = 10\,(h-1) + n$

The $C_i$ factors cannot be evaluated in data. The calculation requires information about the soft-drop mass distribution in the tag double-b-tag region after kinematic event selection, and this was a blinded region in the analysis. However, using the QCD dominated event selection space, one can gain an insight into what the typical $C_i$ values should be (note that there will only be 10 values, rather than 30, because there is only one $H_\mathrm{T}$ bin in the QCD dominated event selection space). The event densities $\rho^{\text{tag}}(m_a, m_b)$ and $\rho^{\text{anti-tag}}(m_a, m_b)$ can be estimated as the 2D product of the fatJetA soft-drop mass distribution, where the fatJetA double-b-tag discriminator is greater than 0.3, and less than 0.3, respectively (due to its triangular shape, the tag double-b-tag region is only roughly represented this way). This approach is legitimate because the event density is invariant under $a \leftrightarrow b$ exchange and because, for QCD, the two soft-drop mass distributions have only a small dependence on each other (see Appendix E.3).

The following power law functional form was fit to the fatJetA soft-drop mass distributions:

$$\begin{aligned} f(m) = \ &p_0 + \frac{p_1}{(m-p_2)} + \frac{p_3}{(m-p_4)^2} + \frac{p_5}{(m-p_6)^3} + \frac{p_7}{(m-p_8)^4} \\ &+ p_9(m-p_{10}) + p_{11}(m-p_{12})^2 + p_{13}(m-p_{14})^3 + p_{15}(m-p_{16})^4 \end{aligned} \tag{6.5}$$

The fits can be seen overlaid on the distributions in Figure 6.13. Using the fits, the $C_i$ factors were calculated by applying Equation 6.4. For 2016, the $C_i$ factors range between 0.94 and 1.02 and for 2017, the $C_i$ factors range between 0.89 and 1.02. For both 2016 and 2017, the $C_i$ factors that deviate furthest from unity correspond to the 2D mass regions with the lowest soft-drop masses. This is because the double-

b-tag requirements have their greatest impact on the soft-drop mass distribution shape at low masses. It should be noted that $C_1$ does not follow this trend due to the unique shape of the $U_1$ and $D_1$ sideband mass regions.

Although these calculations of the $C_i$ factors do not have an associated uncertainty, the results suggest that, for both 2016 and 2017, the QCD estimation method is approximately correct to within 10%.

### 6.3.3   Testing Method with MC

In order to determine whether the QCD MC was correctly simulating the relationship between the soft-drop mass and the double-b-tag discriminator of the AK8 jets, the tests performed on data in Section 6.3.2 were repeated using the QCD MC. The same event selection criteria, which allows the data to be dominated by QCD events, were applied to the QCD MC. The distributions could then be compared to those acquired using real data. Figure 6.14 compares the ratio between the normalised soft-drop mass distributions of fatJetA in the cases where the double-b-tag discriminator of fatJetA is greater than, and less than, 0.3. For 2016, the QCD MC does a good job in emulating the data. There is a slight discrepancy at low soft-drop masses, but it occurs below the lowest mass used in the 2D mass regions. For 2017, the QCD MC does a very good job in emulating the data. These results mean that the QCD MC is a reliable sample for directly testing the QCD estimation method.

The test of the QCD estimation method, using QCD MC event yields as the inputs, is presented in Figure 6.15. The yields in the 30 search regions with tag double-b-tag and $S_n$ masses are compared to the number of events predicted by applying Equation 6.2. To the level of the statistical errors, the method works well for both 2016 and 2017. Note that the QCD MC has event weightings less than unity, so if one used real data, the statistical errors would be larger.

Figure 6.14: Ratio between the normalised soft-drop mass distributions of fatJetA in the cases where the double-b-tag discriminator of fatJetA is greater than, and less than, 0.3. This is done within the QCD dominated event selection space for data (red) and QCD MC (blue). The 2016 case is on the left and the 2017 case is on the right.



Figure 6.15: Comparison of the actual event yield with the predicted event yield, for events meeting the tag double-b-tag requirement and within the $S_n$ mass regions. This is done using the 2016 QCD MC (left) and the 2017 QCD MC (right).

### 6.3.4 Testing Method in a Control Region

In order to further test the QCD estimation method, it was used directly on data to estimate the yields in a control region, rather than the tag double-b-tag region. Equation 6.6 shows this modified version of the QCD estimation calculation.

$$\hat{S}_i^{\text{control PRED}} = F_i \cdot (\hat{U}_i^{\text{control}} + \hat{D}_i^{\text{control}})$$ (6.6)

The control double-b-tag region is defined as the space where the double-b-tag discriminator of one of the selected AK8 jets is between -1.0 and -0.4 and the other is between 0.3 and 0.8. This region is indicated by the grey rectangles in Figure 6.11. The control double-b-tag region exists outside the anti-tag double-b-tag region, but in a space that still has negligible signal contamination and only a small contribution from the other backgrounds. This can be seen in Appendix E.1. Requiring an AK8 jet with a double-b-tag discriminator between 0.3 and 0.8 ensures the control region is comparable to the tag double-b-tag region in one of the dimensions.

The $F_i$ factors belonging to the highest $H_{\text{T}}$ region, when directly evaluated using data, are computed from regions that suffer from low statistics. To better determine the $F_i$ factors, a new method of calculating their values was developed. This method evaluates the $F_i$ factors by performing the following calculation:

$$F_i = \iint_{S_n} \mathrm{d}m_a \mathrm{d}m_b \cdot \rho_h^{\text{anti-tag}}(m_a, m_b) \Big/ \iint_{U_n + D_n} \mathrm{d}m'_a \mathrm{d}m'_b \cdot \rho_h^{\text{anti-tag}}(m'_a, m'_b)$$ (6.7)

$$\text{where:} \quad i = 10\,(h-1) + n$$

The term $\rho_h^{\text{anti-tag}}(m_a, m_b)$ is the event density in the 2D soft-drop mass plane for QCD events after full kinematic cuts, in the $h^{\text{th}}$ $H_{\text{T}}$ bin, that meet the anti-tag double-b-tag requirements. The term $\rho_h^{\text{anti-tag}}(m_a, m_b)$ can be estimated as the 2D product of the corresponding one dimensional soft-drop mass distributions. As was stated in Section 6.3.2, this is legitimate because the event density is invariant under

Figure 6.16: Top: Normalised soft-drop mass distributions of fatJetA, and the fits to these distributions, in the anti-tag double-b-tag region for 2016 data and 2016 QCD MC. The kinematic cuts have been applied and the three figures correspond to $H_{\mathrm{T}} \in$ 1500-2500, 2500-3500, and 3500+ GeV, respectively. Bottom: As above, but for the 2017 data and 2017 QCD MC.

$a \leftrightarrow b$ exchange and because, for QCD, the two soft-drop mass distributions have only a small dependence on each other (see Appendix E.3).

Following the application of the kinematic cuts and the anti-tag double-b-tag requirement, the one dimensional soft-drop mass distributions were obtained from one of the selected AK8 jets, fatJetA. These distributions are shown in Figure 6.16 for data (which is dominated by QCD events, see Appendix E.1) and QCD MC. For each year, there are three graphs corresponding to the three $H_{\mathrm{T}}$ bins in the kinematic event selection. The functional form presented in Equation 6.5 was fit to all the soft-drop mass distributions. The fits can be seen overlaid on the distributions in Figure 6.16. The soft-drop mass range used in the fits covers the masses spanned by the 2D mass regions. The functions acquired from the fits to data were then used to calculate the $F_i$ factors. The exception is the $H_{\mathrm{T}} \in$ 3500+ GeV bin, where due to low statistics at high soft-drop masses the fits to the QCD MC were used instead, for both the 2016 and 2017 datasets.

In Figure 6.17, the new calculations of $F_i$ are compared to the original values evaluated directly from data and QCD MC. It should be noted that for search regions $i = 1$, 11, and 21, the $F_i$ value is around a factor of two larger because it accounts for

Figure 6.17: Top: Comparison of the new 2016 calculations of $F_i$ with the original values evaluated directly from 2016 data (left) and 2016 QCD MC (right). Bottom: As above, but for the 2017 quantities. For 2016 data, $F_{21} = 5.6$. This value, and its error bar, are not contained within the graph.

the unique shape of the mass sidebands. It can be seen, in both 2016 and 2017, how the original values of $F_i$ evaluated directly from data suffer from low statistics in the higher $H_{\mathrm{T}}$ regions. In the first $H_{\mathrm{T}}$ region, the variation of $F_i$ is relatively smooth and the values have small error bars. However, this degrades in the second $H_{\mathrm{T}}$ region, and in the final $H_{\mathrm{T}}$ region the values of $F_i$ vary greatly. The same behaviour is also exhibited by the original values of $F_i$ evaluated directly from the 2016 and 2017 QCD MC, although to a lesser extent as these samples have event weightings that are less than unity. In contrast, the new calculations of $F_i$ vary smoothly for all $H_{\mathrm{T}}$ regions. The $F_i$ values decrease when they represent 2D mass regions closer to the origin. This is because the average QCD event density in the mass sidebands, relative to the signal mass regions, increases closer to the origin.

In the $H_{\mathrm{T}} \in$ 1500-2500 GeV region, the new $F_i$ calculations are given an error of 0.15, which are uncorrelated between the different search regions. With this error,

Figure 6.18: Comparison of the actual event yield with the predicted event yield, for events meeting the control double-b-tag requirement and within the $S_n$ mass regions. This is done using the 2016 data (left) and the 2017 data (right).

the new values of $F_i$ agree with the original values evaluated directly from data and QCD MC, in both 2016 and 2017. This verifies the new method of calculating $F_i$ as, in the first $H_{\mathrm{T}}$ region, the statistical errors of the original $F_i$ factors are small. In the $H_{\mathrm{T}} \in$ 2500-3500 GeV region, the error given to the new $F_i$ factors is also 0.15. For both 2016 and 2017, these new $F_i$ calculations offer an improvement over the original values evaluated directly from data, which suffer from statistical fluctuations. In the $H_{\mathrm{T}} \in$ 3500+ GeV region, the new $F_i$ factors are given an error of 0.30. The size of error is increased due to the lower statistics in the corresponding 1D mass fits. For both years, the new $F_i$ calculations offer a significant improvement over the original values evaluated directly from data, which suffer from large statistical fluctuations. It should be noted that, when comparing the new calculations of the $F_i$ factors, there is very little variation between the 2016 and 2017 values. This demonstrates the stability of the method.

When performing the QCD estimation method on data, the $F_i$ factors derived from the new calculation are used for all three $H_{\mathrm{T}}$ regions. In Figure 6.18, the QCD estimation method is tested on data using the control region. The yields in the 30 search regions with control double-b-tag and $S_n$ mass region type are compared to the number of events predicted by applying Equation 6.6. The QCD estimation method shows good agreement, for both 2016 and 2017. It does so across all the search regions, even though the yields can differ by over a factor of 1000.

### 6.3.5 Summary of QCD Estimation Method

In this section, a data driven method for estimating the QCD yield has been presented. It has be shown to work successfully, using a variety of different approaches. In the $H_\mathrm{T} \in 3500+$ GeV region, the method is limited by low statistics in $\hat{U}_i^\mathrm{tag} + \hat{D}_i^\mathrm{tag}$. This is the most significant uncertainty associated to the background prediction in this analysis. In contrast, the uncertainties on the $F_i$ factors do not have a large impact on the analysis, as can be seen in Appendix E.4.

# 7

---

# Systematic Uncertainties and Final Results

---

## 7.1  Systematic Uncertainties

This analysis uses MC simulation to estimate the signal model yields and the $t\bar{t}$, Z+jets, and W+jets background yields. Associated to these MC samples are a variety of systematic uncertainties. Each systematic represents the uncertainty of a different aspect of the simulation relative to real events in data. As a given aspect of the MC simulations is varied, the event selection efficiency can change and/or the event populations can shift between the different $H_{\mathrm{T}}$ bins and 2D mass regions. This results in a change of signal sensitivity. In the fits used to obtain the final results, the systematic uncertainties are represented in the calculations by nuisance parameters. The information on how these fits work is provided in Section 7.2. There are two versions of each systematic uncertainty. One version corresponds to the 2016 MC and the other corresponds to the 2017 MC. As each systematic is discussed, it will be explicitly stated whether the uncertainty is taken to be correlated, or uncorrelated, between the two years.

In this section, all the systematic uncertainties considered are described, along with information about how they impact the analysis.

## 7.1.1 Double-b-tag Systematics

The CMS b-tagging Physics Object Group (POG) provide $p_\mathrm{T}$ dependent data/MC scale factors for AK8 jets with double-b-tag discriminators greater than a collection of working point values. There are a set of signal scale factors for AK8 jets originating from true $b\bar{b}$ topologies. These are derived using $g\to b\bar{b}$ jets and can be found in Table 7.1 and Table 7.2, for 2016 and 2017, respectively. Additionally, there are a set of mistag scale factors for AK8 jets originating from the merged hadronic decay of a W boson, or the entire top quark, in $t\bar{t}$ events. These can be found in Table 7.3 and Table 7.4, for 2016 and 2017, respectively. Each of the scale factors has an associated error, which leads to systematic uncertainties on the signal model and $t\bar{t}$ event yields. The working points available are 0.3 (loose), 0.6 (medium-1), 0.8 (medium-2), and 0.9 (tight). The tables do not show the scale factors for the tight working point as it is not used in this analysis.

Table 7.1: The 2016 double-b-tag scale factors for AK8 jets arising from true $b\bar{b}$ pairs. The scale factors are provided for three different working points within different $p_\mathrm{T}$ ranges.

| | Double-b-tag working point: | | |
|---|---|---|---|
| AK8 jet $p_\mathrm{T}$ range | 0.3 (Loose) | 0.6 (Med-1) | 0.8 (Med-2) |
| 250-350 GeV | $0.96^{+0.03}_{-0.02}$ | $0.93^{+0.03}_{-0.02}$ | $0.92^{+0.03}_{-0.03}$ |
| 350-430 GeV | $1.00^{+0.04}_{-0.03}$ | $1.01^{+0.03}_{-0.03}$ | $1.01^{+0.03}_{-0.04}$ |
| 430-840 GeV | $1.01^{+0.02}_{-0.04}$ | $0.99^{+0.02}_{-0.04}$ | $0.92^{+0.03}_{-0.05}$ |
| 840+ GeV | $1.01^{+0.04}_{-0.08}$ | $0.99^{+0.04}_{-0.08}$ | $0.92^{+0.06}_{-0.10}$ |

Table 7.2: The 2017 double-b-tag scale factors for AK8 jets arising from true $b\bar{b}$ pairs. The scale factors are provided for three different working points within different $p_\mathrm{T}$ ranges.

| | Double-b-tag working point: | | |
|---|---|---|---|
| AK8 jet $p_\mathrm{T}$ range | 0.3 (Loose) | 0.6 (Med-1) | 0.8 (Med-2) |
| 250-350 GeV | $0.96^{+0.03}_{-0.03}$ | $0.93^{+0.04}_{-0.03}$ | $0.85^{+0.04}_{-0.04}$ |
| 350-840 GeV | $0.95^{+0.06}_{-0.04}$ | $0.90^{+0.08}_{-0.04}$ | $0.80^{+0.07}_{-0.04}$ |
| 840+ GeV | $0.95^{+0.12}_{-0.08}$ | $0.90^{+0.16}_{-0.08}$ | $0.80^{+0.14}_{-0.08}$ |

Table 7.3: The 2016 double-b-tag scale factors for AK8 jets mis-tagged in $t\bar{t}$ events. The scale factors are provided for three different working points within different $p_{\mathrm{T}}$ ranges.

| AK8 jet $p_{\mathrm{T}}$ range | Double-b-tag working point: | | |
|---|---|---|---|
| | 0.3 (Loose) | 0.6 (Med-1) | 0.8 (Med-2) |
| 250-350 GeV | 1.044±0.028 | 1.029±0.034 | 1.050±0.044 |
| 350-430 GeV | 1.074±0.052 | 1.156±0.064 | 1.086±0.078 |
| 430-700 GeV | 1.119±0.079 | 1.156±0.064 | 1.086±0.078 |
| 700+ GeV | 1.119±0.158 | 1.156±0.128 | 1.086±0.156 |

Table 7.4: The 2017 double-b-tag scale factors for AK8 jets mis-tagged in $t\bar{t}$ events. The scale factors are provided for three different working points within different $p_{\mathrm{T}}$ ranges.

| AK8 jet $p_{\mathrm{T}}$ range | Double-b-tag working point: | | |
|---|---|---|---|
| | 0.3 (Loose) | 0.6 (Med-1) | 0.8 (Med-2) |
| 250-350 GeV | $0.939^{+0.026}_{-0.026}$ | $0.922^{+0.027}_{-0.027}$ | $0.875^{+0.030}_{-0.030}$ |
| 350-430 GeV | $1.007^{+0.055}_{-0.054}$ | $0.967^{+0.057}_{-0.056}$ | $0.939^{+0.063}_{-0.063}$ |
| 430+ GeV | $0.996^{+0.080}_{-0.078}$ | $0.902^{+0.083}_{-0.081}$ | $0.893^{+0.091}_{-0.089}$ |

Scale factor weightings are applied to signal and $t\bar{t}$ events in which the two selected AK8 jets meet the tag double-b-tag requirement. The applied weight is not trivially determined because the tag double-b-tag region, which is triangular in shape, is not fully supported by the single AK8 jet scale factors provided. These scale factors only support rectangular shapes in the 2D double-b-tag plane, which correspond to applying the scale factor weighting independently on each AK8 jet. One corner of the rectangle must have co-ordinates composed of double-b-tag working point values and the opposite corner must have the co-ordinates (1,1).

The method devised to determine the scale factor weighting for events in the tag double-b-tag region is illustrated in Figure 7.1. Two rectangular regions, 'Q' and 'Y', for which scale factors can be assigned, are used to approximate the triangular shape of the tag double-b-tag region. Region-Q corresponds to where both double-b-tag discriminators are greater than the medium-1 working point. It is a good approximation as only a small fraction of the space exists outside the tag double-b-tag region. Events falling outside region-Q are given a scale factor corresponding to region-Y, which covers most of the remaining area within the tag double-b-tag

Figure 7.1: Scatter plot of signal events in the double-b-tag plane (where both AK8 jets have double-b-tag discriminators greater than 0.3). The parameters of the signal sample are $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV. Full kinematic cuts are applied with $H_{\mathrm{T}} \in 3500+$ GeV. In addition, both AK8 jets are required to have soft-drop masses greater than 15 GeV. The red triangle represents the tag double-b-tag region. The orange square, region-Q, corresponds to where both double-b-tag discriminators are greater than the medium-1 working point. The yellow rectangles, region-Y, correspond to where one double-b-tag discriminator is between the loose and medium-1 working points whilst the other is greater than the medium-2 working point.

region. Region-Y is the space where one double-b-tag discriminator is between the loose and medium-1 working points whilst the other is greater than the medium-2 working point. This approximation is not as good because a sizeable fraction of region-Y exists outside the tag double-b-tag region. This is not considered a large problem, however, as most of the signal model events are contained in region-Q.

The scale factors corresponding to region-Y cannot be acquired directly, as the rectangles do not have corners with co-ordinates (1,1). In order to evaluate them, two new regions, 'X' and 'Z', are introduced, which have scale factors that can be drawn directly. They are shown, alongside (one half of) region-Y, in Figure 7.2. Region-X(Z) is the space where one double-b-tag discriminator is greater than the medium-1 working point whilst the other is greater than the medium-2 (loose) working point.

Figure 7.2: Left: Region-X and (half of) region-Y in the double-b-tag plane. Right: Region-Z in the double-b-tag plane. For both sub-figures the tag double-b-tag region is indicated by the red triangle.

The scale factor for region-Z, written as a function of the scale factors for region-X and region-Y, is provided in Equation 7.1. The terms $f_X$ and $f_Y$ are the fraction of events, belonging to region-Z, that are within region-X and region-Y, respectively.

$$s_Z = f_X \cdot s_X + f_Y \cdot s_Y \qquad (7.1)$$

The expression can be manipulated so that the scale factor for region-Y is written in terms of the known quantities:

$$s_Y = \frac{1}{f_Y} \cdot (s_Z - f_X \cdot s_X) \qquad (7.2)$$

The fractions $f_X$ and $f_Y$ are attained after applying the kinematic event selection (with $H_T > 1500$ GeV, rather than binning in the quantity). For signal events, the fractions used are $f_X = 0.87$ and $f_Y = 0.13$. The fractions have a slight dependence on the $M_H$ and $M_{SUSY}$ parameters. This variation, however, is only around the percent level and is ignored. For $t\bar{t}$ events, the fractions used are $f_X = 0.68$ and $f_Y = 0.32$. The fractions do change depending on whether the simulation is for 2016 or 2017. The change, however, is minimal and is not accounted for.

An example of the effect that $\pm 1\sigma$ variations, in the double-b-tag systematic uncertainty, have on the signal yield for 2016 MC, can be seen in Figure 7.3. The variations cause a uniform scaling in all the search region bins. A $+1\sigma$ increase in

Figure 7.3: The effect that $\pm 1\sigma$ variations, in the double-b-tag scale factors, have on the 2016 signal yield with parameters $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield, the blue points denote the yield with $+1\sigma$ variation in the scale factors, and the green points denote the yield with $-1\sigma$ variation in the scale factors.

the double-b-tag scale factors corresponds to around a 10% increase in yield. A $-1\sigma$ decrease in the double-b-tag scale factors corresponds to around a 15% decrease in yield. In 2017, the $+1\sigma$ uncertainties, for AK8 jets arising from true $b\bar{b}$ pairs, are significantly larger and correspond to around a 25% increase in signal yield. The uncertainties on the double-b-tag scale factors are taken to be uncorrelated between 2016 and 2017. This is because the CMS pixel detector was upgraded for 2017 operations and because the double-b-tag BDT was re-trained.

To evaluate the impact that the double-b-tag systematic uncertainty has on the analysis, the ordinary expected limits in the $M_\mathrm{H}$-$M_\mathrm{SUSY}$ plane are compared to those where the nuisance parameters corresponding to the systematic uncertainty, for both 2016 and 2017, are frozen in the fit. This comparison is shown in Figure 7.4. In the case where the frozen systematic uncertainties are those originating from the signal topology scale factors, a slight difference can be seen between the expected limits. For $M_\mathrm{H} = 70$ GeV, there is a shift of 5 GeV in the expected $M_\mathrm{SUSY}$ limit. This is small compared to the spread of the expected limit itself. In the other case,

Figure 7.4: A comparison of the ordinary expected limits, using both 2016 and 2017 data, in the $M_{\mathrm{H}}$-$M_{\mathrm{SUSY}}$ plane (black lines) with the expected limits where nuisance parameters corresponding to the double-b-tag systematic uncertainty are frozen in the fit (red lines). Left: The systematic uncertainties originating from the signal topology scale factors, for both 2016 and 2017, are frozen. Right: The systematic uncertainties originating from the $t\bar{t}$ mistag scale factors, for both 2016 and 2017, are frozen.

where the frozen systematic uncertainties are those originating from the $t\bar{t}$ mistag scale factors, the difference in the expected limits is even smaller. For $M_{\mathrm{H}} = 70$ GeV, the shift in the expected $M_{\mathrm{SUSY}}$ limit is less than 1 GeV. Note that an explanation of how the exclusion plots are obtained is provided in Section 7.2.

