

# **Detector optimisation for future linear collider**

Boruo Xu  
of King's College

A dissertation submitted to the University of Cambridge  
for the degree of Doctor of Philosophy



## Abstract

This is my abstract. To be or not to be.



## Declaration

This dissertation is the result of my own work, except where explicit reference is made to the work of others, and has not been submitted for another qualification to this or any other university. This dissertation does not exceed the word limit for the respective Degree Committee.

Boruo Xu



## Acknowledgements

Of the many people who deserve thanks, some are particularly prominent, such as my supervisor . . .



## **Preface**

This will be my preface. Where is Wolly?



# Contents

<b>1 Let's make introduction great again</b>	<b>1</b>
1.1 Future Linear Colliders . . . . .	1
1.2 Motivation . . . . .	1
<b>2 Theoretical overview</b>	<b>3</b>
2.1 Overview of the Standard Model . . . . .	3
2.2 Notations and conventions . . . . .	5
2.3 Quantum electrodynamics . . . . .	5
2.4 Quantum chromodynamics . . . . .	6
2.5 The electroweak interaction . . . . .	6
2.6 Higgs Mechanism . . . . .	8
2.7 Yukawa couplings . . . . .	8
2.8 Standard Model Higgs boson . . . . .	9
2.9 Higgs beyond the Standard Model . . . . .	10
<b>3 Detector and Physics at Future Linear Colliders</b>	<b>15</b>
3.1 ILC . . . . .	15
3.2 CLIC . . . . .	16
3.3 CLIC vs ILC . . . . .	16
3.4 Physics at future linear colliders . . . . .	17
3.5 Impact of physics requirements on the detector design . . . . .	18
3.5.1 Jet energy resolution requirements on the detector design . . . . .	18
3.5.2 Other requirements on the detector design . . . . .	20
3.6 International Large Detector . . . . .	21
3.6.1 ILD vs SiD . . . . .	22
3.7 Overview of ILD sub-detectors . . . . .	22
3.7.1 Vertex Detector . . . . .	23
3.7.2 Tracking Detectors . . . . .	23
3.7.3 Electromagnetic Calorimeter . . . . .	24

3.7.4	Hadronic Calorimeter . . . . .	25
3.7.5	Solenoid . . . . .	26
3.7.6	Yoke and Muon system . . . . .	27
3.7.7	Very Forward Calorimeters . . . . .	27
3.7.8	ILD vs CLIC_ILD . . . . .	29
<b>4</b>	<b>Simulation and Reconstruction</b>	<b>31</b>
4.1	Monte Carlo event generation . . . . .	31
4.2	Event Simulation . . . . .	32
4.3	Event Reconstruction . . . . .	32
4.4	PandoraPFA . . . . .	33
4.4.1	Track selection . . . . .	33
4.4.2	Calorimeter selection . . . . .	34
4.4.3	Cone Clusters Algorithm . . . . .	34
4.4.4	Particle Identification . . . . .	35
4.4.5	Clustering . . . . .	35
4.4.6	Topological cluster association . . . . .	35
4.4.7	Track-cluster association . . . . .	36
4.4.8	Re-clustering . . . . .	36
4.4.9	Fragment removal . . . . .	36
4.4.10	Particle Flow Object Creation . . . . .	36
4.5	MC truth linker . . . . .	37
4.6	CLIC specific simulation and reconstruction . . . . .	37
4.7	Luminosity spectrum . . . . .	37
4.8	Suppression of $\gamma\gamma \rightarrow \text{hadrons}$ backgrounds . . . . .	37
4.9	CLIC simulated particle masses . . . . .	38
4.10	Reconstruction Processors . . . . .	38
4.11	Jet algorithm . . . . .	38
4.11.1	Durham algorithm . . . . .	40
4.11.2	$y$ parameter . . . . .	40
4.11.3	LCFIPlus . . . . .	41
4.11.4	Optimisation and overfitting . . . . .	42
4.11.5	Choice of models . . . . .	43
4.11.6	Optimisation of Boosted Decision Tree . . . . .	46
4.11.7	Multiple classes . . . . .	48
4.12	Event shape variables . . . . .	48