These results show that this analysis is stable against the double-b-tag systematic uncertainty. Due to this stability, the double-b-tag scale factors applied, which only represent an approximation of the tag double-b-tag region, are acceptable to use.

## 7.1.2 Soft-drop Mass Systematics

The soft-drop mass of the AK8 jets has two associated systematic uncertainties; one for the soft-drop mass scale (JMS) and another for the soft-drop mass resolution (JMR). These soft-drop mass uncertainties are applied to all the MC samples in the analysis.

To implement the soft-drop mass scale systematic uncertainty, a factor, corresponding to $\pm 1\sigma$ variations, is applied to the soft-drop masses to shift them up and down. These factors are provided by the CMS JetMET POG. For both 2016 and 2017 MC, the factors are $1.0000 \pm 0.0094$. The factors were derived from boosted $W \rightarrow q\bar{q}$
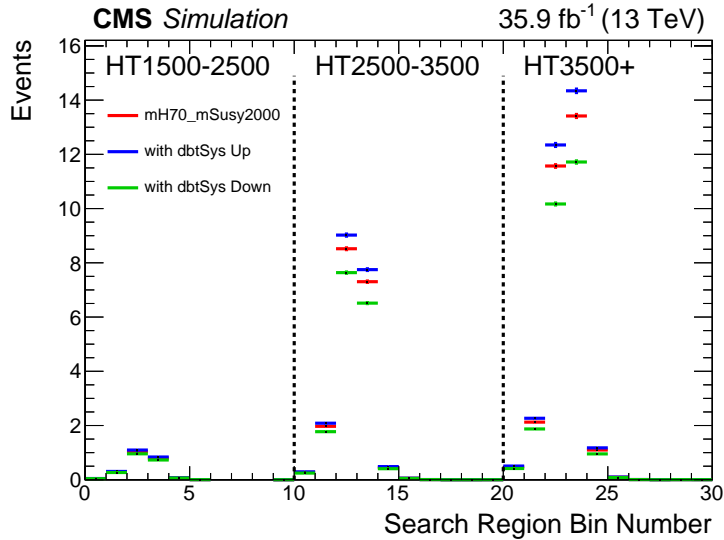
Figure 7.5: The effect that $\pm 1\sigma$ variations in the soft-drop mass scale have on the 2016 signal yield with parameters $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield, the blue points denote the yield with $+1\sigma$ variation in the soft-drop mass scale, and the green points denote the yield with $-1\sigma$ variation in the soft-drop mass scale.

AK8 jets in semileptonic $t\bar{t}$ decays using 2016 data/MC, and the current recommendation is to use the same factor for 2017. Because of this, the uncertainties are taken to be correlated between the two years. The soft-drop mass scale variations slightly change the shape of the final signal yield, an example of which can be seen in Figure 7.5. The total number of events is roughly conserved, but when the scale is shifted up/down the population in the bins corresponding to the higher 2D mass regions (i.e. the mass regions that exist beyond the centre of the primary soft-drop mass distribution) increases/decreases. The impact that the soft-drop mass scale systematic uncertainty has on the analysis is very small, as can be seen in Figure 7.6. When the nuisance parameter corresponding to the systematic uncertainty is frozen in the fit, the shift induced in the expected $M_\mathrm{SUSY}$ limit is less than 1 GeV.

The uncertainty on the resolution of the soft-drop mass is implemented by adding a stochastic smearing term to the nominal soft-drop mass, corresponding to a $1\sigma$ variation. This creates an asymmetrical systematic uncertainty, as the smearing can only degrade the resolution. The smearing value, following the CMS JetMET

Figure 7.6: A comparison of the ordinary expected limits, using both 2016 and 2017 data, in the $M_{\mathrm{H}}$-$M_{\mathrm{SUSY}}$ plane (black lines) with the expected limits where the nuisance parameter corresponding to the soft-drop mass scale systematic uncertainties, for both 2016 and 2017, is frozen in the fit (red lines).

POG instructions, is provided in Equation 7.3. The expression $\mathcal{N}(\mu, \sigma)$ represents a random number drawn from a Gaussian distribution with mean, $\mu$, and standard deviation (SD), $\sigma$.

$$c = \mathcal{N}(0, \sigma_{\mathrm{JMR}}) \cdot \sqrt{s_{\mathrm{JMR}}^2 - 1} \tag{7.3}$$

The soft-drop mass resolution SD, $\sigma_{\mathrm{JMR}}$, takes the value 10.1 GeV and the soft-drop mass resolution scale factor, $s_{\mathrm{JMR}}$, takes the value 1.20. The values of $\sigma_{\mathrm{JMR}}$ and $s_{\mathrm{JMR}}$ were both derived using 2016 data/MC and the current recommendation is to use the same values for 2017. Because of this, the soft-drop mass resolution systematic uncertainties are taken to be correlated between the two years. The value of $\sigma_{\mathrm{JMR}}$ was derived with the primary interest of reconstructing boosted W bosons and is provided as an absolute error. As a consequence, its value of 10.1 GeV is an overestimate when dealing with soft-drop masses significantly below 80 GeV. Taking a conservative approach, however, this value is used for soft-drop masses as low as 30.3 GeV. Below this threshold, $\sigma_{\mathrm{JMR}}$ is set to equal one third of the nominal soft-drop mass. This prevents low soft-drop masses being smeared to negative values.

Figure 7.7 explores how the soft-drop mass smearing changes the signal model populations in the 2D mass regions. When the smearing is applied, the distributions broaden in the 2D mass plane. As a consequence, the population decreases in the $S_n$ mass regions containing the central part of the distribution. Some of the events move into the $S_n$ mass regions containing the tails of the distribution (see Figure 7.8) and some others are moved into the corresponding sideband mass regions. This effect leads to a loss in signal sensitivity. The impact is most significant for signal models with low $M_H$ values because the size of the mass smearing, relative to the width of the primary mass distribution, is larger.

The soft-drop mass resolution uncertainty is the systematic error that has the largest impact on the analysis. In Figure 7.9, the ordinary expected limits are once again compared to those where the nuisance parameter corresponding to the systematic uncertainty is frozen in the fit. The largest differences occur for signal models with low $M_H$ values, for the reasons stated above. For $M_H = 40$ GeV, the shift in the expected $M_{SUSY}$ limit is 32 GeV. When $M_H = 70$ GeV, this shift is reduced to 11 GeV and when $M_H = 125$ GeV, the shift is only 1 GeV. The reason why the two sets of expected limits rejoin at $M_H = 30$ GeV, is that the majority of the soft-drop mass distribution no longer exists within the 2D mass regions.

## 7.1.3   Jet Energy Systematics

This analysis uses both AK4 jets and AK8 jets in the event selection. Jet energy corrections (JECs) are applied to the $p_T$ of both sets of jets. The JEC systematic uncertainty arises as there is an associated error for the corrections applied to each jet. For all the MC samples, these uncertainties are propagated through the analysis. The JECs, and the corresponding uncertainties, are derived by the CMS JetMET POG [40]. They are provided for both AK4 jets and AK8 jets, and for both 2016 and 2017. The same JECs are applied to the AK4 and AK8 jets, therefore the uncertainties between both types of jet are correlated. Furthermore, the JEC uncertainties are also taken to be correlated between 2016 and 2017, because the method used to derive the JECs remained the same between both years.

Figure 7.7: Distribution of signal events in the 2D soft-drop mass plane, for 2016 MC, with the mass regions overlaid. The events have passed the tag double-b-tag requirement and the kinematic cuts have been applied with $H_\mathrm{T} \in 3500+$ GeV. Left column: Nominal soft-drop mass resolution. Right column: Soft-drop mass resolution after smearing, which corresponds to a $+1\sigma$ variation in the systematic uncertainty. Top row: $M_\mathrm{H} = 50$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. Middle row: $M_\mathrm{H} = 90$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. Bottom row: $M_\mathrm{H} = 125$ GeV and $M_\mathrm{SUSY} = 2000$ GeV.

Figure 7.8: The effect that the soft-drop mass smearing has on the 2016 signal yield with parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield and the blue points denote the yield after soft-drop mass smearing, which corresponds to a $+1\sigma$ variation in the systematic uncertainty.



Figure 7.9: A comparison of the ordinary expected limits, using both 2016 and 2017 data, in the $M_{\mathrm{H}}$-$M_{\mathrm{SUSY}}$ plane (black lines) with the expected limits where the nuisance parameter corresponding to the soft-drop mass resolution systematic uncertainties, for both 2016 and 2017, is frozen in the fit (red lines).

Figure 7.10: The effect that $\pm 1\sigma$ variations in the JECs have on the 2016 signal yield with parameters $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield, the blue points denote the yield with $+1\sigma$ variation in the JECs, and the green points denote the yield with $-1\sigma$ variation in the JECs.

An example of the effect that $\pm 1\sigma$ variations, in the JEC systematic uncertainties, have on the final signal yield can be seen in Figure 7.10. There is a slight difference in shape due to the displacement created in the $H_\mathrm{T}$ distribution. When the JECs are shifted up, some events in the $H_\mathrm{T} \in 2500\text{-}3500$ GeV bin move to the $H_\mathrm{T} \in 3500+$ GeV bin, and when the JECs are shifted down, the opposite occurs. For larger values of $M_\mathrm{SUSY}$, the fraction of signal events in the highest $H_\mathrm{T}$ bin increases, and because the bulk of the $H_\mathrm{T}$ distribution is no longer near the bin boundary, the size of this effect diminishes. The variation in the JECs has a minute effect on the requirement of having two AK8 jets and a separated AK4 jet. This is because the average $p_\mathrm{T}$ of these jets is much larger than the cut of 300 GeV applied to them (even more so when $M_\mathrm{SUSY}$ has a high value), so slight shifts in the jet energy scale barely change this aspect of the selection efficiency. The impact that the JEC systematic uncertainties have on the analysis is very small. When the nuisance parameter corresponding to the systematic uncertainty is frozen in the fit, the shift induced in the expected $M_\mathrm{SUSY}$ limit is less than 1 GeV. Due to this small change, the plot is not shown.

Figure 7.11: The effect that $\pm 1\sigma$ variations, in the jet energy resolution, have on the 2016 signal yield with parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield, the blue points denote the yield with $+1\sigma$ variation in the JERs, and the green points denote the yield with $-1\sigma$ variation in the JERs.

There is another source of systematic uncertainty in the $p_{\mathrm{T}}$ of the AK4 jets and AK8 jets. The jet energy resolution (JER) in data is worse than it is in simulation and, consequently, the jets in MC need to be smeared to describe the data. This smearing is applied to each jet according to the CMS JetMET POG prescription. The smearing comes with an associated error, which gives rise to the systematic uncertainty. The uncertainty is correlated between the AK4 and AK8 jets. It is treated as uncorrelated between 2016 and 2017 as the derivation, which remains the same, is statistically limited in both years. The impact that the JER systematic uncertainty has on the signal yields is negligible, as can be seen in Figure 7.11. Consequently, there is no difference in the expected limits when the nuisance parameters corresponding to the JER systematics are frozen in the fit.

### 7.1.4  Prefire Scale Factors and Uncertainty

Prefiring, described in Section 4.5.2, was a problem that occurred in the Level-1 trigger during 2016 and 2017, where events could self veto if they deposited a significant amount of energy in the $2.0 < |\eta| < 3.0$ region of the ECAL. In order to

account for this problem, all the MC events in the analysis are reweighted by a factor representing the probability that they would not prefire if they had occurred in data. This is achieved by looping over all the offline photons and AK4 jets in an event, and multiplying the probabilities of each object not causing an event to prefire, as shown in Equation 7.4. Note that in the case where an offline photon and jet spatially overlap, the maximum prefiring probability is taken.

$$\omega = 1 - P(\text{prefire}) = \prod_{i=\text{photons, jets}} \left(1 - \epsilon_i^{\text{prefire}}(\eta, p_{\text{T}}^{(\text{EM})})\right) \tag{7.4}$$

In Equation 7.4, the term $\epsilon_i^{\text{prefire}}(\eta, p_{\text{T}}^{(\text{EM})})$ represents the probability that an offline photon or AK4 jet would cause an event to prefire. It is parameterized as a function of $\eta$ and the $p_{\text{T}}$ deposited in the ECAL. These probabilities were derived, by the CMS Level-1 Trigger POG, using a subset of events that could not prefire. Such events occur when they are accepted by the Level-1 trigger exactly three bunch crossings after another triggered event. These events cannot self veto because, as was explained in Section 4.5.2, the two bunch crossings following a triggered event are not accepted. This subset of events represent only around 0.25% of the total data set. Figure 7.12 shows the 2D maps derived for $\epsilon_i^{\text{prefire}}(\eta, p_{\text{T}}^{(\text{EM})})$. Note that there is a distinction between the two years. In 2017 there was a greater probability of an object causing an event to prefire because the ECAL trigger primitives became further out of phase as time progressed.

The reduction in the number of signal events due to prefiring is not excessively large. Although the signal model produces a large number of high $p_{\text{T}}$ jets, they are not inherently forward. Figure 7.13 shows the distribution of prefire weights, corresponding to 2017 data taking, for signal samples with low and high $M_{\text{SUSY}}$ values. In both cases, the majority of events have a prefiring weight of unity and the average weighting is around 0.95. The $t\bar{t}$ distribution is also included in Figure 7.13, it has a very similar distribution to the signal samples. These prefiring weight distributions are for all events passing the pre-selection requirement. There is very

Figure 7.12: Probability of an object causing an event to prefire, parameterized as a function of $\eta$ and the $p_{\mathrm{T}}$ deposited in the ECAL. Top left: photons in 2016. Top right: photons in 2017. Bottom left: AK4 jets in 2016. Bottom right: AK4 jets in 2017.



Figure 7.13: Normalised distribution of prefire weights for different 2017 MC samples. Red: signal model with $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 800$ GeV. Blue: signal model with $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2800$ GeV. Green: $t\bar{t}$. All events have passed the pre-selection requirement.

Table 7.5: Average prefire weight for different MC samples in 2016 and 2017.

| | 2016 | 2017 |
|---|---|---|
| Signal: $M_H = 70$ GeV, $M_{\text{SUSY}} = 800$ GeV | 0.969 | 0.949 |
| Signal: $M_H = 70$ GeV, $M_{\text{SUSY}} = 2800$ GeV | 0.975 | 0.957 |
| $t\bar{t}$ | 0.963 | 0.939 |

little difference when the kinematic event selection is applied, however, as the cuts to do not prioritise forward jets. For the 2016 signal samples, the prefire weight averages lie closer to unity. This is because there was a lower probability of a given physics object causing an event to prefire. This can be seen in Table 7.5, which compares the average prefire weights of MC samples between 2016 and 2017. It should be noted that the event selection cut flow tables (Tables 6.2, 6.3, 6.4, and 6.5) do not have the prefire weights applied in order to avoid unnecessary complication.

There is a systematic uncertainty associated to the prefire weight applied to each MC event. This arises from the errors on the $\epsilon_i^{\text{prefire}}(\eta, p_{\text{T}}^{\text{(EM)}})$ factors. The error of each factor is the maximum between the statistical error and 20% of the prefiring probability. These errors change the average prefire weightings by less than 1% and, consequently, the impact on the final signal yield is minute, an example of which can be seen in Figure 7.14. The change in weighting is modest because the majority of the prefire weights are unity, or close to unity, and are composed of the smallest $\epsilon_i^{\text{prefire}}(\eta, p_{\text{T}}^{\text{(EM)}})$ factors, which have the smallest uncertainties. The systematic uncertainty is taken to be uncorrelated between 2016 and 2017. Due to their small size, there is no difference in the expected limits when the nuisance parameters corresponding to the prefire weight uncertainties are frozen in the fit.

For MC events that have more than one object with a significant probability of causing prefiring, the weights applied are slightly too high. This is because the method works by combining the probability that each object independently causes a prefire. It does not take into account di-object triggers, for example double e/$\gamma$, where the $p_{\text{T}}$ thresholds are reduced. This overestimation does not have a significant impact on the analysis. As can be seen in Figure 7.15, the fraction of signal events with two prefire candidates is only around 1%.

Figure 7.14: The effect that $\pm 1\sigma$ variations, in the prefire event weighting, have on the 2016 signal yield with parameters $M_H = 70$ GeV and $M_{SUSY} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield, the blue points denote the yield with $+1\sigma$ variation in prefire weighting, and the green points denote the yield with $-1\sigma$ variation in prefire weighting.
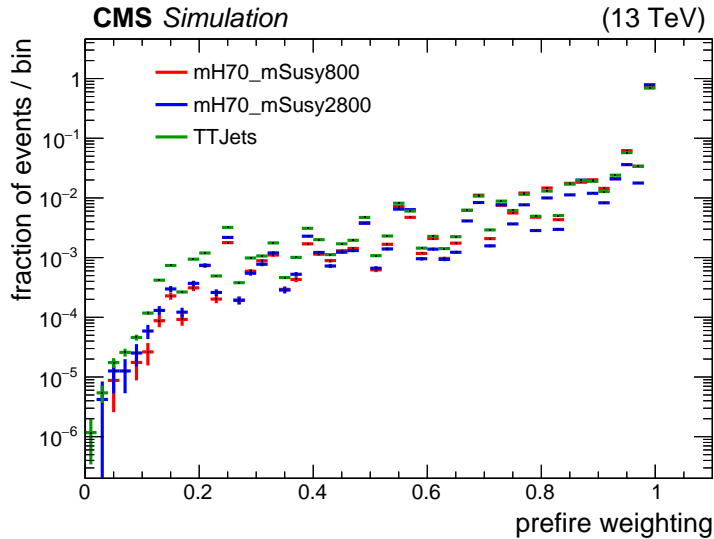


Figure 7.15: Normalised distribution of prefire candidates for different 2017 MC samples. Red: signal model with $M_H = 70$ GeV and $M_{SUSY} = 800$ GeV. Blue: signal model with $M_H = 70$ GeV and $M_{SUSY} = 2800$ GeV. Green: $t\bar{t}$. All events have passed the pre-selection requirement. A prefire candidate is defined as an AK4 jet with $2.25 < |\eta| < 3.00$ which deposits more than 30 GeV in the ECAL.

## 7.1.5   Other Systematics

In this subsection, some of the more common systematic uncertainties are described:

- **Luminosity:** To account for the uncertainty in the luminosity of the data collected, a normalisation uncertainty is applied to all MC events. An uncertainty of 2.5% is applied to the 2016 MC samples and, independently, an uncertainty of 2.3% is applied to the 2017 MC samples. When the nuisance parameters corresponding to the luminosity systematic uncertainties are frozen in the fit, the shift induced in the expected $M_{\mathrm{SUSY}}$ limit is less than 1 GeV. This is a negligible difference compared to the spread of the expected limit itself.

- **MC statistics:** The statistical uncertainties of all the MC samples are considered in this analysis. The signal samples contain around 200,000 (raw) events and, in the most populated search region bins, the statistical error is less than 1%. Consequently, the signal model statistical uncertainties have a negligible impact on the analysis. The MC backgrounds have large sample sizes with event weightings less than unity (except for the 2016 W+jets sample). However, due to a low selection efficiency, there are very few (raw) events that make it into the two highest $H_{\mathrm{T}}$ regions. This means that the relative statistical errors are large. But, because the event yields are low, this does not have a large impact on the analysis.

- **Background MC cross sections:** To account for the uncertainty on their cross sections and for any mismodelling of their selection efficiencies, a normalisation uncertainty of 50% is independently applied to the $t\bar{t}$, Z+jets, and W+jets background MC samples, for both 2016 and 2017. For the 2017 $t\bar{t}$ MC, which is divided up according to the decay mode, this uncertainty is taken to be correlated between each sample. There is no correlation between the 2016 and 2017 samples because, for each process, there is a key difference between the samples, e.g. different $H_{\mathrm{T}}$ thresholds applied during generation.

The cross section uncertainties are conservative estimates. However, despite their large size, the differences in the expected limits, when the corresponding nuisance parameters are frozen in the fit, are relatively small. When freezing the Z+jets and W+jets cross section systematic uncertainties (for both 2016 and 2017 simultaneously), the shift in the expected $M_{\mathrm{SUSY}}$ limit is less than 1 GeV. This is due to the very small yields expected from the Z+jets and W+jets processes. In the case of freezing the $t\bar{t}$ cross section systematic uncertainties, the change in the expected $M_{\mathrm{SUSY}}$ limit varies with the $M_{\mathrm{H}}$ parameter. For $M_{\mathrm{H}} = 70$ GeV, the shift in the expected $M_{\mathrm{SUSY}}$ limit is less than 1 GeV. However, for $M_{\mathrm{H}} = 125$ GeV, the shift increases to 6 GeV. This is because the fraction of background events due to the $t\bar{t}$ process is largest in the search regions corresponding to the highest 2D mass regions. It should be noted that even if the yields from the background MC processes are incorrect, they will be compensated for by the data driven QCD estimation method (albeit not quite correctly, as the $F_i$ factors correspond to the soft-drop mass distributions of QCD events). This helps reduce the impact caused by the uncertainty of the background MC processes.

The largest background estimated directly from MC is the $t\bar{t}$ background. A unique test of the 2016 $t\bar{t}$ MC sample can be found in Appendix F.

- **Initial state radiation (ISR) reweighting:** An ISR correction was derived, by another CMS analysis group, from $t\bar{t}$ events in the fully leptonic final state. The event selection requires two leptons (electrons or muons) and two b-tagged jets, implying that any other jets in the event arise from ISR. The correction factors are 1.000, 0.920, 0.821, 0.715, 0.662, 0.561, and 0.511 for $n_{\mathrm{ISR}} = 0$, 1, 2, 3, 4, 5, and 6+, respectively. These corrections are applied to the signal samples along with a normalisation factor to ensure the overall cross section remains constant. The normalisation factor is typically around 1.08. It has a slight dependence on the signal model parameters and on the simulation year. The systematic uncertainties used for the event weightings are set to be

half the deviation from unity for each correction factor. These uncertainties are treated as correlated between the 2016 and 2017 signal MC. The effect that $\pm 1\sigma$ variations, in the ISR reweighting systematic uncertainty, have on the final signal yield is very small, and diminishes for larger values of the $M_{\mathrm{SUSY}}$ parameter. Consequently, the difference in the expected limits when this systematic is removed from the fit is negligible.

## 7.1.6 Pile-up Reweighting

This analysis has very little dependence on pile-up due to the high energies required to pass the kinematic event selection. Consequently, pile-up reweighting is not performed on the MC samples.