4.13 Miscellaneous . . . . .	49
<b>5 Photon Reconstruction in PandoraPFA</b>	<b>51</b>
5.1 Electromagnetic shower . . . . .	51
5.2 Overview of photon reconstruction in PandoraPFA . . . . .	51
5.3 Photon reconstruction algorithm . . . . .	52
5.3.1 Form photon clusters . . . . .	53
5.3.2 Reconstruct photon candidates . . . . .	53
5.3.3 Photon ID test . . . . .	54
5.3.4 Photon Fragment removal . . . . .	54
5.4 Two dimensional peak finding algorithm for photon candidate . . . . .	54
5.4.1 Candidate close to track projection . . . . .	55
5.4.2 Peak filtering . . . . .	56
5.4.3 Inclusive mode . . . . .	56
5.5 Likelihood classifier for photon ID . . . . .	57
5.5.1 Overview of Projective Likelihood . . . . .	57
5.5.2 Projective Likelihood in PandoraPFA . . . . .	57
5.6 Photon fragment removal algorithm in the ECAL . . . . .	59
5.7 High energy photon fragment recovery algorithm . . . . .	61
5.8 Photon splitting algorithm . . . . .	64
5.9 Photon reconstruction performance improvement . . . . .	65
5.10 Breakdown of photon reconstruction improvement . . . . .	68
5.11 Photon reconstruction performance . . . . .	69
<b>6 Tau Lepton Final State Separation</b>	<b>71</b>
6.1 Introduction . . . . .	71
6.2 Simulation and reconstruction . . . . .	72
6.3 Generator level cut . . . . .	72
6.4 Decay modes . . . . .	73
6.5 Discriminative variables . . . . .	74
6.6 Multivariate Analysis . . . . .	78
6.7 Result . . . . .	79
6.8 Electromagnetic calorimeter optimisaiton . . . . .	80
<b>7 Double Higgs Bosons Production Analysis</b>	<b>85</b>
7.1 Analysis Straggly Overview . . . . .	86
7.2 Monte Carlo Sample Generation . . . . .	87

7.3	Lepton identification . . . . .	88
7.3.1	Electron and muon identification . . . . .	90
7.3.2	Tau identification . . . . .	92
7.3.3	Very forward electron identification . . . . .	94
7.3.4	Lepton identification performance . . . . .	96
7.3.5	Other lepton identification processors . . . . .	97
7.4	Jet reconstruction . . . . .	98
7.4.1	Jet reconstruction optimisation . . . . .	98
7.5	Jet flavour tagging . . . . .	102
7.5.1	Deploying LCFIPlus . . . . .	102
7.6	Jet pairing . . . . .	103
7.7	Pre-selection . . . . .	104
7.7.1	Discriminative pre-selection cuts . . . . .	104
7.7.2	Loose cuts for the MVA . . . . .	105
7.7.3	Mutually exclusive cuts for $\text{HH} \rightarrow b\bar{b}W^+W^-$ and $\text{HH} \rightarrow b\bar{b}b\bar{b}$ . . . . .	108
7.8	Discriminative variables for MVA . . . . .	109
7.9	Multivariate analysis . . . . .	111
7.10	Signal selection results . . . . .	114
7.11	$e^-e^+ \rightarrow \text{HH}\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$ hadronic decay at $\sqrt{s} = 3 \text{ TeV}$ analysis .	114
7.12	$e^-e^+ \rightarrow \text{HH}\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$ semi-leptonic decay at $\sqrt{s} = 3 \text{ TeV}$ analysis	122
7.13	Result interpretation . . . . .	125
7.14	Simultaneous couplings extraction . . . . .	128
	<b>Bibliography</b>	<b>135</b>
	<b>List of figures</b>	<b>139</b>
	<b>List of tables</b>	<b>141</b>

*“Two bags of pork scratchings are worth  
a bag of gold.”*

— Joris the Dutch



# Chapter 1

## Let's make introduction great again

*“Introduction means introduction”*

— Theresa Trump

Introduction

### 1.1 Future Linear Colliders

Basic intro. LHC.