A simple method was devised in order to prove that the analysis has very little dependence on pile-up. Each 2016 signal sample was divided into two new samples based on the number of pile-up interactions in each event. The low pile-up category was defined by $n_{\mathrm{PU}} \leq 23$ and the high pile-up category was defined by $n_{\mathrm{PU}} \geq 24$. Note that each new sample was weighted so that it corresponded to the full 2016 data set luminosity. The final yields of the low pile-up and high pile-up samples could then be compared. Figure 7.16 shows this comparison for the signal model with $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2200$ GeV. The difference between the two distributions is marginal. It has been shown earlier in this section that changes in signal yield of this size have a negligible effect on the expected limits. Furthermore, the differences between the two $n_{\mathrm{PU}}$ distributions in this test are far larger than those that would be produced by pile-up reweighting. Thus, the conclusion is that pile-up reweighting, and its associated errors, do not need to be applied in this analysis.

## 7.1.7 QCD Scale Reweighting

This analysis does not apply the systematic uncertainties due to variations in the QCD scale or the PDF. This approach is validated in this subsection by studying the QCD scale variation (which causes larger changes than variations in the PDF).

Figure 7.16: A comparison of the event yields for low pile-up and high pile-up signal samples, using 2016 MC. The parameters of the signal model are $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2200$ GeV. The events meet the tag double-b-tag requirement and are within the $S_n$ mass regions.

The QCD scale uncertainty was calculated by varying the renormalization and factorization scales, $\mu_R$ and $\mu_F$, by a factor of two. The $+1\sigma$ scale weightings correspond to where $\mu_R = 0.5$ and $\mu_F = 0.5$, and the $-1\sigma$ scale weightings correspond to where $\mu_R = 2.0$ and $\mu_F = 2.0$. Figure 7.17 shows the distributions of these two weightings for the signal model with parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV. It should be noted that the nominal QCD scale weighting is unity.

The impact that the QCD scale uncertainty has on the final signal yield, whilst maintaining the same production cross section, was evaluated as follows. The $\pm 1\sigma$ scale weightings were applied to each event. They were then normalised by dividing by the average of the $\pm 1\sigma$ scale weighting distributions before any cuts were applied. The final yields could then be compared to evaluate the impact of the systematic uncertainty. Figure 7.18 shows this comparison for the 2016 signal sample with parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV. There is only a minimal change in the shape and normalisation of the signal yield, which will have no impact on the expected limits. Consequently, it is reasonable to omit the QCD scale systematic uncertainty from the analysis and simply let it be absorbed it into the theoretical error on the production cross section.

Figure 7.17: Normalised distribution of the $\pm 1\sigma$ scale weightings for the centrally produced 2016 signal sample with parameters $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. No event selection requirements have been applied.



Figure 7.18: The effect that $\pm 1\sigma$ variations, in the QCD scale weighting, have on the (normalised) signal yield. The signal is the 2016 centrally produced sample with parameters $M_\mathrm{H} = 70$ GeV and $M_\mathrm{SUSY} = 2000$ GeV. All the events meet the tag double-b-tag requirement and are within the $S_n$ mass regions. The red points denote the nominal yield, the blue points denote the yield with $+1\sigma$ scale variation, and the green points denote the yield with $-1\sigma$ scale variation.

## 7.1.8 Conclusion of Systematics Study

In this section, all the different systematic uncertainties effecting the MC samples have been described. It has been shown that the expected limits are stable against variation in these systematic uncertainties. This is primarily due to the high rate at which the signal production cross section decreases as $M_{\mathrm{SUSY}}$ increases. For example, using Table 5.6, the production cross section is reduced by a factor of 2.44 between $M_{\mathrm{SUSY}} = 2400$ GeV and $M_{\mathrm{SUSY}} = 2600$ GeV. Consequently, a systematic uncertainty that changes a signal model yield by around 10% has only a small impact on the $M_{\mathrm{SUSY}}$ limit.

## 7.2 Results

This section begins by describing the likelihood function of the analysis, which forms an important part of the test statistic used to extract results. This is followed by an outline of how the upper limits on the production cross section are determined. Then, in accordance with how the analysis was conducted, the expected limits are presented before showing the observed limits.

## 7.2.1 Likelihood Function

The total likelihood function is described by Equation 7.5. It is the product of the individual likelihood functions for both the signal and sideband mass categories, across all 30 search regions, using both 2016 and 2017. The likelihood function is parameterized by $\mu$, the hypothesized signal strength parameter and $\boldsymbol{\theta}$, the collection of nuisance parameters (all of which are described below).

$$\mathcal{L}(\mu, \boldsymbol{\theta}) = \mathcal{L}_{\mathrm{constrain}}(\boldsymbol{\theta}) \cdot \prod_{y=2016}^{2017} \cdot \prod_{i=1}^{30} \cdot \prod_{m=S}^{U+D} \mathrm{Poisson}\Big(n_{y,i,m} \Big| b_{y,i,m}^{\mathrm{total}}(\boldsymbol{\theta}) + \mu \cdot s_{y,i,m}(\boldsymbol{\theta})\Big) \quad (7.5)$$

The term $n_{y,i,m}$ is the number of events observed in data for a given year $(y)$, search region $(i)$, and mass region type $(m)$. For a given signal model, $s_{y,i,m}$ is the expected number of signal events. The term $b_{y,i,m}^{\mathrm{total}}$ is the total background yield.

It is composed of a QCD yield and the expected number of events from the other background processes (which are determined using MC):

$$b_{y,i,m}^{\text{total}} = b_{y,i,m}^{\text{QCD}} + b_{y,i,m}^{\text{MC}} \tag{7.6}$$

The QCD yields are nuisance parameters determined during the fitting procedure. To implement the QCD estimation method, each signal and sideband mass region pair is coupled as follows:

$$b_{y,i,S}^{\text{QCD}} = F_{y,i} \cdot b_{y,i,U+D}^{\text{QCD}} \tag{7.7}$$

where $F_{y,i}$ is the $F_i$ factor from Section 6.3, but with an additional index to denote the year. The $F_{y,i}$ terms are independent nuisance parameters, as each one has an associated uncertainty. They are assigned Gaussian probability density functions which are contained within the $\mathcal{L}_{\text{constrain}}$ part of the likelihood function.

The parameters representing the QCD yields, $b_{y,i,m}^{\text{QCD}}$, fill the gap between the observed yield and the sum of the other background yields plus any potential signal. If the signal process did exist, and there was an excess of events in a given signal mass category, the corresponding QCD yield, $b_{y,i,S}^{\text{QCD}}$, would not absorb these excess events. This is because it is coupled to the QCD yield in the associated sideband mass category, $b_{y,i,U+D}^{\text{QCD}}$, which, by design, will contain far less signal events. This can be seen in Figure 7.19.

In some of the search regions belonging to $H_{\text{T}} \in 3500+$ GeV, no events are observed in at least one of the mass categories. In order to avoid having a QCD prediction of zero in these search regions, the QCD parameter for each sideband mass category, $b_{y,i,U+D}^{\text{QCD}}$, is constrained to have a minimum value of 0.25. The minimum value was determined using the QCD MC, which has event weightings less than unity. To ensure that the QCD MC was trustworthy for this task, in the lower $H_{\text{T}}$ bins, the QCD MC yields were compared to data in the sideband mass regions. The test showed reasonable agreement.

Figure 7.19: A comparison of the number of signal events, meeting the tag double-b-tag requirement, in the $S_n$ mass regions (red) and $U_n + D_n$ sideband mass regions (blue). The signal sample is 2016 MC and has the parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV.

In the likelihood function, the yields derived from MC, $s_{y,i,m}$ and $b_{y,i,m}^{\mathrm{MC}}$, can vary due to the collection of systematic uncertainties associated to the MC modelling, as was described in Section 7.1. To represent the yield variations, a nuisance parameter is assigned to each independent systematic uncertainty. The nuisance parameters follow probability density functions, constraining the extent to which the associated yields can vary about their estimated values. These probability densities form the rest of the $\mathcal{L}_{\mathrm{constrain}}$ term in the likelihood function.

The gamma distribution is used to, independently, describe the statistical uncertainties of all the MC yields [75]. The predicted yield, $Y$, has the following probability density function, where $N$ is the raw number of events in MC and $w$ is the event weighting:

$$p(Y) = \frac{1}{w} \frac{(Y/w)^N}{N!} \exp(-Y/w) \qquad (7.8)$$

The remaining systematic uncertainties follow the log-normal distribution [75]. Within a given year, they are correlated between the different mass categories, search regions, and MC samples (if applicable). In addition, some systematic uncertainties are correlated between 2016 and 2017 (see Section 7.1). Each independent system-

atic uncertainty has a nuisance parameter, $\theta_j$, with the following probability density function:

$$p(\theta_j) = \frac{1}{\sqrt{2\pi}} \exp(-\theta_j^2/2) \tag{7.9}$$

All the associated yields, $Y_k$, which are effected by the systematic uncertainty, are given by:

$$Y_k = \widetilde{Y_k} \cdot (\kappa_{jk})^{\theta_j} \tag{7.10}$$

where $\widetilde{Y_k}$ is the best prior estimate of the yield and $\kappa_{jk}$ characterises the width of the probability density function due to the systematic uncertainty, such that:

$$p(Y_k) = \frac{1}{\sqrt{2\pi}\, Y_k \ln(\kappa_{jk})} \exp\left(-\frac{(\ln(Y_k/\widetilde{Y_k}))^2}{2(\ln(\kappa_{jk}))^2}\right) \tag{7.11}$$

## 7.2.2 Setting Upper Limits

The profile likelihood ratio is used as the basis to test a hypothesized value of $\mu$:

$$\lambda(\mu) = \frac{\mathcal{L}(\mu, \hat{\hat{\boldsymbol{\theta}}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\boldsymbol{\theta}})} \tag{7.12}$$

The terms $\hat{\mu}$ and $\hat{\boldsymbol{\theta}}$ are the values of $\mu$ and $\boldsymbol{\theta}$ that maximize the likelihood function. The term $\hat{\hat{\boldsymbol{\theta}}}_\mu$ is the value of $\boldsymbol{\theta}$ that maximizes the likelihood function for a given value of $\mu$. The profile likelihood ratio exists in the range $0 \leq \lambda(\mu) \leq 1$, with the larger values representing better compatibility between the data and the hypothesized value of $\mu$.

The test statistic used to calculate an upper limit on $\mu$, is as follows:

$$q_\mu = \begin{cases} -2\ln\lambda(\mu) & \hat{\mu} \leq \mu, \\ 0 & \hat{\mu} > \mu, \end{cases} \tag{7.13}$$

The values of the test statistic satisfy $0 \leq q_\mu \leq \infty$. The closer the value is to zero, the better the compatibility between the data and the hypothesized value of $\mu$. As can be seen in Equation 7.13, the test statistic is set to zero when $\hat{\mu} > \mu$. This is

because, when setting upper limits, one does not want upwards fluctuations in data to represent an incompatibility with the hypothesized value of $\mu$.

To quantify the level of agreement between the data and the hypothesized value of $\mu$, the p-value is calculated as follows:

$$p_\mu = \int\limits_{q_{\mu,\mathrm{obs}}}^{\infty} dq_\mu \; f(q_\mu|\mu) \tag{7.14}$$

The term $q_{\mu,\mathrm{obs}}$ is the value of the test statistic observed using the data. The function $f(q_\mu|\mu)$ is the probability density function of $q_\mu$ for a hypothesized value of $\mu$. It is determined using the asymptotic formulae found in Ref. [76].

The 95% confidence level upper limit on the signal strength parameter is given by the value of $\mu$ which satisfies $p_\mu = 0.05$. This represents, for a given signal model, the 95% confidence level upper limit on $\sigma/\sigma_{\mathrm{theory}}$. For simplicity, it will henceforth be referred to as the upper limit. A signal model is considered to be excluded if the upper limit of $\sigma/\sigma_{\mathrm{theory}}$ is less than unity.

There are two types of upper limits obtained in this analysis; observed limits and expected limits. The observed upper limits are obtained as described above, using the yields observed in data to determine $q_{\mu,\mathrm{obs}}$. For the expected upper limits, the value of $q_{\mu,\mathrm{obs}}$ is estimated, assuming $\mu = 0$, without (directly) using the data. The expected upper limits allow one to assess the sensitivity of the analysis before looking at the observed limits. The median expected upper limits are obtained by using the median value of $f(q_\mu|\mu' = 0)$ as a best estimate of $q_{\mu,\mathrm{obs}}$ [76]. The $\pm 1\sigma$ variations are obtained by using the 16[th] and the 84[th] percentile points of $f(q_\mu|\mu' = 0)$. These represent how the observed upper limits are expected to vary due to statistical fluctuations in the data.

### 7.2.3   Expected Limits

The expected upper limits of $\sigma/\sigma_{\mathrm{theory}}$ were calculated for all the different signal models. Using linear interpolation, the limits in the whole $M_{\mathrm{H}}$-$M_{\mathrm{SUSY}}$ plane were

Figure 7.20: Expected signal model exclusions, as a function of $M_{\mathrm{H}}$ and $M_{\mathrm{SUSY}}$, for 2016 and 2017 data. This is for the scenario where the gluino mass is 1% higher than that of the squarks. The solid black line indicates the median expected excluded region. The dashed black lines indicate the expected excluded regions with $\pm 1\sigma$ in experimental uncertainty. The colour scale indicates the median expected 95% confidence level upper limit of $\sigma/\sigma_{\mathrm{theory}}$.

then calculated, as can be seen in Figure 7.20. The solid black line is the contour where the median expected upper limit of $\sigma/\sigma_{\mathrm{theory}}$ equals unity. It represents the expected boundary where the signal model is excluded, as a function of $M_{\mathrm{H}}$ and $M_{\mathrm{SUSY}}$. For a given $M_{\mathrm{H}}$ value on such a contour, the $M_{\mathrm{SUSY}}$ value is referred to as the $M_{\mathrm{SUSY}}$ limit. The dashed black lines are the contours where the $\pm 1\sigma$ variations in the expected upper limit of $\sigma/\sigma_{\mathrm{theory}}$ equal unity. For a given $M_{\mathrm{H}}$ value, the distance between the dashed lines in the $M_{\mathrm{SUSY}}$ dimension represents the $\pm 1\sigma$ spread that the observed $M_{\mathrm{SUSY}}$ limit should follow if the signal model did not exist.

In Figure 7.20, the expected limits of the $M_{\mathrm{SUSY}}$ parameter are approximately uniform within the range $40 < M_{\mathrm{H}} < 110$ GeV. The median expected $M_{\mathrm{SUSY}}$ limits are around 2500 GeV, and the $\pm 1\sigma$ variations give values of approximately 2400 GeV and 2600 GeV. Most of the sensitivity to these high values of $M_{\mathrm{SUSY}}$ comes from the $H_{\mathrm{T}} \in 3500+$ GeV region. In this $H_{\mathrm{T}}$ region, the yields observed in data are very low. It is these low statistics which provide the main contribution to the distribution width of the expected $M_{\mathrm{SUSY}}$ limit.

There are multiple effects occurring that change the expected $M_{\mathrm{SUSY}}$ limits as the $M_{\mathrm{H}}$ parameter is varied:

- As $M_{\mathrm{H}}$ decreases from 40 GeV, the primary distribution of signal events in the 2D soft-drop mass plane starts to move outside the 2D mass regions. This results in a rapid reduction in signal sensitivity and, thus, in the expected $M_{\mathrm{SUSY}}$ limit.

- As $M_{\mathrm{H}}$ increases, the $\mathrm{H} \to \mathrm{b\bar{b}}$ branching ratio decreases (see Table 5.5). This reduces the number of expected signal events and, therefore, reduces the expected $M_{\mathrm{SUSY}}$ limit. This is most noticeable for $110 < M_{\mathrm{H}} < 125$ GeV, where the expected $M_{\mathrm{SUSY}}$ limit decreases by around 150 GeV. This region corresponds to where the rate of change in the branching ratio is greatest.

- As $M_{\mathrm{H}}$ decreases, the tag double-b-tag requirement is achieved more efficiently (see Table 6.1). This increases the fraction of signal events selected and, thus, increases the signal sensitivity. The effect is relatively small and cannot be clearly seen in the expected $M_{\mathrm{SUSY}}$ limits. It may contribute, along with the changing $\mathrm{H} \to \mathrm{b\bar{b}}$ branching ratio, to the slight reduction in the expected $M_{\mathrm{SUSY}}$ limit in the range $40 < M_{\mathrm{H}} < 90$ GeV.

- As $M_{\mathrm{H}}$ increases, the signal events occupy the higher 2D mass regions. Here, there are less background events, leading to an increase in signal sensitivity. This effect cannot be directly seen in the expected $M_{\mathrm{SUSY}}$ limits. It does, however, oppose the two effects listed above, leading to the approximately uniform expected $M_{\mathrm{SUSY}}$ limits within the $40 < M_{\mathrm{H}} < 110$ GeV range.

### 7.2.4   Observed Limits

The yields observed in data did not contain any significant excesses with respect to the SM background predictions. This can be seen in Figure 7.21, which compares the data with the SM expectations, following a background only ($\mu = 0$) maximisation of the likelihood function.

Figure 7.21: Yields observed in data with tag double-b-tag and $S_n$ mass regions for 2016 (top) and 2017 (bottom). This is compared to the SM expectations, following a background only fit. The data points have Poisson error bars. The shaded region represents the systematic uncertainty on the total background. The red line represents the simulated yield of the signal model with $M_{\mathrm{H}} = 110$ GeV and $M_{\mathrm{SUSY}} = 2400$ GeV.

Table 7.6: The 2016 SM expectations, following a background only fit, and the 2016 observed yields in search region 4, for both the signal and sideband mass regions.

|  | $S$ mass region | $U + D$ mass region |
|---|---|---|
| QCD | $172.68 \pm 14.27$ | $208.46 \pm 15.17$ |
| $t\bar{t}$ | $9.07 \pm 3.64$ | $9.38 \pm 3.59$ |
| $Z \to q\bar{q}$ | $3.94 \pm 2.41$ | $7.47 \pm 4.24$ |
| $W \to q\bar{q}$ | $8.19 \pm 6.13$ | $0$ |
| total background | $193.87 \pm 13.20$ | $225.31 \pm 14.64$ |
| observed | $189$ | $230$ |

Table 7.7: The 2017 SM expectations, following a background only fit, and the 2017 observed yields in search region 4, for both the signal and sideband mass regions.

|  | $S$ mass region | $U + D$ mass region |
|---|---|---|
| QCD | $210.03 \pm 14.77$ | $223.86 \pm 17.45$ |
| $t\bar{t}$ | $13.11 \pm 5.85$ | $13.81 \pm 5.77$ |
| $Z \to q\bar{q}$ | $5.18 \pm 2.67$ | $9.78 \pm 4.79$ |
| $W \to q\bar{q}$ | $7.42 \pm 4.17$ | $5.81 \pm 3.15$ |
| total background | $235.74 \pm 13.69$ | $253.26 \pm 16.20$ |
| observed | $235$ | $254$ |

A more detailed evaluation of the SM expectations following the background only fit, using the fourth search region as an example, is provided in Table 7.6 and Table 7.7 for 2016 and 2017, respectively. The tables provide the expected yields, and the corresponding uncertainties, for each of the backgrounds considered in both the signal and sideband mass regions. It should be noted that the uncertainty on the total background is smaller than the uncertainty on the expected QCD yield. This is because there is an anti-correlation between the uncertainties on the QCD yield and the uncertainties of the other backgrounds, due to the data driven QCD estimation method.

In order to corroborate the values in Table 7.6 and Table 7.7, the pre-fit event yield information is provided in Table 7.8 and Table 7.9 for 2016 and 2017, respectively. For each systematic uncertainty, the factors by which the MC yields change under $\pm 1\sigma$ variations are provided. If the uncertainty is asymmetrical, the factors corresponding to both the $+1\sigma$ and $-1\sigma$ variations are given. It should be noted that the systematic uncertainties arising from the statistical error of the MC yields are not provided.

Table 7.8: Search region 4 event yield information (pre-fit), for both the signal and sideband mass regions, using 2016 data and MC. For each systematic uncertainty, the factors by which the MC yields change under $\pm 1\sigma$ variations are provided. The example signal sample has parameters $M_H = 70$ GeV and $M_{\mathrm{SUSY}} = 1200$ GeV. At the bottom of the table the $F_i$ factor, which is used for the QCD estimation, is also provided.

| Mass type | $S$ mass region | | | |
|---|---|---|---|---|
| Observed yield | 189 | | | |
| MC process | signal | TTJets | ZJets | WJets |
| Yield | 594.02 | 13.83 | 4.13 | 9.99 |
| Systematic | | | | |
| isrWeight$201^6_7$ | 0.984/1.016 | - | - | - |
| luminosity2016 | 1.025 | 1.025 | 1.025 | 1.025 |
| XS_TTJets2016 | - | 1.500 | - | - |
| XS_ZJets2016 | - | - | 1.500 | - |
| XS_WJets2016 | - | - | - | 1.500 |
| jecUnc$201^6_7$ | 1.017/0.984 | 0.955/1.130 | 1.000/1.056 | 1.000/1.000 |
| jerUnc2016 | 1.002/0.998 | 1.000/1.032 | 1.035/1.000 | 1.000/1.000 |
| jmsUnc$201^6_7$ | 0.941/1.047 | 1.025/0.880 | 1.092/0.958 | 1.000/1.000 |
| jmrUnc$201^6_7$ | 1.000/0.874 | 1.000/1.098 | 1.000/1.151 | 1.000/1.000 |
| SigDbtTag2016 | 0.922/1.050 | - | - | - |
| TtDbtTag2016 | - | 0.890/1.116 | - | - |
| prefire2016 | 1.005/0.995 | 1.009/0.991 | 1.004/0.996 | 1.000/1.000 |
| Mass type | $U + D$ mass region | | | |
| Observed yield | 230 | | | |
| MC process | signal | TTJets | ZJets | WJets |
| Yield | 8.12 | 14.36 | 8.32 | 0.00 |
| Systematic | | | | |
| isrWeight$201^6_7$ | 0.998/1.002 | - | - | - |
| luminosity2016 | 1.025 | 1.025 | 1.025 | 1.025 |
| XS_TTJets2016 | - | 1.500 | - | - |
| XS_ZJets2016 | - | - | 1.500 | - |
| XS_WJets2016 | - | - | - | 1.500 |
| jecUnc$201^6_7$ | 1.002/0.998 | 0.922/1.097 | 0.970/1.076 | - |
| jerUnc2016 | 1.016/0.994 | 0.984/0.984 | 1.024/1.046 | - |
| jmsUnc$201^6_7$ | 1.011/1.003 | 1.082/0.966 | 0.976/1.041 | - |
| jmrUnc$201^6_7$ | 1.000/2.644 | 1.000/1.070 | 1.000/0.834 | - |
| SigDbtTag2016 | 0.925/1.054 | - | - | - |
| TtDbtTag2016 | - | 0.890/1.116 | - | - |
| prefire2016 | 1.005/0.995 | 1.007/0.993 | 1.011/0.989 | - |
| $F_i$ factor | $0.946 \pm 0.150$ | | | |

Table 7.9: Search region 4 event yield information (pre-fit), for both the signal and sideband mass regions, using 2017 data and MC. For each systematic uncertainty, the factors by which the MC yields change under $\pm 1\sigma$ variations are provided. The example signal sample has parameters $M_H = 70$ GeV and $M_{\mathrm{SUSY}} = 1200$ GeV. At the bottom of the table the $F_i$ factor, which is used for the QCD estimation, is also provided.

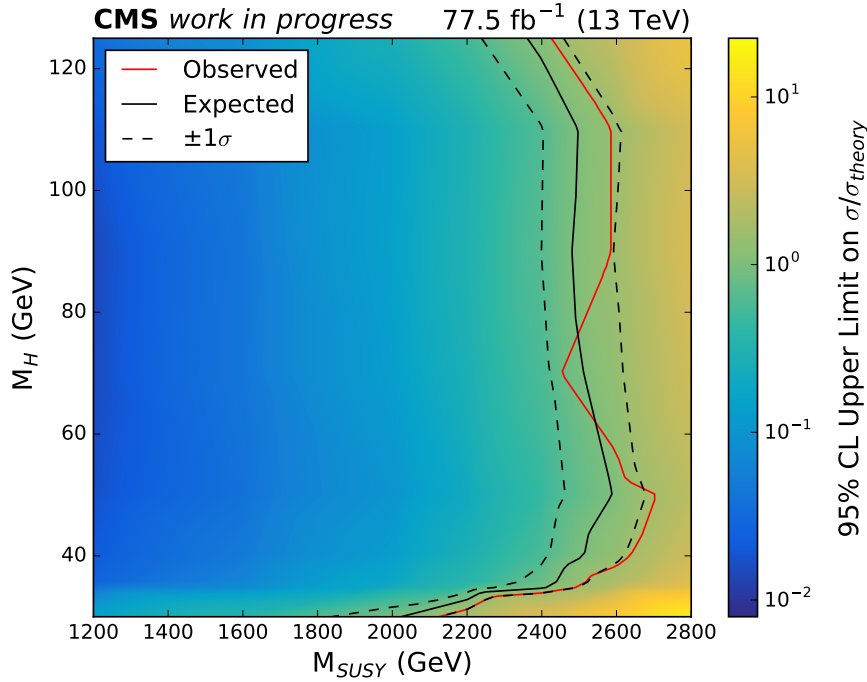| Mass type | $S$ mass region | | | |
|---|---|---|---|---|
| Observed yield | 235 | | | |
| MC process | signal | TTJets | ZJets | WJets |
| Yield | 626.67 | 12.67 | 4.77 | 6.89 |
| Systematic | | | | |
| isrWeight201$_7^6$ | 0.986/1.014 | - | - | - |
| luminosity2017 | 1.023 | 1.023 | 1.023 | 1.023 |
| XS_TTJets2017 | - | 1.500 | - | - |
| XS_ZJets2017 | - | - | 1.500 | - |
| XS_WJets2017 | - | - | - | 1.500 |
| jecUnc201$_7^6$ | 1.019/0.979 | 0.941/1.076 | 1.000/1.056 | 1.000/1.025 |
| jerUnc2017 | 1.007/0.992 | 1.016/1.008 | 1.035/1.000 | 1.000/0.975 |
| jmsUnc201$_7^6$ | 0.948/1.049 | 1.019/0.950 | 1.092/0.958 | 1.050/0.975 |
| jmrUnc201$_7^6$ | 1.000/0.874 | 1.000/1.075 | 1.000/1.151 | 1.000/0.847 |
| SigDbtTag2017 | 0.912/1.169 | - | - | - |
| TtDbtTag2017 | - | 0.883/1.127 | - | - |
| prefire2017 | 1.008/0.992 | 1.010/0.990 | 1.004/0.996 | 1.012/0.988 |
| Mass type | $U + D$ mass region | | | |
| Observed yield | 254 | | | |
| MC process | signal | TTJets | ZJets | WJets |
| Yield | 10.38 | 13.66 | 9.61 | 5.32 |
| Systematic | | | | |
| isrWeight201$_7^6$ | 0.995/1.005 | - | - | - |
| luminosity2017 | 1.023 | 1.023 | 1.023 | 1.023 |
| XS_TTJets2017 | - | 1.500 | - | - |
| XS_ZJets2017 | - | - | 1.500 | - |
| XS_WJets2017 | - | - | - | 1.500 |
| jecUnc201$_7^6$ | 1.008/0.977 | 0.929/1.049 | 0.970/1.076 | 0.908/1.000 |
| jerUnc2017 | 1.024/0.992 | 0.978/1.006 | 1.024/1.046 | 1.000/1.000 |
| jmsUnc201$_7^6$ | 1.013/1.071 | 0.974/0.964 | 0.976/1.041 | 1.000/0.991 |
| jmrUnc201$_7^6$ | 1.000/2.345 | 1.000/0.946 | 1.000/0.834 | 1.000/0.978 |
| SigDbtTag2017 | 0.914/1.151 | - | - | - |
| TtDbtTag2017 | - | 0.880/1.130 | - | - |
| prefire2017 | 1.008/0.992 | 1.010/0.990 | 1.011/0.989 | 1.010/0.990 |
| $F_i$ factor | $0.954 \pm 0.150$ | | | |

Figure 7.22: Observed and expected signal model exclusions, as a function of $M_H$ and $M_{SUSY}$, for 2016 and 2017 data. This is for the scenario where the gluino mass is 1% higher than that of the squarks. The red line indicates the observed excluded region. The solid black line indicates the median expected excluded region. The dashed black lines indicate the expected excluded regions with $\pm 1\sigma$ in experimental uncertainty. The colour scale indicates the median expected 95% confidence level upper limit of $\sigma/\sigma_{\text{theory}}$.

As there were no excesses in the data, the observed upper limits of $\sigma/\sigma_{\text{theory}}$ were calculated for all the different signal models. These were then used to determine the observed exclusion contour in the $M_H$-$M_{SUSY}$ plane, which can be seen in Figure 7.22. Across all the $M_H$ parameter space, the observed $M_{SUSY}$ limits agree well with the expected limits and, other than at $M_H \approx 50$ GeV, they are contained within the expected $\pm 1\sigma$ bands.

The shape of the observed exclusion contours, relative to the median expected contours, can be understood by considering the statistical fluctuations of the data yields in the $H_T \in 3500+$ GeV bin. In Figure 7.21, the observed 2016 yield is zero in search regions 21 and 22, and the observed 2017 yield is zero in search region 22. This deficit of events is what causes the observed $M_{SUSY}$ limit to be greater than expected in the range $30 < M_H < 50$ GeV. The slight excess of the observed 2016 yield in search
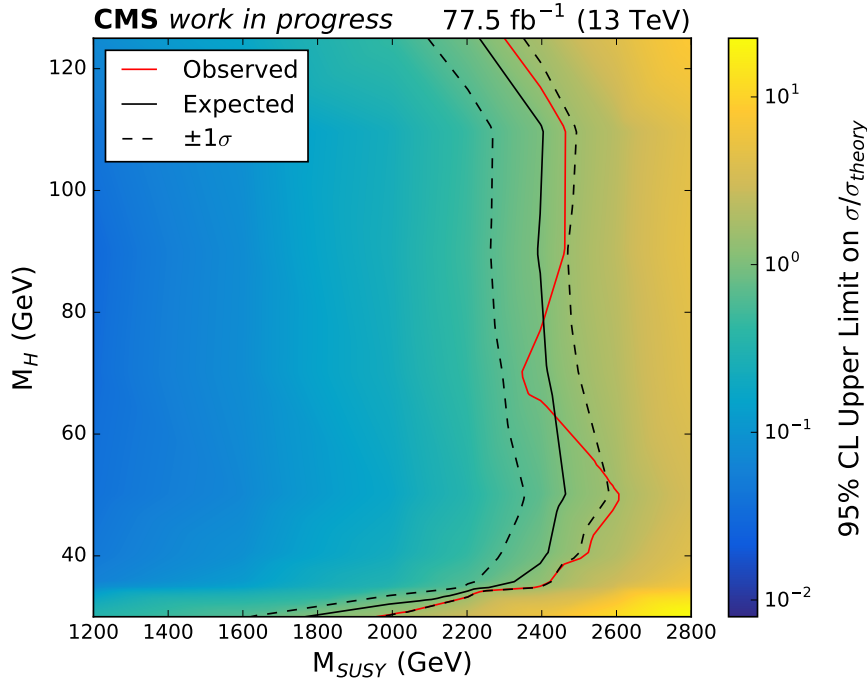
Figure 7.23: Observed and expected signal model exclusions, as a function of $M_\mathrm{H}$ and $M_\mathrm{SUSY}$, for 2016 and 2017 data. This is for the scenario where only the squarks are kinematically accessible. The red line indicates the observed excluded region. The solid black line indicates the median expected excluded region. The dashed black lines indicate the expected excluded regions with $\pm 1\sigma$ in experimental uncertainty. The colour scale indicates the median expected 95% confidence level upper limit of $\sigma/\sigma_\mathrm{theory}$.

region 23, and of the observed 2017 yield in search regions 23 and 24, lead to the $M_\mathrm{SUSY}$ limit being less than expected for $M_\mathrm{H} \approx 70$ GeV. Finally, the higher than expected $M_\mathrm{SUSY}$ limit in the range $90 < M_\mathrm{H} < 125$ GeV is driven by the observed 2016 yields of zero in search regions 26-30, and the observed 2017 yields of zero in search regions 26 and 28.

Figure 7.23 shows the signal model exclusion for the scenario where the gluino mass is much larger than that of the squarks (described in Section 5.1.1). The shapes of the exclusion contours are very similar to those in Figure 7.22, but shifted to around 100 GeV less in the $M_\mathrm{SUSY}$ dimension. This is because the production modes through gluino states are no longer kinematically accessible, resulting in lower production cross sections and, thus, less signal sensitivity.

In order to demonstrate the validity of the statistical methods employed, the observed limits were recalculated using data that was artificially injected with signal. This study is presented in Appendix G.

## 7.2.5 Comparison to other SUSY Searches

A recent CMS analysis [77] has set limits on light-flavour squark pair production, for the scenario where both squarks decay as $\tilde{q}_i \rightarrow q_i + \tilde{\chi}_1^0$, using the entire Run-II data set (137 fb$^{-1}$). For low LSP masses, the squarks are excluded for masses less than 1800 GeV, at 95% confidence level. These results can be loosely compared to the $M_{\mathrm{SUSY}}$ limits obtained in this search, as the signal model has the same production mechanism when the gluino mass is decoupled. In this search, the $M_{\mathrm{SUSY}}$ limit is greater than 2400 GeV across a broad range of $M_{\mathrm{H}}$ parameter space (see Figure 7.23). Despite using less data, significantly higher squark mass exclusions are achieved due to the attempted reconstruction of the two light Higgs bosons.

Due to the unique signal model, the results from this analysis cannot be directly compared to any other analyses at CMS or ATLAS. The most similar analysis [78] was conducted by CMS. It searched for SM Higgs bosons in SUSY decay cascades following the production of gluino pairs. The gluinos decayed as $\tilde{g} \rightarrow q\bar{q} + \tilde{\chi}_2^0$ and the NLSP decayed as $\tilde{\chi}_2^0 \rightarrow H_{\mathrm{SM}} + \tilde{\chi}_1^0$. The LSP mass was fixed at 1 GeV and the NLSP mass was set to be 50 GeV less than the gluino mass. Consequently, the final state had two high energy SM Higgs bosons and two high energy LSPs. This meant, in addition to the reconstruction of the boosted SM Higgs bosons, the analysis could identify its signal events from the considerable amount of $E_{\mathrm{T}}^{\mathrm{miss}}$ arising from the LSPs. Using the 2016 data set (36 fb$^{-1}$), an upper limit of 0.001 pb was set on the $\tilde{g}\tilde{g}$ production cross section, at 95% confidence level. The result, in the context of this analysis, corresponds to a $M_{\mathrm{SUSY}}$ limit of 2400 GeV when the gluino mass is decoupled (see Table 5.6). Despite using less data, this limit is slightly higher than the $M_{\mathrm{SUSY}}$ limit obtained in this analysis for $M_{\mathrm{H}}$=125 GeV (see Figure 7.23). This is because the other analysis could use $E_{\mathrm{T}}^{\mathrm{miss}}$ in its event selection and was not required to reconstruct Higgs bosons with variable masses.

# 8

# Conclusion

In this thesis, an entirely new Beyond the Standard Model physics search, conducted using the CMS experiment, has been presented. The analysis has determined that the NMSSM decay cascades described in Section 5.1.1, configured so that the LSPs receive very little momentum, do not exist for squark and gluino masses less than around 2500 GeV.

The analysis is built around the reconstruction of the two light Higgs bosons produced in the decay cascades. Due to the suppression of the LSP momenta, the Higgs bosons have very high energies. They are reconstructed in the $b\bar{b}$ decay mode, which dominates at low Higgs boson masses. Because the Higgs bosons are boosted objects, the angular separation between their daughter b-quarks is small. Consequently, each $b\bar{b}$ pair is reconstructed in a single fat-jet.

Before the analysis was created, it was not clear whether searching for a rare, all hadronic final state would result in being inundated by the QCD multi-jet background. However, using jet grooming algorithms, a novel double-b-tagger, and aggressive $H_\mathrm{T}$ binning, the QCD background is highly suppressed whilst a significant fraction of signal events are retained. The event selection focusses on the 2D spaces formed from the variables of the two fat-jets representing the Higgs bosons. These

2D spaces allow for a data driven QCD estimation that can operate, independently, on a set of search regions spanning a broad range of Higgs boson masses.

There are a few ways in which future iterations of this analysis could gain additional signal sensitivity. A new double-b-tag discriminator has recently been developed by CMS, which is trained using a deep neural network. The initial studies suggest that it achieves greater signal acceptance for an equivalent QCD mistag rate. This, of course, would need to be validated in the context of this analysis. Additionally, the QCD estimation method would need to be verified again. Greater signal sensitivity could also be achieved following the development of a new jet grooming algorithm. Currently, the soft-drop algorithm is used. It allows for strong QCD suppression, however, a sizeable fraction of the signal model fat-jets, which reconstruct the decay products of both b-hadrons, also have a soft-drop mass of around 0 GeV. Another way to increase the signal sensitivity is to simply use more data in the analysis. Currently around 60 fb$^{-1}$ of data, collected by CMS during the 2018 run, is yet to be analysed.

It is not clear, however, whether this analysis should be repeated in the future. There is a complimentary way to identify (or rule out) the targeted NMSSM configuration; to search for the direct production of the light CP-even neutral Higgs boson. This approach is not dependent on producing squarks with arbitrarily high masses (low cross sections) and, having set lower bounds on the squark mass of around 2500 GeV, might now be the better search strategy. It should be noted that despite its lower mass, the light Higgs boson is more difficult to discover than the SM Higgs boson. This is because, when the light Higgs boson is primarily composed of the gauge singlet, it has significantly smaller couplings than the SM Higgs boson [14, 57]. Consequently, the SM Higgs boson searches at LEP [79] did not exclude the light Higgs boson. If the light Higgs boson did exist, however, one might expect some tension to arise in the 'searches' for the different Higgs boson decay modes conducted at CMS and ATLAS. In accordance, the targeted NMSSM configuration should be rigorously checked against the most recent CMS and ATLAS publications

(e.g. Ref. [80, 81] for the H $\to$ b$\bar{\text{b}}$ decay mode). If this results in the exclusion of the light Higgs boson (or demonstrates a strong potential to do so), a future iteration of the analysis could instead focus solely on the scenario where the Higgs bosons in the final state are SM Higgs bosons, as described in Ref. [82].

## 8.1 Closing Words

In one sense, the analysis presented in this thesis is like every other Beyond the Standard Model physics search; it did not reveal any new fundamental particles or interactions. At CMS and ATLAS, a myriad of supersymmetry hypotheses have been tested and, to date, none have been proven correct. These results do not disprove supersymmetry. Due to its rich parameter space, one can never rule out supersymmetry definitively. It is, however, becoming less convincing. The supersymmetric solution to the fine-tuning problem is being undermined by the lower bounds on the squark and gluino masses which, under many different scenarios, are currently well within the TeV domain.

It may be that Beyond the Standard Model physics will always remain a mystery to us and that supersymmetric theories will never be more than brilliant pieces of mathematics.

# A

# LSP Momentum Calculations

The relativistic kinematic equations used in this appendix can be found in many relativity textbooks, such as Ref. [83]. In the rest frame of the NLSP, the decay $\widetilde{\chi}_2^0 \to H + \widetilde{\chi}_1^0$ yields an LSP with energy:

$$E_{\widetilde{\chi}_1^0}^{\mathrm{rest}} = \frac{M_{\widetilde{\chi}_2^0}^2 + M_{\widetilde{\chi}_1^0}^2 - M_{\mathrm{H}}^2}{2M_{\widetilde{\chi}_2^0}} \tag{A.1}$$

and momentum:

$$|\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\,\mathrm{rest}}|^2 = (E_{\widetilde{\chi}_1^0}^{\mathrm{rest}})^2 - M_{\widetilde{\chi}_1^0}^2 \tag{A.2}$$

Boosting into the lab frame, the LSP has a momentum within the following range:

$$|\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\,\mathrm{lab}}| = \left| \beta\gamma E_{\widetilde{\chi}_1^0}^{\mathrm{rest}} \pm \gamma |\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\,\mathrm{rest}}| \right| \tag{A.3}$$

where the variation arises from the alignment of $\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\,\mathrm{rest}}$ with the boost direction. Consequently, the $|\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\,\mathrm{lab}}|$ distribution is approximately sinusoidal in shape with the central value being the most probable and the extrema having a probability of zero (assuming isotropic decays). The values of $\beta$ and $\gamma$ correspond to the NLSP in the

lab frame, and can be written in terms of the NLSP momentum and mass, such that Equation A.3 becomes:

$$|\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\text{lab}}| = \left| \frac{|\overrightarrow{p}_{\widetilde{\chi}_2^0}^{\text{lab}}|}{M_{\widetilde{\chi}_2^0}} E_{\widetilde{\chi}_1^0}^{\text{rest}} \pm \frac{(|\overrightarrow{p}_{\widetilde{\chi}_2^0}^{\text{lab}}|^2 + M_{\widetilde{\chi}_2^0}^2)^{\frac{1}{2}}}{M_{\widetilde{\chi}_2^0}} |\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\text{rest}}| \right| \tag{A.4}$$

Equations A.2 and then A.1 can be substituted into Equation A.4 to evaluate $|\overrightarrow{p}_{\widetilde{\chi}_1^0}^{\text{lab}}|$ in terms of the $M_{\widetilde{\chi}_1^0}$, $M_{\widetilde{\chi}_2^0}$, $M_{\text{H}}$, and $|\overrightarrow{p}_{\widetilde{\chi}_2^0}^{\text{lab}}|$ (which depends on $M_{\text{SUSY}}$). For a given value of $M_{\text{H}}$ and $|\overrightarrow{p}_{\widetilde{\chi}_2^0}^{\text{lab}}|$, this allows one to see how the LSP momentum varies with the signal model parameters $\Delta_M \equiv M_{\widetilde{\chi}_2^0} - M_{\text{H}} - M_{\widetilde{\chi}_1^0}$ and $R_M \equiv M_{\text{H}} / M_{\widetilde{\chi}_2^0}$.

Interpreting the resultant equation is not trivial and, therefore, a graphical solution is provided in Figure A.1. It is configured with $M_{\text{H}} = 70$ GeV and $|\overrightarrow{p}_{\widetilde{\chi}_2^0}^{\text{lab}}| = 1000$ GeV (corresponding to an initial squark, at rest, with mass 2000 GeV). Using the graphical solution, it can be seen how $R_M$ primarily controls the average momentum transferred to the LSP. The spread of the LSP momenta is controlled by $\Delta_M$, however, the size of this parameter is constrained by the value of $R_M$. As $R_M$ tends towards unity, the LSP momentum tends towards zero, which is the key feature of the signal model.
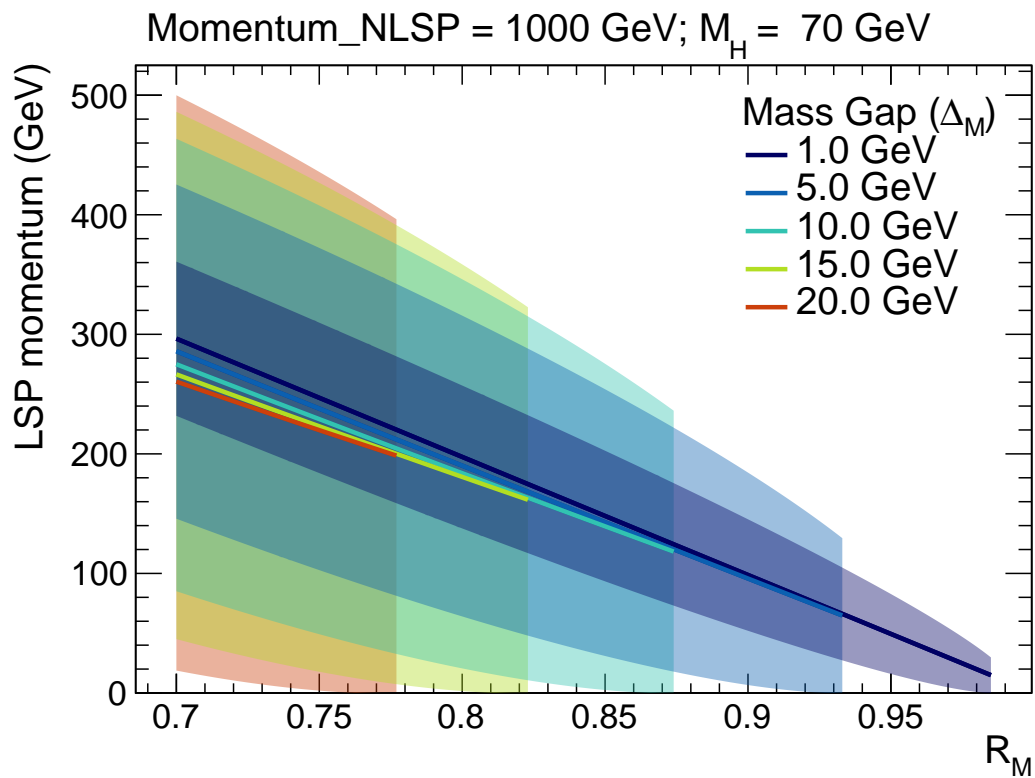
Figure A.1: Graphical representation of the LSP momentum as a function of $R_M$ in different $\Delta_M$ categories (see legend). The solid lines represent the distribution modes and the shaded regions represent the distribution spread. The NLSP momentum is 1000 GeV and $M_{\mathrm{H}} = 70$ GeV. Note that as $R_M \to 1$, the larger $\Delta_M$ configurations cease to exist.