Next challenge

Future Options. FCC vs LC

LC options

### 1.2 Motivation

Photon - passage through matter. Photon electromagnetic shower

Since Higgs discovery in the LHC in 2012, Higgs

Ha there is a higgs.

We found higgs. Higgs is cool. It explains mass.

Why double higgs. Double higgs coupling is unique to linear collider. It can reveal much about the BSM models.

Generator level study has performed. ILC has done this this and that. gHHH in CLIC before

Here we do things differently. First subchannels, then extract both couplings simultaneously.

# Chapter 2

## Theoretical overview

*“ILC will be built next year”*

— Mysterious person

This chapter provides a theoretical overview which would be used in the subsequent chapters. A short review of the Standard Model of Particle Physics, the current best particle theory, is provided, with an emphasis on the Higgs mechanism and the Higgs boson. A general parametrisation of the Higgs theory, beyond the Standard Model, is discussed, and supplies the theoretical background for the physics analysis in the chapter [7](#).

### 2.1 Overview of the Standard Model

The Standard Model (SM) is a quantum field theory concerning three fundamental interactions of nature: the electromagnetic, weak and strong integrations. The SM also describes the interactions between the sub-atomic particles. The deployment and the experimental verification of the SM throughout the second half of the 20th century is one of the greatest triumph of the particle physics. The most recent discovery of the Higgs boson in 2012 [1] further verified the theory. This chapter summarise the SM based on the summaries of the SM [2–5].

The fundamental particles in the SM consist of three categories: force exchange bosons, leptons and neutrinos, and quarks. In the SM, the force exchange bosons carries the fundamental forces between particles. For example, photon is the force carrier of the

electromagnetic force.  $W^+$ ,  $W^-$ , and  $Z$  are the force carriers of the weak force. Gluon,  $g$ , is the force carrier of the strong force. These bosons will be discussed further in section 2.5. There is also the Higgs boson from the spontaneous symmetry breaking of the Higgs field, which is discussed in section 2.8

Another category of fundamental particles contains leptons and neutrinos. These particles are fermions. For each fermion in the SM, there is an anti-fermion with same mass and spin, but opposite charge. Leptons and neutrinos have three generations. Each generation has same interaction properties, but different masses. Although neutrinos could not be directly detected, measurements of the  $Z$  decay width strongly suggested the three generations of neutrinos [6]. Leptons and neutrinos experience weak forces as well as electromagnetic forces, which will be further discussed in section 2.5.

The last category of fundamental particles are quarks, which are also fermions and have three generations. Each generation has a positively charged up type quark and a negatively charged down type quark. Quarks experience all three fundamental forces described by the SM.

The SM has enjoyed great success with theoretical predictions being experimentally verified. Some highlights included the discovery of the top quark in 1995 cite, the tau neutrino in 2000 cite and the Higgs bosons in 2012 citeAad:2012tfa. However, there are observations which are not explained by the SM. One issue is that the SM does not incorporate the gravitational force. There have been attempts to modify the SM but no conclusive theory yet. Another issue is that the SM does not allow neutrino masses and mixings. There have been many theories beyond the Standard Model (BSM). One such example is the generalisation of the Higgs theory to allow non SM coupling strengths. This will be discussed in section 2.9.

Notations and conventions will be introduced. The overview of the Standard Model starts with the quantum electrodynamics, and its generalisation to quantum chromodynamics. The unification of electromagnetism and weak interaction, electroweak gauge theory will be discussed. Afterwards, Higgs mechanism and Yukawa couplings will be introduced to explain bosons and fermions masses whilst preserving the Lagrangian symmetry. This will be followed by a detailed discussion on the Standard Model Higgs boson, its mass and interactions with other particles. The chapter finishes with explanation for possible Higgs theories beyond the Standard Model, the Lagrangian of the Higgs interaction, and observables, which forms the theoretical background for the analysis on the double Higgs production in the chapter 7.

## 2.2 Notations and conventions

The natural unit is used in this document,  $\hbar = c = 1$ . The metric is mostly-minus,  $\eta^{\mu\nu} = \text{diag}(1, -1, -1, -1)$ . The Dirac gamma matrices are represented with  $\gamma^\mu$ , with  $\mu$  goes from 0 to 3.  $\gamma^5 = i\gamma^0\gamma^1\gamma^2\gamma^3$ .  $\bar{\psi} = \psi^\dagger\gamma^0$ . Einstein summation convention is used as well this document.

This set of notations allows a contracted pair to be a Lorentz invariant. For a Weyl spinor,  $\psi_\alpha$ , the mass term in the lagrangian is of the form  $\psi^\alpha\psi_\alpha$ , which is the Majorana mass term. The contracted pair between two different Weyl spinors would form a Dirac mass term.

## 2.3 Quantum electrodynamics

The natural starting point to introduce the SM is the quantum electrodynamics (QED). The QED is a quantum field theory explaining electromagnetic interactions. The theory involves a spin-half Dirac (electron) field  $\psi$  and a vector (photon) field  $A_\mu$ . When the local (gauge) symmetry is imposed, which is equivalent to the Lagrangian invariance under transformations,

$$\psi \rightarrow e^{ie\phi(x)}, A_\mu \rightarrow A_\mu - \partial^\mu\phi(x), \quad (2.1)$$

the Lagrangian is fixed to be

$$\mathcal{L}_{\text{QED}} = \bar{\psi} (i\gamma^\mu D_\mu - m) \psi - \frac{1}{4} F_{\mu\nu} F^{\mu\nu}, \quad (2.2)$$

if up to cubic terms are allowed in the fields. There are two free parameters in the QED,  $m$  the electron mass, and  $e$  the electron charge. The mass term for the photon,  $\nabla^2 A_{\mu\nu} A^\nu$ , is forbidden by gauge invariance.

The QED has been verified experimentally. One of the greatest prediction is the spin magnetic dipole moment of the electron, defined as  $\vec{\mu} = g_s \frac{Qe}{2m} \vec{s}$ . The  $g_s$  is predicted to be 2 by the Dirac equations. The small corrections to the value comes from the electron's interaction with virtual photons, so called higher "loop" corrections in Feynman diagrams. The precise agreement of the theoretical prediction and the experimental value is a success of the QED.

## 2.4 Quantum chromodynamics

Like the QED, the quantum chromodynamics (QCD) a quantum field theory explaining strong interactions. There are eight gauge bosons, gluons, coupling to nine fermions, quarks. Unlike the QED, the theory is invariant under local non-Abelian SU(3) transformations. Gluons can interact with other gluons, and carry colour charge (red, green, and blue). Nine quarks transform as colour triplets. The QCD Lagrangian is

$$\mathcal{L}_{\text{QCD}} = \sum_{f \in u, d, s, c, b, t} \bar{\psi} \left( i\gamma^\mu \partial_\mu - g_s \gamma^\mu G_\mu^a \frac{\lambda^a}{2} - m_f \right) \psi - \frac{1}{4} G_{\mu\nu}^a G^{a\mu\nu}, \quad (2.3)$$

where  $g_s$  is the strong coupling constant.  $a$  is the colour charge.  $\lambda$  is the Gell-Mann matrices.  $G_{\mu\nu}^a$  is the gluon field strength, given by

$$G_{\mu\nu}^a = \partial_\mu \gamma_\nu^a - \partial_\nu \gamma_\mu^a - g_s f_{abc} G_\mu^b G_\nu^c. \quad (2.4)$$

The last extra term comparing to QED indicates the non-Abelian nature of the QCD.

## 2.5 The electroweak interaction

The electroweak interaction can be thought of a extension to the QED to incorporate the weak force, and to explain different coupling strength to left-handed and right-handed fermions. However, the Lagrangian does not allow the massive electroweak force exchange bosons and fermions, which are explained by the Higgs mechanism and the Yukawa interactions.