# B

# Signal Model Dependence on the $R_M$ and $\Delta_M$ Parameters

The signal samples generated for the main analysis have the parameters $R_M = 0.99$ and $\Delta_M = 0.1$ GeV. In this appendix, the effect of relaxing $R_M$ away from unity, and then being able to increase the value of $\Delta_M$, is explored. All the signal samples used in this appendix have the parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 1200$ GeV. Instead, it is the $R_M$ and $\Delta_M$ parameters that are varied. These signal samples are produced in a different way to those used in the main analysis. The simulation of the initial state squarks and gluinos are still generated at leading order using MAD-GRAPH5_aMC@NLO 2.3.3, however, there is no longer an associated additional parton. Another difference is that PYTHIA version 6, rather than version 8, is used. The final difference is that the detector response is simulated in DELPHES, rather than using the full CMS simulation, because it is much quicker. These differences are acceptable because this appendix is only used to illustrate properties of the signal model.

As was explained in Section 5.1, the principle motivation behind searching for this signal model is that it can have low $E_{\mathrm{T}}^{\mathrm{miss}}$. This is achieved when the parameter
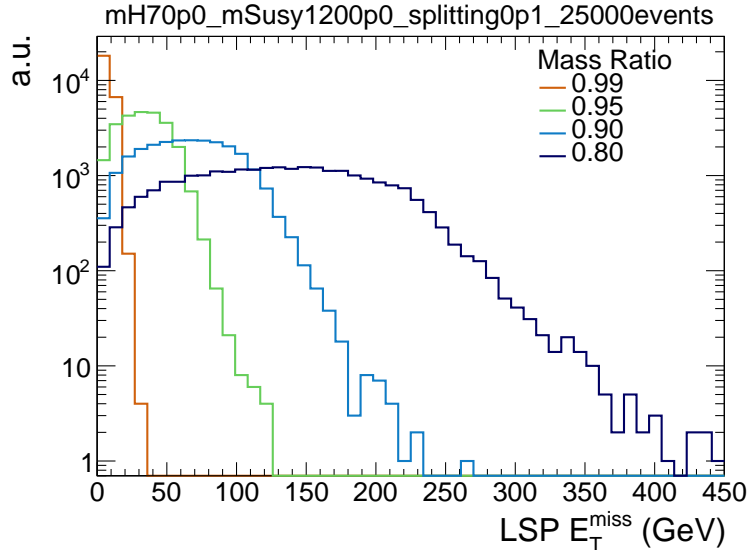
Figure B.1: Distribution of the magnitude of the sum of the two LSP transverse momenta vectors (i.e. the $E_T^{\mathrm{miss}}$ due to the LSPs). This is done for a variety of $R_M$ values, as indicated by the legend. The other parameters of the signal samples are $M_H = 70$ GeV, $M_{\mathrm{SUSY}} = 1200$ GeV, and $\Delta_M = 0.1$ GeV.

$R_M$ tends towards unity, as it suppresses the fraction of the energy that the LSP receives from the NLSP in the $\widetilde{\chi}_2^0 \rightarrow H + \widetilde{\chi}_1^0$ decays. Figure B.1 demonstrates this point. As $R_M$ decreases from unity, the $E_T^{\mathrm{miss}}$ due to the LSPs increases rapidly. This is because there is a lot of energy flowing through the decay cascade arms and only a small fraction of this energy has to leak into the LSP sector to create a considerable amount of $E_T^{\mathrm{miss}}$. The signal sample used has $M_{\mathrm{SUSY}} = 1200$ GeV, which is a relatively low mass scale in this analysis. For the signal samples with greater $M_{\mathrm{SUSY}}$ values, there will be more energy flowing through the decay cascade arms and thus, to constrain $E_T^{\mathrm{miss}}$ to the same degree, the $R_M$ parameter will have to be closer to unity.

The other degree of freedom, $\Delta_M$, controls the amount of phase space available to the products of the $\widetilde{\chi}_2^0 \rightarrow H + \widetilde{\chi}_1^0$ decays. As $R_M$ tends towards unity, the $\Delta_M$ parameter is constrained to be small. Figure B.2 shows how the $E_T^{\mathrm{miss}}$ distribution, due to the LSPs, broadens as the $\Delta_M$ parameter increases. The $E_T^{\mathrm{miss}}$ scale does not change significantly like it does when the $R_M$ parameter is varied. Consequently, the $\Delta_M$ parameter only plays a secondary role in controlling the signal model $E_T^{\mathrm{miss}}$ distribution.
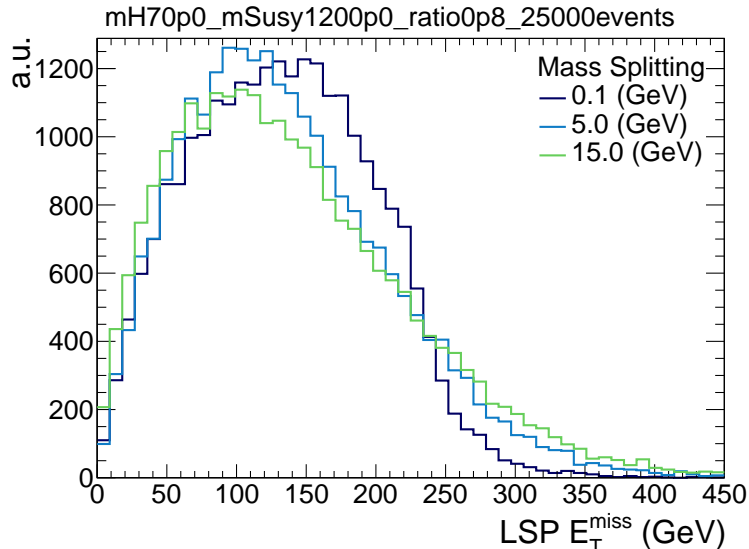
Figure B.2: Distribution of the magnitude of the sum of the two LSP transverse momenta vectors (i.e. the $E_T^{miss}$ due to the LSPs). This is done for a variety of $\Delta_M$ values, as indicated by the legend. The other parameters of the signal samples are $M_H = 70$ GeV, $M_{SUSY} = 1200$ GeV, and $R_M = 0.80$.

The signal model ceases to be interesting when the $E_T^{miss}$ due to the LSPs becomes too large. This is because the conventional jet$+E_T^{miss}$ SUSY searches would become sensitive to it. The $E_T^{miss}$ thresholds in these searches are typically around 200 GeV. Consequently, in the example provided in Figure B.1, the signal model stops being of primary interest as the $R_M$ parameter reaches 0.80. At this point, around 20% of the signal events have $E_T^{miss}$, due to the LSPs, above the 200 GeV threshold. This is defined as the transition point. For signal samples with greater values of $M_{SUSY}$, the transition point will be attained for values of $R_M$ closer to unity.

An important part of this analysis is reconstructing the two Higgs bosons in the two selected AK8 jets. As the signal model is brought from the $E_T^{miss}$ suppressed state to the transition point, the energies of these Higgs bosons will decrease. This is because they will receive a smaller fraction of the energy in the $\widetilde{\chi}_2^0 \rightarrow H + \widetilde{\chi}_1^0$ decays. The remainder of this appendix considers the impact on the analysis as the signal model is brought towards the transition point.

If the Higgs bosons have lower $p_T$, then so too will the two selected AK8 jets. In the kinematic event selection, these AK8 jets are required to have $p_T > 300$ GeV.
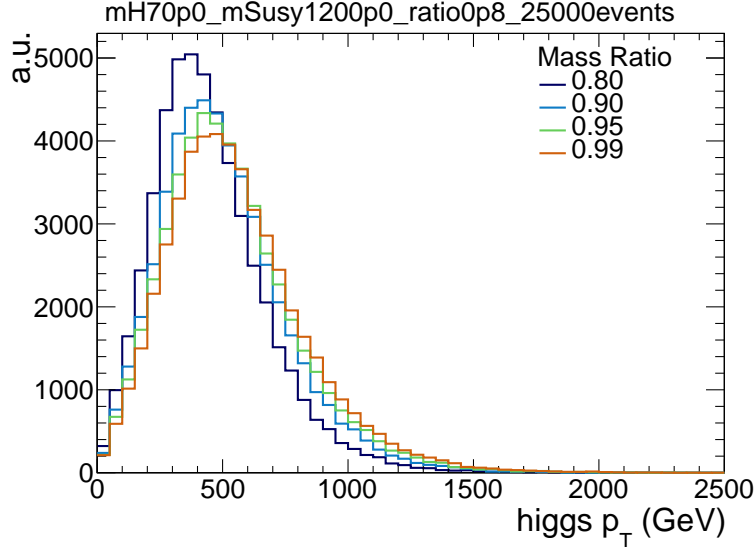
Figure B.3: Distribution of the Higgs boson $p_T$ in the decay cascades. This is done for a variety of $R_M$ values, as indicated by the legend. The other parameters of the signal samples are $M_H = 70$ GeV, $M_{SUSY} = 1200$ GeV, and $\Delta_M = 0.1$ GeV.

Figure B.3 shows how the $p_T$ distribution of the Higgs bosons changes as the $R_M$ parameter is relaxed from unity. At the transition point, the average $p_T$ of the Higgs bosons is reduced by around 100 GeV. This is an intuitive result, as the transition point corresponds to where the $E_T^{miss}$ can reach over 200 GeV. For signal models with low $M_{SUSY}$ values, this reduction in the Higgs boson $p_T$ will result in a noticeable loss in event selection efficiency. However, in the $M_{SUSY} = 1200$ GeV case, the majority of the Higgs bosons still have $p_T > 300$ GeV and, due to the large production cross sections, there will still be large excesses in the final yields. For signal models with larger $M_{SUSY}$ values, the loss in efficiency will be much smaller because the Higgs boson $p_T$ distributions are centred at values significantly larger than the 300 GeV threshold.

Another consequence of reducing the $p_T$ of the Higgs bosons, is that it will lead to an increase in the angular separation of the $b\bar{b}$ pairs. If the change was significant, and the b-quarks became too widely separated, it could become inefficient to try and reconstruct them both in a single AK8 jet. Figure B.4 shows how the $\Delta R(b\bar{b})$ distribution changes as the $R_M$ parameter is relaxed away from unity. The change in the $\Delta R(b\bar{b})$ distribution shape, as the signal model tends towards the transition
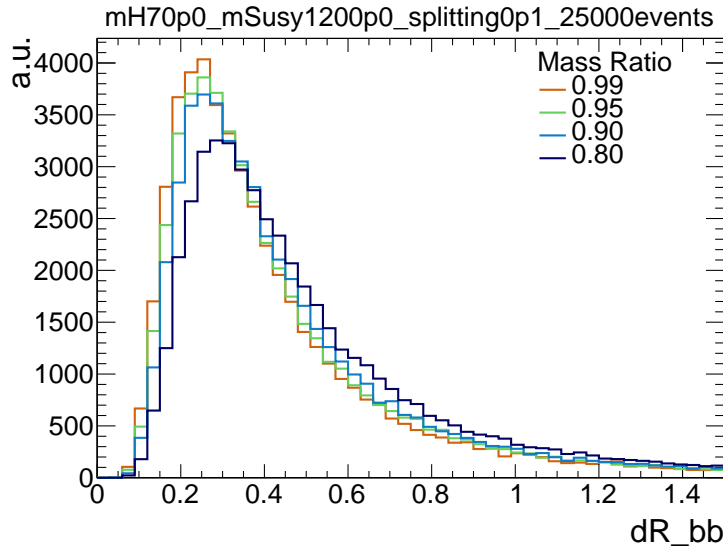
Figure B.4: Distribution of the $\Delta R(b\bar{b})$ separation resulting from the H $\rightarrow$ b$\bar{b}$ decays in the decay cascades. This is done for a variety of $R_M$ values, as indicated by the legend. The other parameters of the signal samples are $M_\mathrm{H} = 70$ GeV, $M_\mathrm{SUSY} = 1200$ GeV, and $\Delta_M = 0.1$ GeV.

point, is not large. The differences incurred by changing the $M_\mathrm{H}$ and $M_\mathrm{SUSY}$ parameters are much more considerable. For signal models with larger $M_\mathrm{SUSY}$ values, the fractional change in the Higgs boson $p_\mathrm{T}$ is smaller and, thus, the change in the $\Delta R(b\bar{b})$ distribution will be less than that in Figure B.4.

The final consequence of reducing the $p_\mathrm{T}$ of the Higgs bosons, is that this will reduce the $H_\mathrm{T}$ in the events. In the kinematic event selection, $H_\mathrm{T}$ is a binned quantity with a minimum value of 1500 GeV. Figure B.5 shows how the $H_\mathrm{T}$ distribution changes as the $R_M$ parameter is relaxed away from unity. At the transition point, there is an average reduction in $H_\mathrm{T}$ of around 200 GeV. The reduction in $H_\mathrm{T}$ will cause a slight shift in the shape of the final yields, with the populations in the lower $H_\mathrm{T}$ regions increasing. This will lead to a small loss in sensitivity because the background yields are higher in these regions. For signal models with low $M_\mathrm{SUSY}$ values, there will also be a reduction in the total efficiency because some events will drop below the $H_\mathrm{T}$ threshold of 1500 GeV. As was stated above, the low $M_\mathrm{SUSY}$ signal models can compensate for these losses due to their large production cross sections.
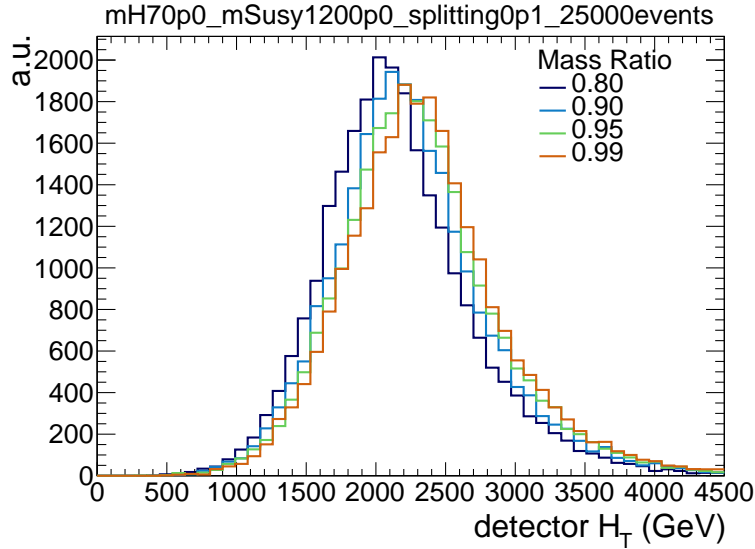
Figure B.5: The $H_T$ distribution for a variety of $R_M$ values, as indicated by the legend. The other parameters of the signal samples are $M_H = 70$ GeV, $M_{SUSY} = 1200$ GeV, and $\Delta_M = 0.1$ GeV.

In conclusion, the search strategy used in this analysis will maintain sensitivity to the signal model as the $E_T^{miss}$ suppression is relaxed. The introduction of $E_T^{miss}$ reduces the average $p_T$ of the Higgs bosons, which in turn reduces the event selection efficiency. The impact is greatest on signal models with low $M_{SUSY}$ values, however, these signal models have large production cross sections that can compensate for the losses. Significant signal sensitivity should still be present at the point where more conventional jet$+E_T^{miss}$ SUSY searches start to gain sensitivity.

# C

# Additional Information about the Signal Model

This appendix presents some additional properties of the signal model.

Figure C.1 compares the $p_{\mathrm{T}}$ distributions of the Higgs boson and the light-flavour quark from the same decay cascade arm. Figure C.2 compares the angular separation of the $b\bar{b}$ pair with the Higgs boson mass divided by its $p_{\mathrm{T}}$. This is in order to test the relationship $\Delta R_{b\bar{b}} \approx 2M_{\mathrm{H}}/(p_{\mathrm{T}})_H$. Figure C.3 shows the soft-drop mass distribution of the AK8 jet spatially matched ($\Delta R < 0.4$) with a Higgs boson. This is done for the cases where the corresponding $b\bar{b}$ pair have $\Delta R$ separation greater, and less than, 0.8. It shows that the primary reason why the selected AK8 jets, in the signal model, can have low soft-drop masses is due to the soft-drop grooming algorithm. This typically occurs when one of the b-quarks receives less than 10% of the Higgs boson $p_{\mathrm{T}}$, as can be seen in Figure C.4. It is demonstrated, in Figure C.5, how there are no search region gaps between the signal yields for samples with adjacent $M_{\mathrm{H}}$ values. Consequently, it is legitimate to use linear interpolation to create the exclusion contours in the $M_{\mathrm{H}}$-$M_{\mathrm{SUSY}}$ plane from the upper limits of $\sigma/\sigma_{\mathrm{theory}}$ calculated for all the different signal models.
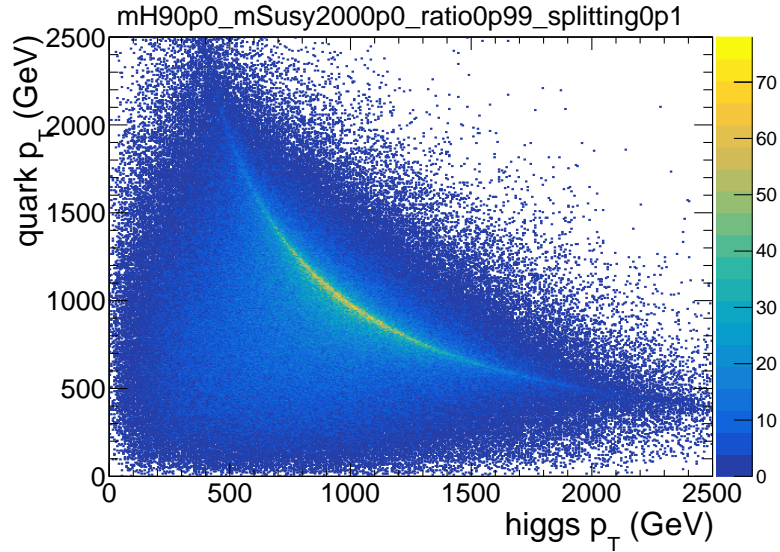
Figure C.1: Comparison of the $p_{\mathrm{T}}$ of the Higgs boson and the light-flavour quark from the same decay cascade arm. The signal sample has parameters $M_{\mathrm{H}} = 90$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV.
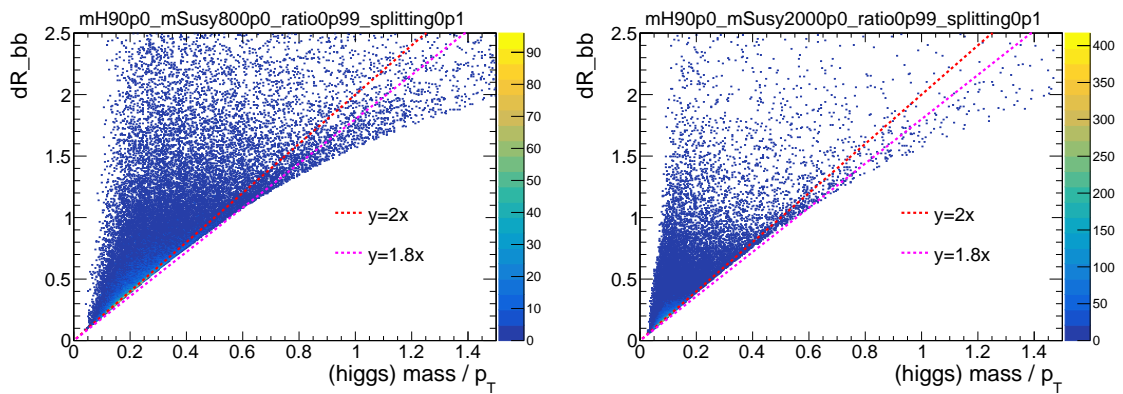


Figure C.2: Comparison of the angular separation of the $b\bar{b}$ pair with the Higgs boson mass divided by its $p_{\mathrm{T}}$. Left: $M_{\mathrm{H}} = 90$ GeV and $M_{\mathrm{SUSY}} = 800$ GeV. Right: $M_{\mathrm{H}} = 90$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV.
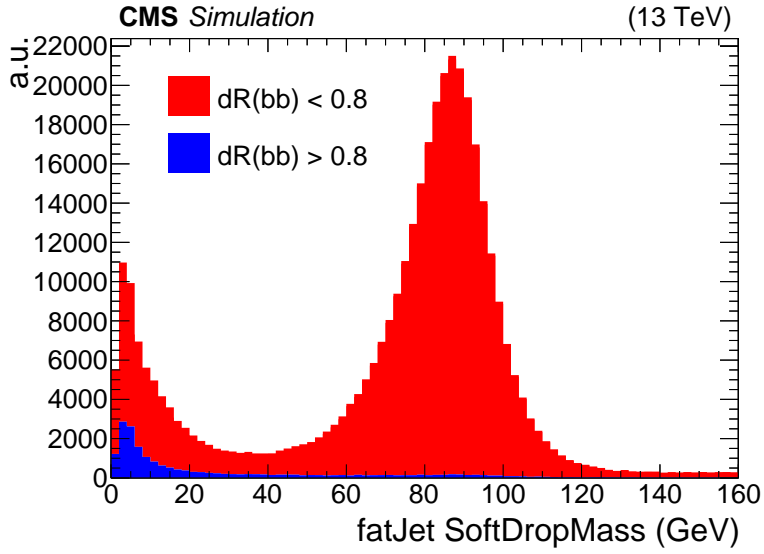
Figure C.3: Soft-drop mass distribution of the AK8 jet matched with a Higgs boson ($\Delta R < 0.4$). Red: The corresponding $b\bar{b}$ pair are separated by $\Delta R < 0.8$. Blue: The corresponding $b\bar{b}$ pair are separated by $\Delta R > 0.8$. The AK8 jet is required to have $p_{\mathrm{T}} > 170$ GeV. The signal sample has parameters $M_{\mathrm{H}} = 90$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV.
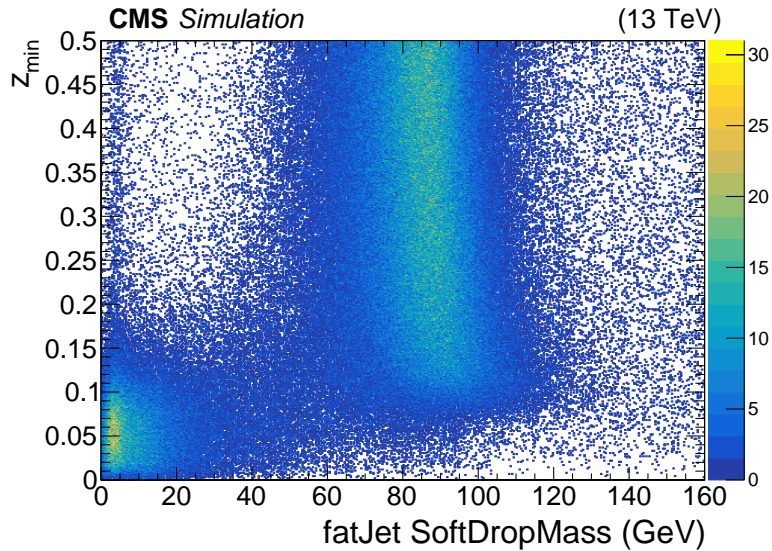


Figure C.4: Soft-drop mass of an AK8 jet matched with a Higgs boson ($\Delta R < 0.4$) plotted against $z_{\mathrm{min}}$, the fraction of the Higgs boson $p_{\mathrm{T}}$ received by the lowest $p_{\mathrm{T}}$ b-quark. Note that only $b\bar{b}$ pairs with angular separation $\Delta R < 0.8$ are used. The signal sample has parameters $M_{\mathrm{H}} = 90$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV.
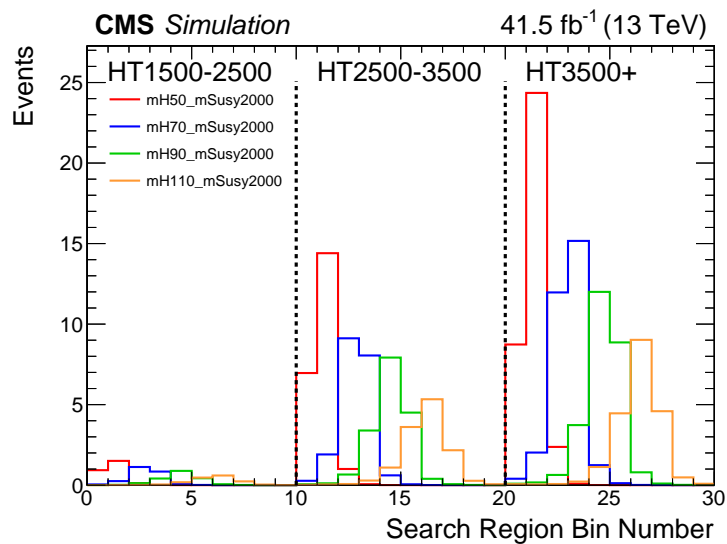
Figure C.5: Simulated yields for 2017 signal models parameterized by $M_{\mathrm{SUSY}} = 2000$ GeV and $M_{\mathrm{H}} = 50$, 70, 90, and 110 GeV. The tag double-b-tag region and the $S_n$ mass regions are used.