There are four vector boson fields in the theory, 3  $W$  and 1  $B$  field. The Lagrangian can be divided into two parts: the bosonic self interaction and the couplings to the fermions.

$$\mathcal{L}_{\text{Electroweak}} = \mathcal{L}_{\text{Boson}} + \mathcal{L}_{\text{Fermion}} \quad (2.5)$$

The bosonic self interaction Lagrangian,  $\mathcal{L}_{\text{Boson}}$ , is given by

$$\mathcal{L}_{\text{Boson}} = -\frac{1}{4} W_{\mu\nu}^i W^{i\mu\nu} - \frac{1}{4} B_{\mu\nu} B^{\mu\nu}, \quad (2.6)$$

where

$$W_{\mu\nu}^i = \partial_\nu W_\mu^i - \partial_\mu W_\nu^i - g \varepsilon^{ijk} W_\mu^j W_\nu^k \quad (2.7)$$

$$B_{\mu\nu} = \partial_\nu B_\mu - \partial_\mu B_\nu \quad (2.8)$$

$B$  field is invariant under  $U(1)$ .  $W$  field is invariant under non-Abelian  $SU(2)$  transformations.  $g$  is the coupling strength of the  $W$  field. The indices,  $i$ ,  $j$ , and  $k$  indicates 3  $W$  fields, going from 1 to 3.

The fermionic part of the Lagrangian,  $\mathcal{L}_{\text{Fermion}}$ , has different components for the left-handed and right-handed fermions, given by

$$\mathcal{L}_{\text{Fermion}} = \sum_{\psi \in \text{fermions}} \bar{\Psi}_L \gamma^\mu D_\mu^L \Psi_L + \bar{\Psi}_R \gamma^\mu D_\mu^R \Psi_R \quad (2.9)$$

$D_\mu^L$  and  $D_\mu^R$  are defined as

$$D_\mu^L = \partial_\mu + ig \frac{\tau_i}{2} W_\mu^i + ig' Y_\psi B_\mu \quad (2.10)$$

$$D_\mu^R = \partial_\mu + ig' Y_\psi B_\mu \quad (2.11)$$

This Lagrangian allows  $W$  and  $B$  field to couple with left-handed fermions, but only  $B$  field couples to right-handed fermions. The  $\tau_i$  matrices are the generators of the  $SU(2)$ . Pauli spin matrices are one of the representations.  $Y_\psi$  is the hypercharge associating with the fermion field  $\psi$ .  $g'$  is the  $B$  field strength.

Physical bosons  $W^+$ ,  $W^-$  only couples to left-handed fermions.  $Z$  and  $\gamma$  couples to both left-handed and right-handed fermions. Hence  $W^1$  and  $W^2$  are associated to  $W^+$  and  $W^-$ . Mass eigenstates for  $Z$  and  $\gamma$ ,  $Z_\mu$  and  $A_\mu$  are mixture of  $W_\mu^3$  and  $B_\mu$ .

$$Z_\mu = \cos(\theta_W) W_\mu^3 - \sin(\theta_W) B_\mu \quad (2.12)$$

$$A_\mu = \sin(\theta_W) W_\mu^3 + \cos(\theta_W) B_\mu \quad (2.13)$$

$\theta_W$  is the Weinberg mixing angle [7], which is determined experimentally. So far the  $SU(2) \otimes U(1)$  gauge theory explains the parity violating nature of the weak interaction.

The explicit fermion mass are not allowed in the gauge symmetry. The higgs mechanism via spontaneous symmertry breaking would introduce mass terms for fermions.

## 2.6 Higgs Mechanism

A complex scalar Higgs field,  $\Phi_H$ , is added to the electroweak Lagrangian.  $\Phi_H$  transforms as a doublet of SU(2) with hypercharge  $Y = \frac{1}{2}$

$$\mathcal{L}_{\text{Higgs}} = (D_\mu \Phi_H)^\dagger (D^\mu \Phi_H) - \mu^2 \Phi_H^\dagger \Phi_H - \lambda (\Phi_H^\dagger \Phi_H)^2 \quad (2.14)$$

with

$$D_\mu \Phi_H = \left( \partial_\mu + ig \frac{\tau_i}{2} W_\mu^i + ig' \frac{1}{2} B_\mu \right) \Phi_H \quad (2.15)$$