# D

---

# Definition of 2D Mass Regions

---

This appendix contains the information required to construct the thirty 2D mass regions used in the analysis. The 2D mass regions can be seen in Figure 6.6. All the construction lines are straight lines in the 2D soft-drop mass plane. The equations of the lines are listed below. The symbols $A$ and $B$ represent the soft-drop mass of fatJetA and fatJetB, respectively.

The two positive gradient lines for all $S_n$ mass regions:

$B = 1.360042 \cdot (A - 17.600000) + 40.000000$

$B = 0.735272 \cdot (A - 40.000000) + 17.600000$

The positive gradient line for the $U_n$ mass regions ($n > 1$) that does not correspond to the edge of the $S_n$ mass regions:

$B = 1.878173 \cdot (A - 6.400000) + 51.200000$

The positive gradient line for the $D_n$ mass regions ($n > 1$) that does not correspond to the edge of the $S_n$ mass regions:

$B = 0.532432 \cdot (A - 51.200000) + 6.400000$

The negative gradient line for the lower bound of the $S_1$ mass region:

$B = -1.000000 \cdot (A - 40.000000) + 17.600000)$

The negative gradient line for the lower bound of the $U_1$ mass region:

$B = -6.129856 \cdot (A - 17.600000) + 40.000000$

The negative gradient line for the lower bound of the $D_1$ mass region:

$B = -0.163136 \cdot (A - 40.000000) + 17.600000$

The negative gradient lines for the upper bounds of the $S_n$, $U_n$, and $D_n$ mass regions.

For $n = 1, 2, \dots 10$, respectively:

$B = -1.000000 \cdot (A - 51.900000) + 26.349732)$

$B = -1.000000 \cdot (A - 65.100000) + 36.055318)$

$B = -1.000000 \cdot (A - 78.300000) + 45.760903)$

$B = -1.000000 \cdot (A - 91.500000) + 55.466488)$

$B = -1.000000 \cdot (A - 104.700000) + 65.172073)$

$B = -1.000000 \cdot (A - 117.900000) + 74.877659)$

$B = -1.000000 \cdot (A - 131.100000) + 84.583244)$

$B = -1.000000 \cdot (A - 144.300000) + 94.288829)$

$B = -1.000000 \cdot (A - 157.500000) + 103.994415)$

$B = -1.000000 \cdot (A - 170.700000) + 113.700000)$

# E

---

# Additional Information for the

# QCD Estimation Method

---

This appendix contains the additional information that was referred to in the discussion of the data driven QCD estimation method in Section 6.3.

## E.1 Signal and Background Yields in Different Regions

The QCD estimation method predicts the QCD yield in the $S_n$ mass regions with tag double-b-tag. To do so, it uses the data yields from the different 2D mass and double-b-tag regions. In this part of the appendix, the signal and background yields in these spaces are examined using MC simulation. The 2017 MC samples are used because, for the background samples, the event weightings are much smaller and thus smoother distributions are attained.

In Figure E.1, the yields are shown in the 30 search regions with anti-tag double-b-tag and $S_n$ mass regions. The background is dominated by QCD events and the

signal processes are suppressed. The yields for the $U_n + D_n$ mass regions with anti-tag double-b-tag are shown in Figure E.2. Again, the background is dominated by QCD events. The signal processes are even further suppressed because of the use of the sideband mass regions. As was discussed in Section 6.3.4, the $F_i$ factors are calculated from the 1D soft-drop mass distributions in the anti-tag double-b-tag region. For completeness, the signal and background yields for these mass distributions, in all three $H_{\mathrm{T}}$ regions, are shown in Figure E.3. These distributions are also dominated by QCD.

In Figure E.4, the yields are shown in the 30 search regions with tag double-b-tag and $U_n + D_n$ mass regions. The $t\bar{t}$ process contributes to a sizeable fraction of the total background in the highest mass regions. In addition, there is a significant amount of signal contamination. These signal yields, however, are very small compared to the corresponding signal yields in the $S_n$ mass regions, as was shown in Figure 7.19. The fit used to obtain results accounts for the other background processes, plus any potential signal contamination, as was explained in Section 7.2.1.

In Section 6.3.4, the QCD method was implemented on data using the control double-b-tag region. It was claimed that this control region was a good space in which to test the methodology because it is dominated by QCD events and free of any significant signal contamination. This is validated in Figure E.5 and Figure E.6, which show the yields with control double-b-tag for the $S_n$ and $U_n+D_n$ mass regions, respectively.
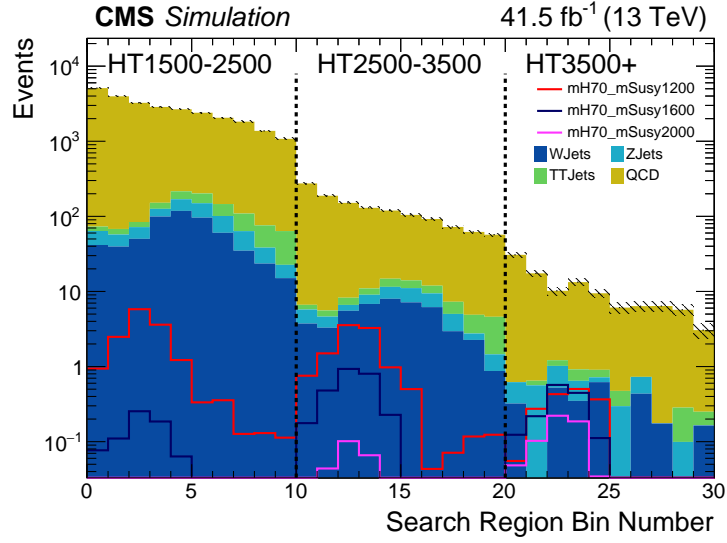
Figure E.1: Simulated yields for the 2017 signal and background. The signal samples are parameterized by $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 1200$, 1600, and 2000 GeV. The anti-tag double-b-tag region and the $S_n$ mass regions are used.
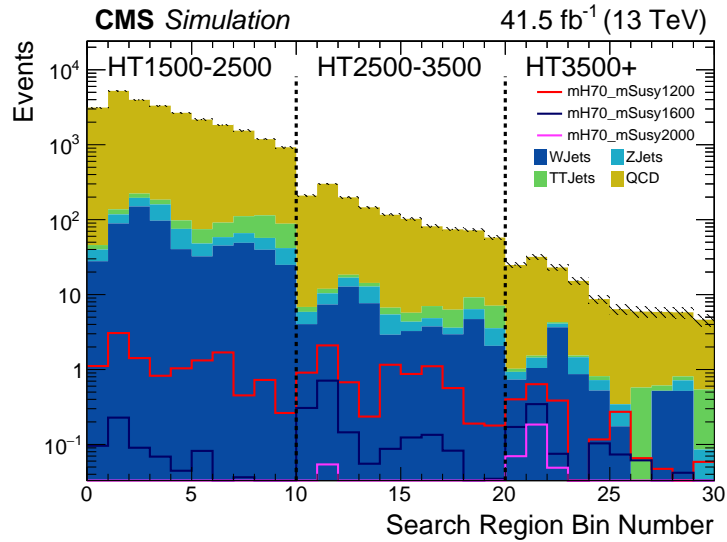


Figure E.2: Simulated yields for the 2017 signal and background. The signal samples are parameterized by $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 1200$, 1600, and 2000 GeV. The anti-tag double-b-tag region and the $U_n + D_n$ mass regions are used.
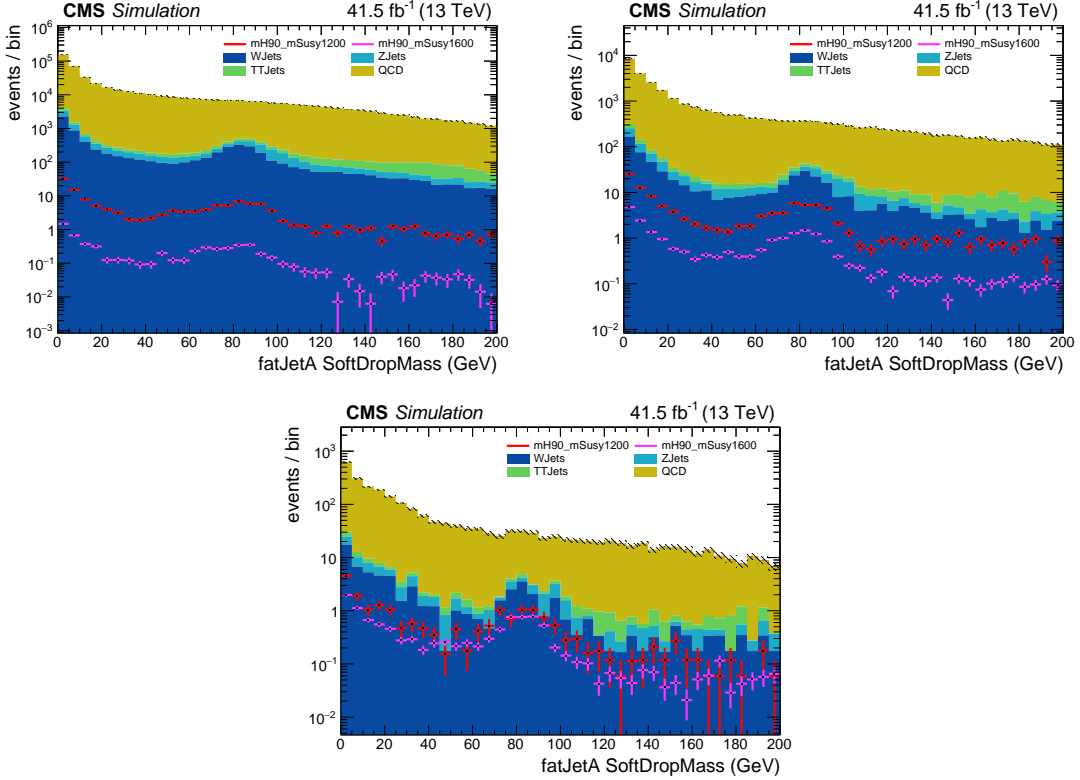
Figure E.3: Simulated soft-drop mass distribution of fatJetA for the 2017 signal and background. The signal samples are parameterized by $M_H = 90$ GeV and $M_{SUSY} = 1200$ and $1600$ GeV. All the events meet the anti-tag double-b-tag requirement and pass the kinematic cuts. Top left: $H_T \in 1500\text{-}2500$ GeV. Top right: $H_T \in 2500\text{-}3500$ GeV. Bottom: $H_T \in 3500+$ GeV.



Figure E.4: Simulated yields for the 2017 signal and background. The signal samples are parameterized by $M_H = 70$ GeV and $M_{SUSY} = 1200$, $1600$, and $2000$ GeV. The tag double-b-tag region and the $U_n + D_n$ mass regions are used.
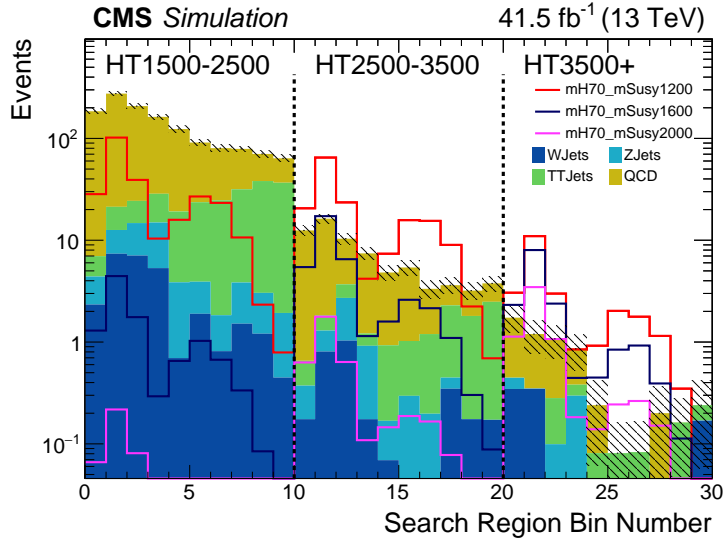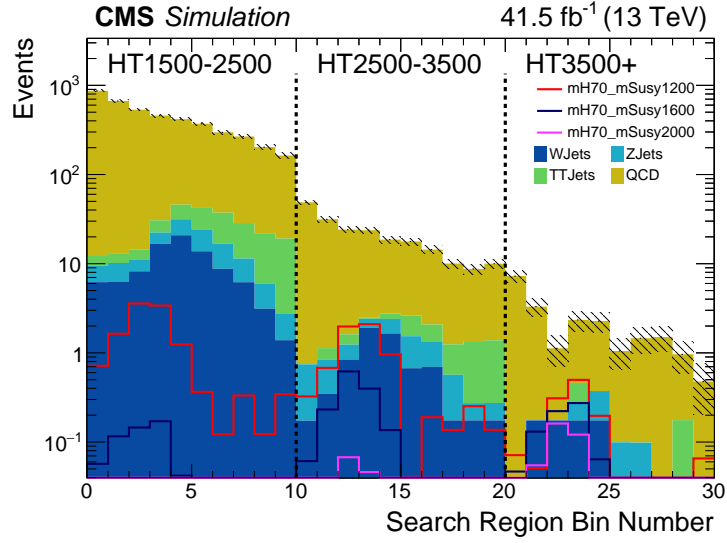
Figure E.5: Simulated yields for the 2017 signal and background. The signal samples are parameterized by $M_\text{H} = 70$ GeV and $M_\text{SUSY} = 1200$, $1600$, and $2000$ GeV. The control double-b-tag region and the $S_n$ mass regions are used.
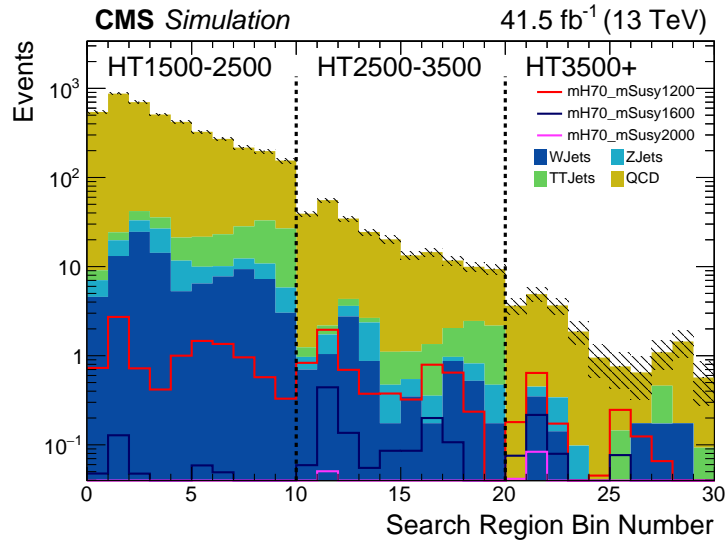


Figure E.6: Simulated yields for the 2017 signal and background. The signal samples are parameterized by $M_\text{H} = 70$ GeV and $M_\text{SUSY} = 1200$, $1600$, and $2000$ GeV. The control double-b-tag region and the $U_n + D_n$ mass regions are used.

# E.2   QCD Dominated Event Selection Space

In Section 6.3.2, the QCD dominated event selection space was defined. The space was designed to primarily select QCD events from data, whilst resembling the main kinematic event selection. It does so without placing any double-b-tag discriminator requirements on one of the selected AK8 jets, fatJetA. Using this AK8 jet, one can then study the relationship between the soft-drop mass and double-b-tag discriminator with real QCD events. In this part of the appendix, the tests that were applied to fatJetA in data are applied to the 2017 signal and background MC. This is done in order to validate that the distributions acquired from the data primarily consist of QCD events and have negligible signal contamination.

The first test looked at the distribution shape of the fatJetA double-b-tag discriminator in different fatJetA soft-drop mass bins. Using MC, this is repeated for the 45-50 GeV mass bin in Figure E.7. This is the highest mass bin used and, thus, the least favourable for the QCD yields. The signal samples used have $M_{\mathrm{H}} = 50$ GeV to replicate the worst case scenario, where the signal events are contained within the fatJetA mass bin. In this case, the $M_{\mathrm{SUSY}} = 1200$ GeV sample causes a 7% increase in the highest double-b-tag discriminator bin. If a $M_{\mathrm{SUSY}} = 1200$ GeV signal model did indeed exist, however, the excesses one would see in the search regions would be extremely large (see Figure 6.9). In this situation, one would not be concerned with the subtleties of validating the QCD estimation method. As for the background processes, the QCD events dominate across the whole double-b-tag discriminator spectrum.

The second test looked at the soft-drop mass distribution shape of fatJetA whilst requirements were placed on its double-b-tag discriminator value. The two cases considered were for double-b-tag discriminator values above and below 0.3. These distributions are recreated, using the MC samples, in Figure E.8. In the case where the double-b-tag discriminator is less than 0.3, the non-QCD backgrounds make only a small fraction of the total background and the signal contributions are negligible. The distribution shape of the non-QCD backgrounds roughly follows the QCD shape.
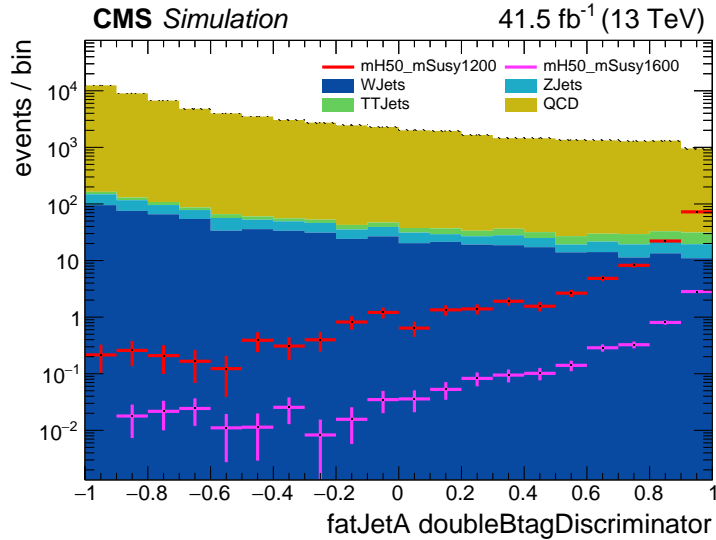
Figure E.7: Simulated double-b-tag discriminator distribution of fatJetA, where the soft-drop mass of fatJetA is between 45 and 50 GeV. This is done in the QCD dominated event selection space, for the 2017 signal and background. The signal samples are parameterized by $M_\mathrm{H} = 50$ GeV and $M_\mathrm{SUSY} = 1200$ and 1600 GeV.

It does, however, have a bump around 80-90 GeV, corresponding to events where the AK8 jet has reconstructed a W or Z boson. For $t\bar{t}$ events, this bump merges into a second, broader bump, which peaks at around 170 GeV. This second bump corresponds to events where the AK8 jet has reconstructed an entire top quark. In the case where the double-b-tag discriminator is greater than 0.3, the non-QCD backgrounds make up a larger fraction of the total background. This is especially true for the bumps mentioned above, because AK8 jets with substructure are more likely to have higher double-b-tag discriminator values. Despite these increases in the non-QCD backgrounds, they still do not have a significant impact on the total soft-drop mass distribution shape. Furthermore, the signal contributions still remain insignificant, despite being enhanced when the double-b-tag discriminator is greater than 0.3.
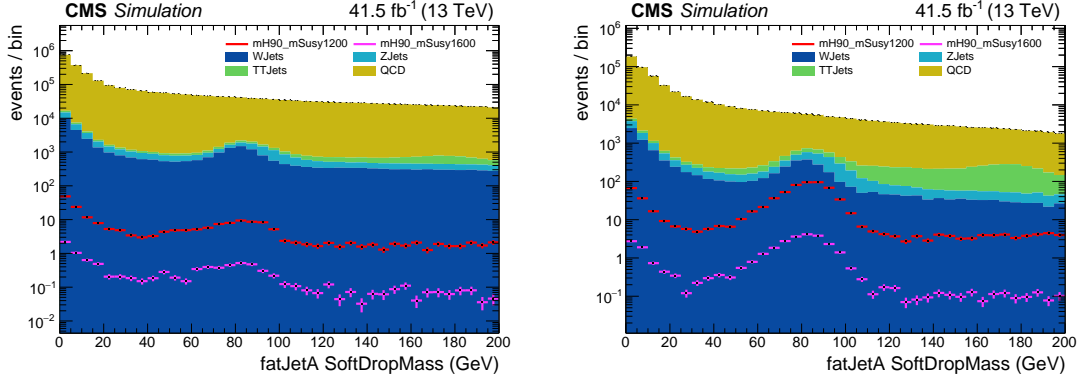
Figure E.8: Simulated soft-drop mass distribution of fatJetA, where the double-b-tag discriminator of fatJetA is less than 0.3 (left) and greater than 0.3 (right). This is done in the QCD dominated event selection space, for the 2017 signal and background. The signal samples are parameterized by $M_H = 90$ GeV and $M_{SUSY} = 1200$ and 1600 GeV.

## E.3 Independence of QCD AK8 Jets

In this part of the appendix, the independence of the two selected AK8 jets, in QCD events, is tested. The tests are performed on 2016 QCD MC. A minimal version of the kinematic event selection is used. The events have $H_T > 1500$ GeV and both selected AK8 jets are required to have $p_T > 300$ GeV.

The studies conducted on the QCD dominated event selection space assume that requiring one of the selected AK8 jets, fatJetB, to have a double-b-tag discriminator less than 0.3, does not effect the double-b-tag discriminator shape of the other AK8 jet, fatJetA. This assumption is tested in Figure E.9. There is a correlation between the double-b-tag discriminators of the two selected AK8 jets. When fatJetB has a double-b-tag discriminator less than 0.3, the double-b-tag discriminator of fatJetA is more likely to take a low value. This correlation, however, is very small and does not effect the study that was performed.

In Section 6.3.4, the $F_i$ factors were calculated from the 1D soft-drop mass distributions. It was claimed that this is possible because, in QCD events, the soft-drop mass distributions of the two selected AK8 jets only have a small dependence on each other. This is demonstrated in Figure E.10. The normalised soft-drop mass distributions of fatJetA are compared in the cases where there is, and is not, a soft-
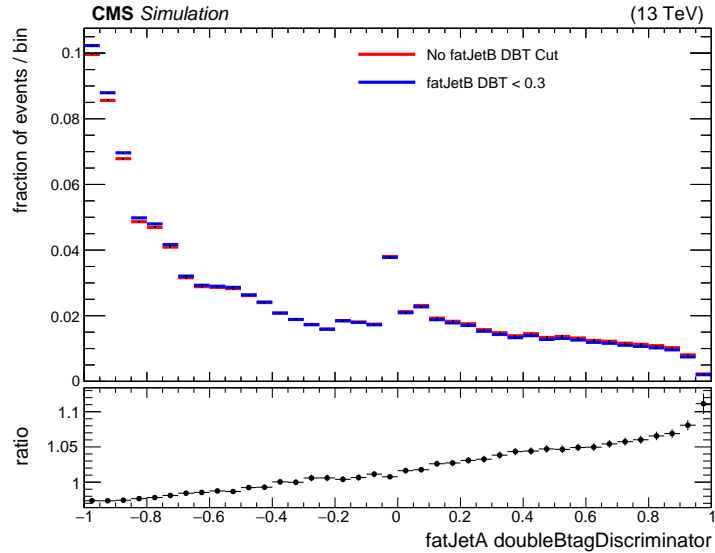
Figure E.9: Normalised fatJetA double-b-tag discriminator distribution for the case where there are no fatJetB requirements (red), and the case where the double-b-tag discriminator of fatJetB is required to be less than 0.3 (blue). This is performed on 2016 QCD MC using the events that pass pre-selection, have $H_\mathrm{T} > 1500$ GeV, and have both selected AK8 jets with $p_\mathrm{T} > 300$ GeV.

drop mass cut applied to fatJetB. The mass cuts applied to fatJetB do not have a significant effect on the fatJetA soft-drop mass distribution shape. The ratios of the normalised distributions are always within 5%, which is very small given that the yields, for soft-drop masses between 15 and 200 GeV (the masses spanned by the 2D mass regions), change by an order of magnitude. In the case of the W+jets and Z+jets samples, the soft-drop mass distributions of the two selected AK8 jets do have a dependence on each other. If one of the selected AK8 jets is revealed to have a soft-drop mass of around 80 GeV, it is likely that this jet is a reconstructed W or Z boson. Consequently, the other selected AK8 jet is now more likely to be from ISR and the probability of it having a mass of around 80 GeV is reduced.

Figure E.10: Normalised soft-drop mass distributions of fatJetA. The red points represent the distribution when there are no fatJetB mass requirements. The blue points represent the distribution when the fatJetB soft-drop mass is within a given mass window. Top left: 20-50 GeV. Top right: 70-110 GeV. Bottom: 120-200 GeV. This is performed on 2016 QCD MC using the events that pass pre-selection, have $H_{\mathrm{T}} > 1500$ GeV, and have both selected AK8 jets with $p_{\mathrm{T}} > 300$ GeV.

# E.4 Analysis Sensitivity to $F_i$ Uncertainties

To evaluate the impact that the $F_i$ uncertainties (see Section 6.3.4) have on the analysis, the ordinary expected limits in the $M_\mathrm{H}$-$M_\mathrm{SUSY}$ plane are compared to those where the $F_i$ factors, for both 2016 and 2017, are frozen to their central value. This comparison is shown in Figure E.11. The differences in the expected $M_\mathrm{SUSY}$ limits are small, especially when compared to the spread of the expected limits.



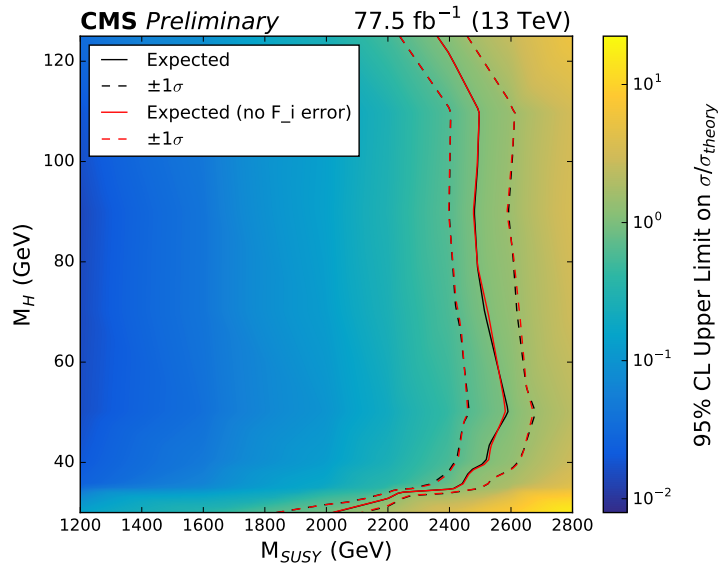Figure E.11: A comparison of the ordinary expected limits, using both 2016 and 2017 data, in the $M_\mathrm{H}$-$M_\mathrm{SUSY}$ plane (black lines) with the expected limits where the $F_i$ factors are frozen to their central values in the fit (red lines).

# F

---

# Test of 2016 $t\bar{t}$ MC

---

The $t\bar{t}$ process is the second largest background in this analysis. It is estimated using MC simulation. A test was devised in order to determine how accurate the $t\bar{t}$ MC is. The test compares a data driven estimate of the $t\bar{t}$ process with the 2016 $t\bar{t}$ MC. It does so in a new event selection space which is similar to that used in the main analysis. The test, and the results, are presented in this appendix.

The test requires a new event selection space that meets the following criteria:

- The new event selection space should be similar to that used in the main analysis.

- The new event selection space should have no potential signal contamination.

- The $t\bar{t}$ yields in the new event selection space should be significantly larger than the statistical fluctuations of the QCD yields.

- The $t\bar{t}$ yields in the new event selection space should be significantly larger than the W+jets and Z+jets contributions.

The criteria above are achieved by using a new set of 2D mass regions. These regions can be seen in Figure F.1. The new 2D mass regions use the same naming
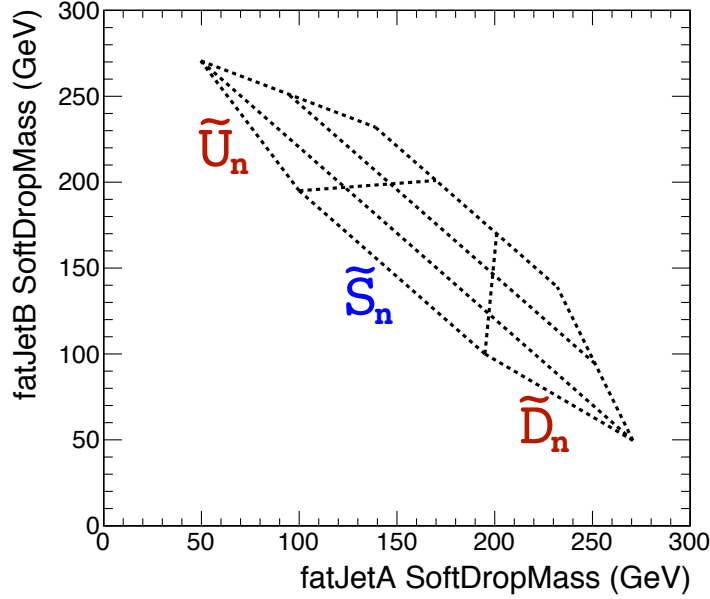
Figure F.1: The new mass regions used in the 2D soft-drop mass plane for the $t\bar{t}$ MC study.

convention as those used in the main analysis, but with a tilde symbol to mark the distinction. Figure F.2 shows the distribution of $t\bar{t}$ events in the 2D soft-drop mass plane, with the new 2D mass regions overlaid. The $\widetilde{S}_n$ mass regions are designed to select $t\bar{t}$ events where both of the selected AK8 jets have a soft-drop mass around that of the top quark mass. The large soft-drop masses prevent QCD events from dominating after the event selection, as can be seen in Figure F.3. In addition, there is no significant signal contamination because the new 2D mass regions exist beyond the mass regions used in the main analysis. The shape of the new 2D mass regions are such that the $\widetilde{S}_n$ mass regions have a high $t\bar{t}$ event density, whilst the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions do not. In order to increase statistics, the area of the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions are twice as large, relative to the corresponding $\widetilde{S}_n$ mass region, as they are in the original 2D mass regions. For $n = 1$, the area of the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions is roughly equal to the area of the $\widetilde{S}_n$ mass region. For $n > 1$, the area of the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions is equal to twice the area of the $\widetilde{S}_n$ mass region.

To further boost the statistics in the new event selection, the shape of the tag double-b-tag region is modified. In order to mark the distinction, the modified region is
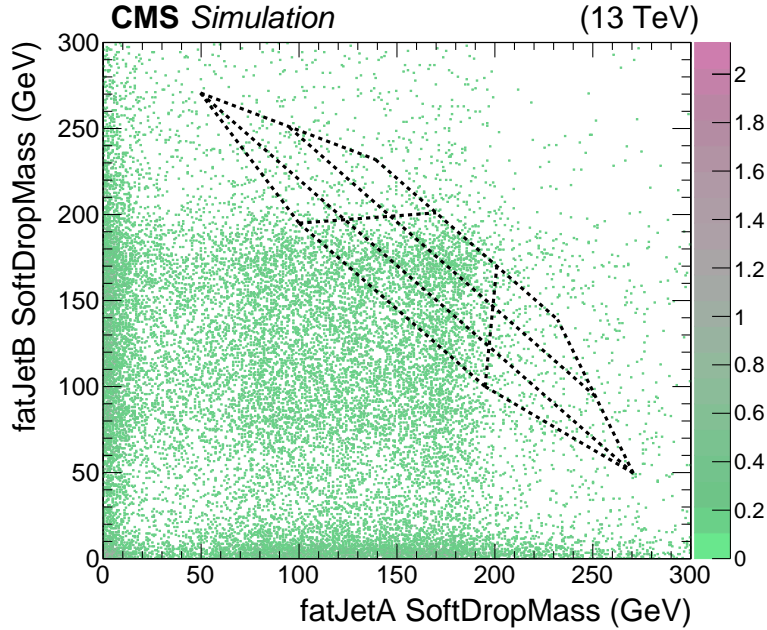
Figure F.2: Distribution of $t\bar{t}$ events in the 2D soft-drop mass plane, with the new mass regions overlaid. The events have passed the tag' double-b-tag requirement and the new kinematic cuts.

labelled the tag' double-b-tag region. The tag' double-b-tag region corresponds to where both the selected AK8 jets have double-b-tag discriminators greater than 0.3. It is similar to the original tag double-b-tag region, but the area is twice as large. It should be noted that there is no change in the anti-tag double-b-tag region.

The kinematic cuts applied in the new event selection are as follows:

- Both of the selected AK8 jets must have $p_T > 300$ GeV.

- $H_T > 1500$ GeV.

These kinematic cuts are similar to those applied in the main analysis. One difference is that the AK4 jet $p_T$ cut is removed. This is done in order to increase statistics. Another difference is that the $H_T$ binning is removed. This is because there are too few events in the $H_T \in 2500$-$3500$ and $3500+$ GeV bins to draw any conclusions. Due to the removal of the $H_T$ binning, the search region index, $i$, is always the same as the 2D mass region index, $n$.

The symbol $\hat{\tilde{S}}_{i,\ \text{t}\bar{\text{t}}}^{\ \text{tag'}}$ represents the number of t$\bar{\text{t}}$ events passing the new event selection in the $\widetilde{S}_i$ mass region with tag' double-b-tag. It is this quantity that needs estimating using the data. This is done by subtracting the equivalent V+jets (W+jets and Z+jets) and QCD yields away from the yield in data:

$$\hat{\tilde{S}}_{i,\ \text{t}\bar{\text{t}}\ \text{estimate}}^{\ \text{tag'}} = \hat{\tilde{S}}_{i,\ \text{data}}^{\ \text{tag'}} - \hat{\tilde{S}}_{i,\ \text{V}}^{\ \text{tag'}} - \hat{\tilde{S}}_{i,\ \text{QCD}}^{\ \text{tag'}} \tag{F.1}$$

The V+jets yields in Equation F.1 are estimated using MC. This means that the data driven t$\bar{\text{t}}$ estimate has a dependency on the V+jets MC. The dependence is small, however, because the t$\bar{\text{t}}$ yields are significantly larger than the V+jets contributions, as can be seen in Figure F.3. The QCD yields in Equation F.1 are predicted using the data driven QCD estimation method described in Section 6.3. The method uses the QCD yields in the different 2D mass regions and 2D double-b-tag regions:

$$\hat{\tilde{S}}_{i,\ \text{QCD}}^{\text{tag'}} = \widetilde{F}_i \cdot (\hat{\tilde{U}}_{i,\ \text{QCD}}^{\ \text{tag'}} + \hat{\tilde{D}}_{i,\ \text{QCD}}^{\ \text{tag'}}) \tag{F.2}$$

where,

$$\widetilde{F}_i \equiv \frac{\hat{\tilde{S}}_{i,\ \text{QCD}}^{\ \text{anti-tag}}}{\hat{\tilde{U}}_{i,\ \text{QCD}}^{\ \text{anti-tag}} + \hat{\tilde{D}}_{i,\ \text{QCD}}^{\ \text{anti-tag}}} \tag{F.3}$$

The equations F.2 and F.3 can then be written in terms of the yields observed in data:

$$\hat{\tilde{S}}_{i,\ \text{QCD}}^{\text{tag'}} = \widetilde{F}_i \cdot \left( \hat{\tilde{U}}_{i,\ \text{data}}^{\ \text{tag'}} + \hat{\tilde{D}}_{i,\ \text{data}}^{\ \text{tag'}} - (\hat{\tilde{U}}_{i,\ \text{t}\bar{\text{t}}}^{\ \text{tag'}} + \hat{\tilde{D}}_{i,\ \text{t}\bar{\text{t}}}^{\ \text{tag'}}) - (\hat{\tilde{U}}_{i,\ \text{V}}^{\ \text{tag'}} + \hat{\tilde{D}}_{i,\ \text{V}}^{\ \text{tag'}}) \right) \tag{F.4}$$

where,

$$\widetilde{F}_i = \frac{\hat{\tilde{S}}_{i,\ \text{data}}^{\ \text{anti-tag}} - \hat{\tilde{S}}_{i,\ \text{t}\bar{\text{t}}}^{\ \text{anti-tag}} - \hat{\tilde{S}}_{i,\ \text{V}}^{\ \text{anti-tag}}}{(\hat{\tilde{U}}_{i,\ \text{data}}^{\ \text{anti-tag}} + \hat{\tilde{D}}_{i,\ \text{data}}^{\ \text{anti-tag}}) - (\hat{\tilde{U}}_{i,\ \text{t}\bar{\text{t}}}^{\ \text{anti-tag}} + \hat{\tilde{D}}_{i,\ \text{t}\bar{\text{t}}}^{\ \text{anti-tag}}) - (\hat{\tilde{U}}_{i,\ \text{V}}^{\ \text{anti-tag}} + \hat{\tilde{D}}_{i,\ \text{V}}^{\ \text{anti-tag}})} \tag{F.5}$$

The equations F.4 and F.5 have a dependency on the t$\bar{\text{t}}$ and V+jets processes. These yields are estimated using MC simulation. This means that the data driven estimate
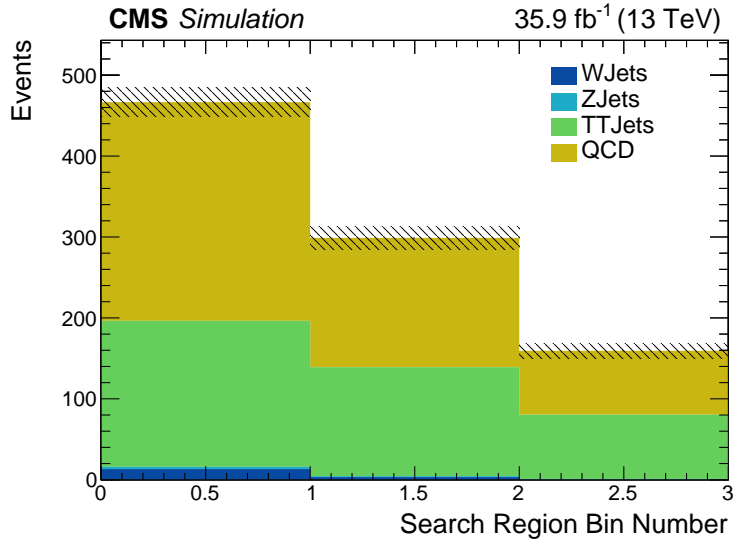
Figure F.3: Simulated yields for the 2016 Standard Model processes in the three new search regions with tag' double-b-tag and $\widetilde{S}_n$ mass regions.

of the $t\bar{t}$ yield, in the $\widetilde{S}_i$ mass region with tag' double-b-tag, has a dependence on the $t\bar{t}$ MC yields in the different 2D mass and 2D double-b-tag regions. The $\widetilde{F}_i$ factor is comprised of yields in the anti-tag double-b-tag region. The $t\bar{t}$ and V+jets contributions are very small in this double-b-tag region, as can be seen in Figure F.5 and Figure F.6. Consequently, the calculation of the $\widetilde{F}_i$ factor has very little dependence on the $t\bar{t}$ and V+jets MC. The largest dependency that the data driven $t\bar{t}$ estimate has on the $t\bar{t}$ MC, arises from the yields corresponding to the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions with tag' double-b-tag. Although this dependence is not ideal, it is limited by the design of the new 2D mass regions. This is why the shape is such that the $t\bar{t}$ event density in the $\widetilde{S}_n$ mass regions is a lot higher than in the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions. Consequently, the fraction of $t\bar{t}$ events is much smaller in the $\widetilde{U}_n + \widetilde{D}_n$ sideband mass regions compared to the $\widetilde{S}_n$ mass regions. This can be seen by comparing Figure F.4 with Figure F.3. Inaccuracies in the $t\bar{t}$ MC will have a factor six less impact on the data driven $t\bar{t}$ estimate. Thus, comparing this data driven $t\bar{t}$ estimate with the equivalent yields in $t\bar{t}$ MC is still a worthwhile comparison.

Figure F.4: Simulated yields for the 2016 Standard Model processes in the three new search regions with tag' double-b-tag and $\widetilde{U}_n + \widetilde{D}_n$ mass regions.



Figure F.5: Simulated yields for the 2016 Standard Model processes in the three new search regions with anti-tag double-b-tag and $\widetilde{S}_n$ mass regions.
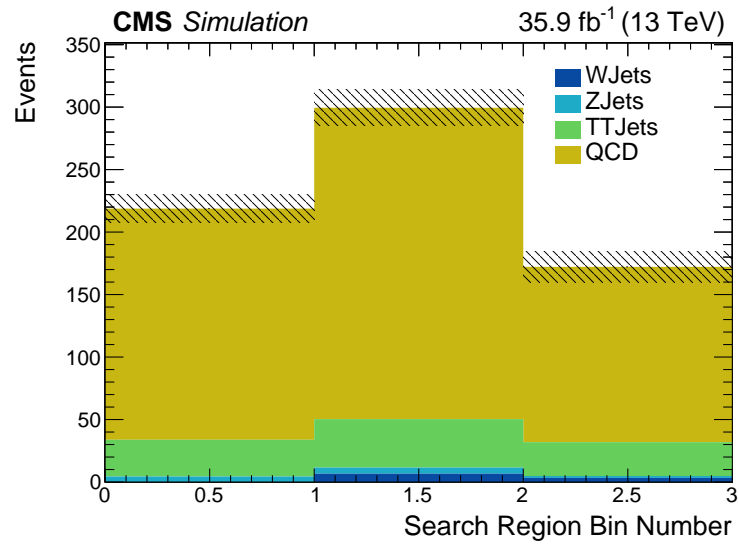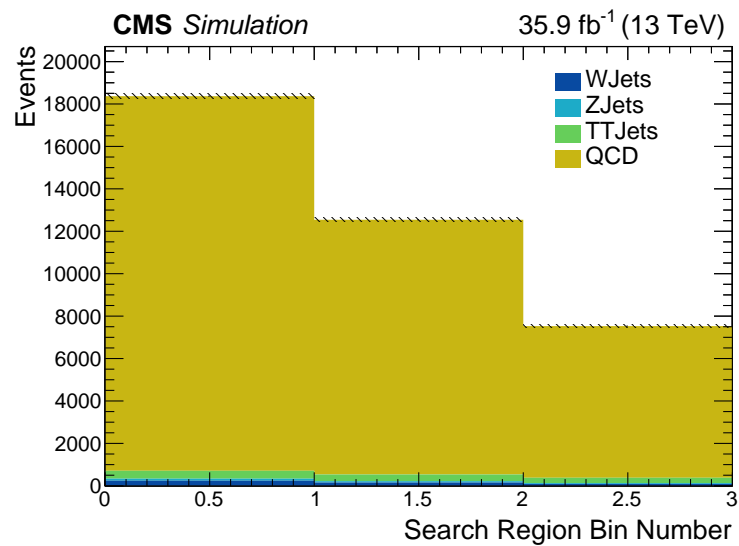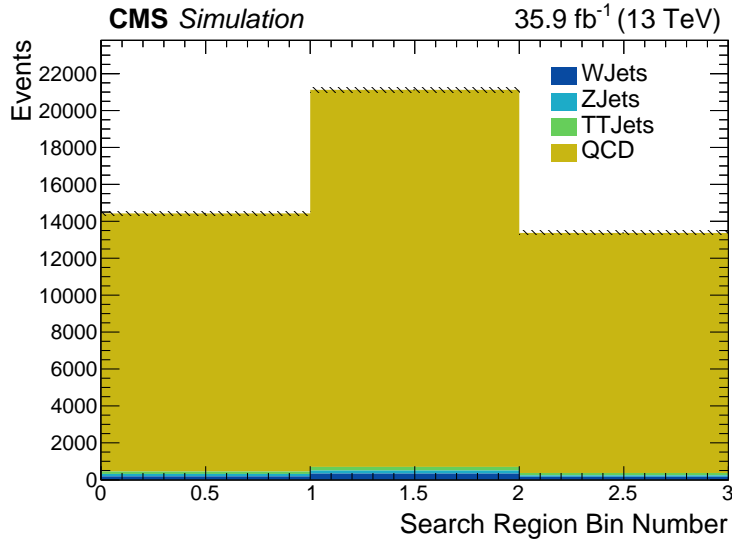
Figure F.6: Simulated yields for the 2016 Standard Model processes in the three new search regions with anti-tag double-b-tag and $\widetilde{U}_n + \widetilde{D}_n$ mass regions.