For negative  $\mu^2$ , the Higgs filed potential

$$\mu^2 \Phi_H^\dagger \Phi_H + \lambda (\Phi_H^\dagger \Phi_H)^2 \quad (2.16)$$

is minimised with a Higgs vacuum potential  $\frac{v}{\sqrt{2}} = \sqrt{\frac{v^2}{2\lambda}}$ . After the symmetry breaking, the  $\mathcal{L}_{\text{Higgs}}$  provides the mass terms for  $W^+$ ,  $W^-$ ,  $Z$  and  $\gamma$  via terms in the Lagrangian:

$$\frac{(gv)^2}{4} W_\mu^+ W^{-\mu} + \frac{(g^2 + g'^2) \mu^2}{8} Z_\mu Z^\mu \quad (2.17)$$

This provides equal mass for  $W^+$  and  $W^-$  with massless photon.

## 2.7 Yukawa couplings

Section 2.6 explains the Higgs mechanism for gauge bosons gaining masses. The fermions gain masses in a similar fashion. Consider a Higgs field transforming as a doublet of SU(2) with hypercharge  $Y = \frac{1}{2}$ , the Yukawa couplings is given

$$\mathcal{L}_{\text{Yukawa}} = -\lambda^u \bar{q}_L \Phi_H^c u_R - -\lambda^d \bar{q}_L \Phi_H d_R - \lambda^e \bar{l}_L \Phi_H e_R + \text{h.c.} \quad (2.18)$$

$\Phi_H^c \equiv i\sigma^2 H^*$  is an SU(2) doublet field with hypercharge  $Y = -\frac{1}{2}$  and  $\sigma$  is the Pauli spin matrix.  $u$ ,  $d$ , and  $e$  are fields for up-type quark, down-type quark and leptons. The Lagrangian is summed over all possible quarks and leptons. The interactions terms in the  $\mathcal{L}_{Yukawa}$  become mass terms when the Higgs vacuum expectation value is substituted. The fermion masses are given by

$$m_u = \frac{\lambda^{u\nu}}{\sqrt{2}}, \quad m_d = \frac{\lambda^{d\nu}}{\sqrt{2}}, \quad m_e = \frac{\lambda^{e\nu}}{\sqrt{2}} \quad (2.19)$$

## 2.8 Standard Model Higgs boson

So far, interactions between different fields in the Standard Model, as well as the mass obtaining mechanism have been discussed. Only left for discussion is the Higgs bosons, and its interactions with other fields.

For the Higgs doublet complex field in the SM, there are four real scalar degrees of freedom. By choosing the unitary gauge, three degree of freedoms are manifestly eaten. The Higgs field becomes

$$H(x) = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ \nu + h(x) \end{pmatrix} \quad (2.20)$$

$h(x)$  is real scalar field for the Higgs boson. It is not charged under electromagnetism as it is real. The Higgs boson interaction terms in the Lagrangian with other particles can be shown by replacing  $\nu$  with  $\nu + h(x)$  in previous expressions. For fermions field,  $\psi_i$ , the Higgs boson interaction term is given by

$$\mathcal{L} \supset -\frac{m_i}{\nu} h \bar{\psi}_i \psi_i \quad (2.21)$$

From the equation 2.17, the Higgs boson interaction terms for bosons can be shown as

$$\mathcal{L} \supset m_W^2 \left( \frac{2h}{\nu} + \frac{h^2}{\nu^2} \right) W_\mu^+ W^{-\mu} + \frac{m_Z^2}{2} \left( \frac{2h}{\nu} + \frac{h^2}{\nu^2} \right) Z_\mu Z^\mu \quad (2.22)$$

The Higgs boson self interactions are obtained from the Higgs field potential

$$\mathcal{L} \supset \frac{\mu^2}{2} (\nu + h)^2 - \frac{\lambda}{4} (\nu + h)^4 \supset -\lambda\nu^2 h^2 - \lambda\nu h^3 - \frac{\lambda}{4} h^4 = -\frac{m_h^2}{2} h^2 - \frac{m_h^2}{2\nu} h^3 - \frac{m_h^2}{8\nu^2} h^4. \quad (2.23)$$

The Higgs boson mass,  $m_h$  is  $2\lambda\nu^2$ . The triple and quartic self interaction strengths are  $-\frac{m_h^2}{2\nu}$  and  $\frac{m_h^2}{8\nu^2}$ . Once the  $m_h$  is determined,  $\lambda$  can be worked out. The Higgs boson decay width and branching fraction can be roughly worked out. For example, figure 2.1a and figure 2.1b show partial decay width and the branching ratios as a function of  $m_h$ .

The Higgs boson decaying to a pair of heavier particles, such as  $W^+W^-$  or  $Z Z$  is forbidden kinematically. However, figure 2.1 shows that the Higgs decaying to  $W^+W^-$  dominates before the mass threshold  $m_H = 2m_W \sim 160$  GeV. This is because one of the  $W^\pm$  gauge bosons is virtual and not on the mass shell, which is allowed by the quantum field theory. The virtual gauge boson subsequently decays to real on-shell particles.



**Figure 2.1:** figure 2.1a shows Standard Model Higgs boson partial widths as a function of its mass,  $M_H$ . The total width is the black curve figure 2.1b shows selected Standard Model Higgs boson branching ratios as a function of its mass,  $M_H$ . Both plots are taken from [8]

## 2.9 Higgs beyond the Standard Model

Since the SM like Higgs boson discovery in 2012, it becomes important to understand of role of the Higgs boson in the electroweak spontaneous symmetry breaking. In the absence of the Higgs boson, the coupling strength of the longitudinally polarised vector

bosons grows with energy and becomes strong at TeV scale. The SM Higgs bosons moderates the interacting strength, allowing the extraction of the weak coupling at short distances. In this scenario, the SM Higgs couplings are constrained and predicted in one parameter only, the Higgs mass. But other alternative scenarios could allow the behaviour of the SM Higgs at low energy. One such example, motivated by the hierarchy problem and the electroweak data, is that the light and narrow Higgs-scalar is a composite bound state of some strongly interacting sector at the TeV scale. The couplings of the Higgs to fermions and bosons would be different to those in the SM. If the composite Higgs is the pseudo Nambu-Goldstone boson from a spontaneous global symmetry breaking, the Higgs can be naturally light [9]. Another scenario is that a composite dilaton, the pseudo Nambu-Goldstone boson arose from a spontaneous scale invariance breaking, partially behaves like a light Higgs [10]. In both scenarios, the interaction of Higgs becomes strong at high energy. The coupling of the Higgs would deviate to those in the SM.

An important physics channel for testing the Higgs theory is the double Higgs production via vector boson fusion at high energy [11–13]. For the composite Higgs scenario, the scattering amplitude increases with the energy. For the dilaton scenario, no energy dependence on the scattering amplitude is expected. It is difficult for the Large Hadron Collider to measure the cross section due to the large SM background rate [12]. However, a multi-TeV linear electron position collider, such as Compact Linear Collider, would be able to precisely measure the cross section [14].

Following the assumption made in the [12, 13], the self interaction of the light scalar Higgs,  $h$ , and its coupling to other SM bosons can be described by the following Lagrangian. The notation in the [13] is followed. After the electroweak symmetry breaking, the bosonic part of the Lagrangian reads:

$$\mathcal{L} = \frac{1}{2}(\partial_\mu h) - V(h) + \left(m_W^2 W_\mu^+ W^{-\mu} + \frac{m_Z^2}{2} Z_\mu Z^\mu\right) \left[1 + 2a \frac{h}{v} + b \frac{h^2}{v^2} + \dots\right], \quad (2.24)$$

where  $V(h)$  is the  $h$  field potential,

$$V(h) = \frac{1}{2} m_h^2 h^2 + d_3 \left(\frac{m_h^2}{2v}\right) h^3 + d_4 \left(\frac{m_h^2}{8v^2}\right) h^4 + \dots \quad (2.25)$$

$a$ ,  $b$ ,  $d_3$  and  $d_4$  are arbitrary dimensionless parameters. Higher order terms in  $h$  are omitted.  $a$  and  $b$  are proportional to the coupling strength of the  $VWh$  and  $VWhh$  vertices, where  $V = W^\pm, Z$ .  $d_3$  and  $d_4$  are proportional to the triple and quartic  $h$  self coupling strength. Comparing with equation 2.22 and equation 2.23, the SM Higgs

suggests  $a = b = d_3 = d_4 = 1$  and all higher order terms vanish. The dilaton scenario imposes the relation,  $a = b^2$ .