There are various uncertainties considered when calculating the the data driven $t\bar{t}$ estimate. These include the statistical uncertainties of all the yields that enter the calculation. There is also a 10% error assigned to the data driven QCD estimation method. This is uncorrelated between each search region. Finally, the uncertainties on the $t\bar{t}$ double-b-tag mistag scale factors are considered, which lead to a correlated 20% error in the $t\bar{t}$ MC yield. However, as was stated above, the calculation only has a small dependence on the $t\bar{t}$ MC. The combined uncertainty is quite large, especially in the first search region, and this ultimately limits what can be concluded from the test. The size of the errors are driven by the statistical uncertainties. Unfortunately, there is no other way to increase the statistics whilst remaining similar to the main analysis.

Figure F.7 compares the data driven $t\bar{t}$ estimate, in the three $\widetilde{S}_n$ mass regions with tag' double-b-tag, with the yield drawn directly from the $t\bar{t}$ MC. For the MC yield, the majority of the error arises from the uncertainty on the $t\bar{t}$ double-b-tag mistag scale factors. Consequently, these errors are strongly correlated between the three search regions. The comparison is reasonable and suggests that the $t\bar{t}$ MC is providing a broadly accurate description in the event selection spaces used in this
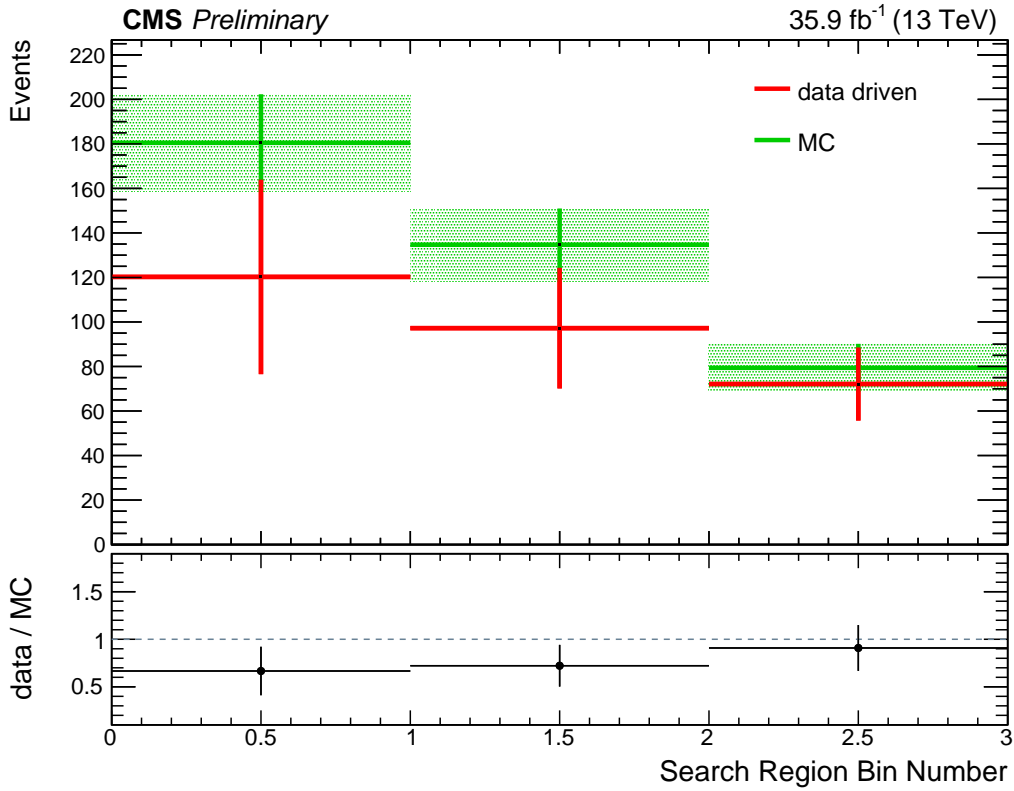
Figure F.7: Comparison of the data driven tt̄ estimate (red) with the equivalent yield from tt̄ MC (green). The is done for the three search regions with with tag' double-b-tag and $\widetilde{S}_n$ mass regions. Note that the errors for the MC yields are strongly correlated between the three search regions.

analysis. Given that a systematic uncertainty of 50% is used for the tt̄ cross section in the main analysis (see Section 7.1.5), it seems that the tt̄ MC is okay to use.

It should be noted that the first iteration of this test was performed using an old training of the double-b-tag BDT. The old training was not as good at suppressing the SM backgrounds and, therefore, provided greater statistics for the test. With the smaller associated uncertainties, a much better agreement was shown between the data driven tt̄ estimate and the yield drawn directly from the tt̄ MC.

# G

# Signal Injection Test

To demonstrate the capability of the statistical methods employed in the analysis, the observed upper limits of $\sigma/\sigma_{\mathrm{theory}}$ were recalculated using data that was artificially injected with signal. Specifically, the data yields, $n_{y,i,m}$ (adopting the notation introduced in Section 7.2.1), were modified so that they had the value:

$$n'_{y,i,m} = n_{y,i,m} + s_{y,i,m} \tag{G.1}$$

Three different signal models were injected, with parameters $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 1600$, 2000, and 2400 GeV.

Figure G.1 shows the resultant exclusion contours in the $M_{\mathrm{H}}$-$M_{\mathrm{SUSY}}$ plane. In all three signal injection scenarios, the $M_{\mathrm{SUSY}}$ limit decreases as $M_{\mathrm{H}}$ approaches 70 GeV. The minimum $M_{\mathrm{SUSY}}$ limit is always less than the $M_{\mathrm{SUSY}}$ value of the injected signal model. This is what one would expect from this demonstration, thus, validating the statistical methods employed in the analysis.

Figure G.1: Observed signal model exclusions, as a function of $M_{\mathrm{H}}$ and $M_{\mathrm{SUSY}}$, for different signal injection scenarios. Black contour: No signal injection. Blue contour: Signal with $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 1600$ GeV injected. Red contour: Signal with $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2000$ GeV injected. Green contour: Signal with $M_{\mathrm{H}} = 70$ GeV and $M_{\mathrm{SUSY}} = 2400$ GeV injected.

# Glossary

**BDT**: Boosted Decision Tree

**BSM**: Beyond the Standard Model

**CHS**: Charged Hadron Subtraction

**CMS**: Compact Muon Solenoid

**CSC**: Cathode Strip Chamber

**DM**: Dark Matter

**DT**: Drift Tube

**EB**: ECAL Barrel

**ECAL**: Electromagnetic Calorimeter

**EE**: ECAL Endcap

**ES**: ECAL Preshower

**FPGA**: Field-Programmable Gate Array

**GCT**: Global Calorimeter Trigger

**HB**: HCAL Barrel

**HCAL**: Hadronic Calorimeter

**HE**: HCAL Endcap

**HF**: HCAL Forward

**HLT**: High-Level Trigger

**HO**: HCAL Outer

**ISR**: Initial State Radiation

**JECs**: Jet Energy Corrections

**JER**: Jet Energy Resolution

**LHC**: Large Hadron Collider

**LSP**: Lightest Supersymmetric Particle

**LUT**: Lookup Table

**MC**: Monte Carlo

**MSSM**: Minimal Supersymmetric Standard Model

**NLSP**: Next-to Lightest Supersymmetric Particle

**NMSSM**: Next-to Minimal Supersymmetric Standard Model

**PDF**: Parton Distribution Function

**PF**: Particle-Flow

**POG**: Physics Object Group

**PS**: Proton Synchrotron

**PSB**: Proton Synchrotron Booster

**PUPPI**: Pile-up Per Particle Identification

**QCD**: Quantum Chromodynamics

**RPC**: Resistive Plate Chamber

**RCT**: Regional Calorimeter Trigger

**SM**: Standard Model

**SPS**: Super Proton Synchrotron

**SUSY**: Supersymmetry / Supersymmetric

**TEC**: Tracker Endcap

**TIB**: Tracker Inner Barrel

**TID**: Tracker Inner Disk

**TMT**: Time-Multiplexed Trigger

**TOB**: Tracker Outer Barrel

# References

[1] S. Heinemeyer et al. *Handbook of LHC Higgs Cross Sections: 3. Higgs Properties: Report of the LHC Higgs Cross Section Working Group.* CERN Yellow Reports: Monographs. Jul 2013.

[2] T. Sakuma and T. McCauley. Detector and Event Visualization with SketchUp at the CMS Experiment. *J. Phys. Conf. Ser.*, 513:022032, 2014.

[3] CMS Collaboration. The CMS experiment at the CERN LHC. *JINST*, 3(08):S08004–S08004, 2008.

[4] CMS Collaboration. *CMS Physics: Technical Design Report Volume 1: Detector Performance and Software.* Technical Design Report CMS. CERN, 2006.

[5] CMS Collaboration. *CMS Technical Design Report for the Level-1 Trigger Upgrade.* Technical Design Report CMS. CERN/LHCC, 2013.

[6] CMS Collaboration. L1 calorimeter trigger upgrade: jet and energy sum performance and commissioning status. Technical Report CMS-DP-2015-051, CERN, Oct 2015.

[7] M. Tanabashi et al. Review of particle physics. *Phys. Rev. D*, 98:030001, Aug 2018.

[8] M. Thomson. *Modern particle physics.* Cambridge University Press, New York, 2013.

[9] M. K. Gaillard, P. D. Grannis, and F. J. Sciulli. The Standard model of particle physics. *Rev. Mod. Phys.*, 71:S96–S111, 1999.

[10] M. E. Peskin and D. V. Schroeder. *An Introduction to quantum field theory.* Addison-Wesley, Reading, USA, 1995.

[11] K. Garrett and G. Duda. Dark Matter: A Primer. *Adv. Astron.*, 2011:968283, 2011.

[12] S. P. Martin. A Supersymmetry primer. pages 1–98, 1997. [Adv. Ser. Direct. High Energy Phys.18,1(1998)].

[13] I. J. R. Aitchison. *Supersymmetry in Particle Physics. An Elementary Introduction.* 2007.

[14] U. Ellwanger, C. Hugonie, and A.M. Teixeira. The Next-to-Minimal Supersymmetric Standard Model. *Phys. Rept.*, 496:1–77, 2010.

[15] Y. Baconnier et al. *LHC: the Large Hadron Collider accelerator project.* CERN, 1993.

[16] T. Pettersson and P. Lefevre. The Large Hadron Collider: conceptual design. Technical Report CERN-AC-95-05-LHC, Oct 1995.

[17] G. Aad et al. The ATLAS Experiment at the CERN Large Hadron Collider. *JINST*, 3:S08003, 2008.

[18] A. Alves et al. The LHCb Detector at the LHC. *JINST*, 3:S08005, 2008.

[19] The ALICE Collaboration. The ALICE experiment at the CERN LHC. *JINST*, 3:S08002, 2008.

[20] CMS Collaboration. CMS: letter of intent by the CMS Collaboration for a general purpose detector at LHC. Technical Report CERN-LHCC-92-003, CERN, 1992.

[21] CMS Collaboration. CMS Physics: Technical Design Report Volume 2: Physics Performance. *J. Phys. G*, 34(CERN-LHCC-2006-021), 2007.

[22] CMS Collaboration. Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC. *Phys. Lett. B*, (CMS-HIG-12-028), Jul 2012.

[23] CMS Collaboration. Precise determination of the mass of the Higgs boson and tests of compatibility of its couplings with the standard model predictions using proton collisions at 7 and 8 TeV. *Eur. Phys. J. C*, 75(CMS-HIG-14-009), Dec 2014.

[24] H. Schopper. *LEP: The lord of the collider rings at CERN 1980-2000: The making, operation and legacy of the world's largest scientific instrument.* 2009.

[25] CMS Collaboration. *The CMS magnet project: Technical Design Report.* Technical Design Report CMS. CERN, 1997.

[26] V. Karimaki et al. *The CMS tracker system project: Technical Design Report.* Technical Design Report CMS. CERN, 1997.

[27] CMS Collaboration. *The CMS tracker: addendum to the Technical Design Report.* Technical Design Report CMS. CERN, 2000.

[28] C. Grupen and B. Shwartz. *Particle Detectors.* Cambridge University Press, 2008.

[29] A. Dominguez et al. CMS Technical Design Report for the Pixel Detector Upgrade. Technical Report CERN-LHCC-2012-016. CMS-TDR-11, Sep 2012.

[30] CMS Collaboration. *The CMS electromagnetic calorimeter project: Technical Design Report.* Technical Design Report CMS. CERN, 1997.

[31] P. Adzic et al. Energy resolution of the barrel of the CMS electromagnetic calorimeter. *JINST*, 2:P04004, 2007.

[32] CMS Collaboration. *The CMS hadron calorimeter project: Technical Design Report.* Technical Design Report CMS. CERN, 1997.

[33] S. Abdullin et al. The CMS barrel calorimeter response to particle beams from 2-GeV/c to 350-GeV/c. *Eur. Phys. J.*, C60:359–373, 2009.

[34] CMS Collaboration. *The CMS muon project: Technical Design Report.* Technical Design Report CMS. CERN, 1997.

[35] S. Agostinelli et al. GEANT4: A Simulation toolkit. *Nucl. Instrum. Meth.*, A506:250–303, 2003.

[36] CMS Collaboration. Particle-flow reconstruction and global event description with the CMS detector. *JINST*, 12(CMS-PRF-14-001), Jun 2017.

[37] S. Chatrchyan et al. Description and performance of track and primary-vertex reconstruction with the CMS tracker. *JINST*, 9(10):P10009, 2014.

[38] G. Salam. Towards Jetography. *Eur. Phys. J.*, C67:637–686, 2010.

[39] CMS Collaboration. Pileup Removal Algorithms. Technical Report CMS-PAS-JME-14-001, CERN, 2014.

[40] CMS Collaboration. Jet energy scale and resolution in the CMS experiment in pp collisions at 8 TeV. *JINST*, 12(02):P02014, 2017.

[41] D. Bertolini, P. Harris, M. Low, and N. Tran. Pileup Per Particle Identification. *JHEP*, 10:059, 2014.

[42] J. Shelton. Jet Substructure. In *Proceedings, Theoretical Advanced Study Institute in Elementary Particle Physics: Searching for New Physics at Small and Large Scales (TASI 2012): Boulder, Colorado, June 4-29, 2012*, pages 303–340, 2013.

[43] S. D. Ellis, J. Huston, K. Hatakeyama, P. Loch, and M. Tonnesmann. Jets in hadron-hadron collisions. *Prog. Part. Nucl. Phys.*, 60:484–551, 2008.

[44] M. Dasgupta, A. Fregoso, S. Marzani, and G. Salam. Towards an understanding of jet substructure. *Journal of High Energy Physics*, 2013(9):29, Sep 2013.

[45] A. Larkoski, S. Marzani, G. Soyez, and J. Thaler. Soft Drop. *JHEP*, 05:146, 2014.

[46] G. Brandenburg et al. Charged track multiplicity in $B$ meson decay. *Phys. Rev.*, D61:072002, 2000.

[47] CMS Collaboration. Identification of b quark jets at the CMS Experiment in the LHC Run 2. Technical Report CMS-PAS-BTV-15-001, CERN, 2016.

[48] CMS Collaboration. Identification of double-b quark jets in boosted event topologies. Technical Report CMS-PAS-BTV-15-002, CERN, 2016.

[49] CMS Collaboration. *CMS TriDAS project: Technical Design Report, Volume 1: The Trigger Systems.* Technical Design Report CMS. CERN/LHCC, 2000.

[50] CMS Collaboration. *CMS The TriDAS Project : Technical Design Report, Volume 2: Data Acquisition and High-Level Trigger.* Technical Design Report CMS. CERN/LHCC, 2002.

[51] A. Zabi et al. Triggering on electrons, jets and tau leptons with the CMS upgraded calorimeter trigger for the LHC RUN II. *JINST*, 11:C02008, 2016.

[52] L. Mastrolorenzo. The CMS Level-1 Tau identification algorithm for the LHC Run II. Technical Report CMS-CR-2014-309, CERN, Oct 2014.

[53] J. Sauvan and the CMS Collaboration. Performance and upgrade of the CMS electron and photon trigger for Run 2. *Journal of Physics: Conference Series*, 587(1):012021, 2015.

[54] M. Cacciari, J. Rojo, G. P. Salam, and G. Soyez. Jet reconstruction in heavy ion collisions. *The European Physical Journal C*, 71(1):1539, Jan 2011.

[55] A. Tapper and the CMS Collaboration. The CMS Level-1 Trigger for LHC Run II. *ICHEP 2016*, 282(242), 2016.

[56] CMS Collaboration. Level-1 jets and energy sums trigger performance with full 2017 dataset. Technical Report CMS-DP-2018-004, CERN, Feb 2018.

[57] U. Ellwanger and A. Teixeira. NMSSM with a singlino LSP: possible challenges for searches for supersymmetry at the LHC. *JHEP*, 10:113, 2014.

[58] CMS Collaboration. Search for supersymmetry in multijet events with missing transverse momentum in proton-proton collisions at 13 TeV. *Phys. Rev. D*, 96:032003, 2017.

[59] CMS Collaboration. Search for new phenomena with the $M_{T2}$ variable in the all-hadronic final state produced in proton-proton collisions at $\sqrt{s} = 13$ TeV. *The European Physical Journal C*, page 77, 2017.

[60] CMS Collaboration. Search for natural and split supersymmetry in proton-proton collisions at $\sqrt{s}$= 13 TeV in final states with jets and missing transverse momentum. *JHEP*, 05:025, 2018.

[61] ATLAS Collaboration. Search for squarks and gluinos in final states with jets and missing transverse momentum using 36 fb$^{-1}$ of $\sqrt{s} = 13$ TeV pp collision data with the ATLAS detector. *Phys. Rev. D*, 97:112001, 2018.

[62] U. Ellwanger, J. F. Gunion, and C. Hugonie. NMHDECAY: A Fortran code for the Higgs masses, couplings and decay widths in the NMSSM. *JHEP*, 02:066, 2005.

[63] J. Alwall, M. Herquet, F. Maltoni, O. Mattelaer, and T. Stelzer. MadGraph 5 : Going Beyond. *JHEP*, 06:128, 2011.

[64] S. Alioli et al. A general framework for implementing NLO calculations in shower Monte Carlo programs: the POWHEG BOX. *Journal of High Energy Physics*, 2010(6):43, Jun 2010.

[65] R. D. Ball et al. Parton distributions for the LHC run II. *Journal of High Energy Physics*, 2015(4):40, Apr 2015.

[66] R. D. Ball et al. Parton distributions from high-precision collider data. *Eur. Phys. J.*, C77(10):663, 2017.

[67] T. Sjostrand, S. Mrenna, and P. Z. Skands. A Brief Introduction to PYTHIA 8.1. *Comput. Phys. Commun.*, 178:852–867, 2008.

[68] CMS Collaboration. Event generator tunes obtained from underlying event and multiparton scattering measurements. *Eur. Phys. J.*, C76(3):155, 2016.

[69] CMS Collaboration. Investigations of the impact of the parton shower tuning in Pythia 8 in the modelling of $t\bar{t}$ at $\sqrt{s} = 8$ and 13 TeV. Technical Report CMS-PAS-TOP-16-021, CERN, Geneva, 2016.

[70] CMS Collaboration. Extraction and validation of a new set of CMS PYTHIA8 tunes from underlying-event measurements. Technical Report CMS-GEN-17-001, 2019.

[71] M. Czakon and A. Mitov. Top++: A Program for the Calculation of the Top-Pair Cross-Section at Hadron Colliders. *Comput. Phys. Commun.*, 185:2930, 2014.

[72] A. Djouadi, J. Kalinowski, and M. Spira. HDECAY: A Program for Higgs boson decays in the standard model and its supersymmetric extension. *Comput. Phys. Commun.*, 108:56–74, 1998.

[73] W. Beenakker, R. Hopker, and M. Spira. PROSPINO: A Program for the production of supersymmetric particles in next-to-leading order QCD. 1996.

[74] R. D. Ball et al. Parton distributions with LHC data. *Nucl. Phys.*, B867:244–289, 2013.

[75] ATLAS and CMS Collaborations. Procedure for the LHC Higgs boson search combination in Summer 2011. Technical Report CMS-NOTE-2011-005, CERN, Geneva, Aug 2011.

[76] Cowan, G. and Cranmer, K. and Gross, E. and Vitells, O. Asymptotic formulae for likelihood-based tests of new physics. *The European Physical Journal C*, page 71, 2011.

[77] Searches for new phenomena in events with jets and high values of the $M_{\mathrm{T}2}$ variable, including signatures with disappearing tracks, in proton-proton collisions at $\sqrt{s}$ = 13 TeV. Technical Report CMS-PAS-SUS-19-005, CERN, Geneva, 2019.

[78] CMS Collaboration. Search for Physics Beyond the Standard Model in Events with High-Momentum Higgs Bosons and Missing Transverse Momentum in Proton-Proton Collisions at 13 TeV. *Phys. Rev. Lett.*, 120(24):241801, 2018.

[79] R. Barate et al. Search for the standard model Higgs boson at LEP. *Phys. Lett.*, B565:61–75, 2003.

[80] CMS Collaboration. Observation of Higgs boson decay to bottom quarks. *Phys. Rev. Lett.*, 121(12):121801, 2018.

[81] ATLAS collaboration. Observation of $H \to b\bar{b}$ decays and $VH$ production with the ATLAS detector. *Phys. Lett.*, B786:59–86, 2018.

[82] U. Ellwanger and A. Teixeira. Excessive higgs pair production with little MET from squarks and gluinos in the NMSSM. *JHEP*, 04:172, 2015.

[83] A. Steane. *Relativity Made Relatively Easy.* Oxford University Press, 2012.