The scattering amplitude for  $V_L V_L \rightarrow hh$  can be written as

$$A = a^2(A_{SM} + A_1\delta_b + A_2\delta_{d_3}), \quad (2.26)$$

where  $A_{SM}$  is the SM amplitude and

$$\delta_b \equiv 1 - \frac{b}{a^2}, \quad (2.27)$$

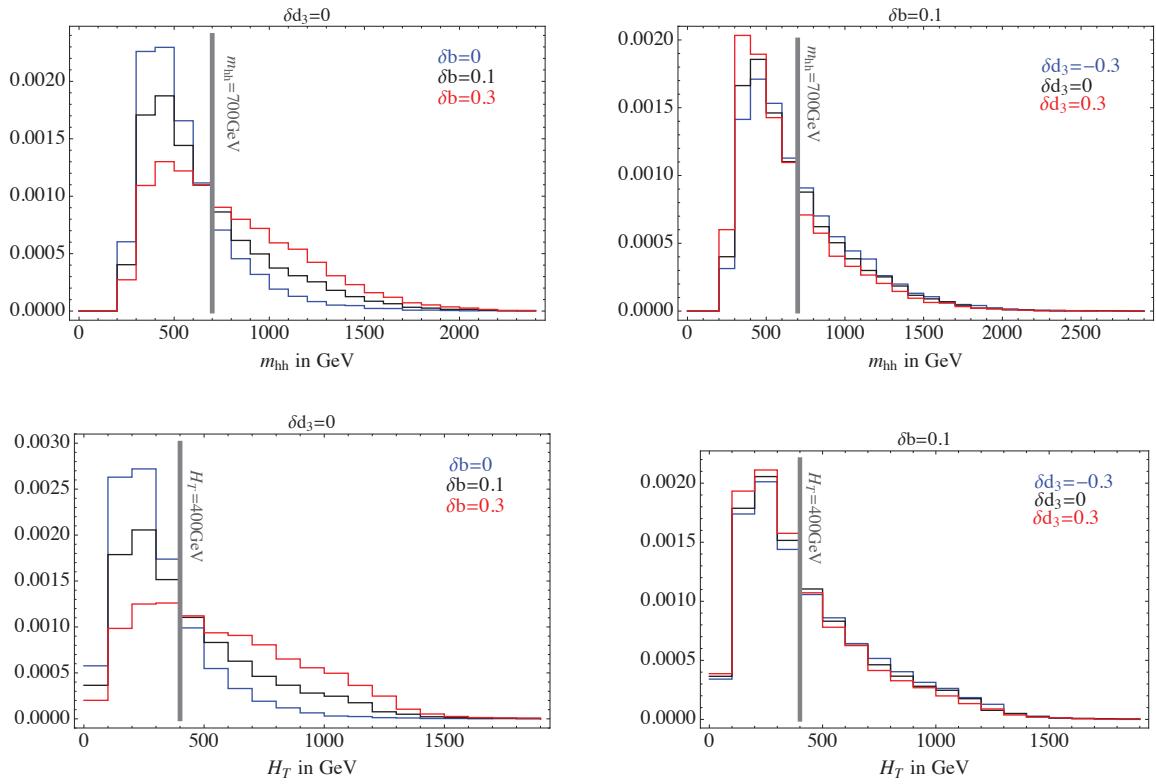
$$\delta_{d_3} \equiv 1 - \frac{d_3}{a}. \quad (2.28)$$

$A_1$  grows like energy squared at large center-of-mass energy,  $E \gg m_V$ .  $A_{SM}$  and  $A_2$  has no energy depended. Therefore,  $\delta_b$  controls the magnitude of the scattering amplitude increasing as a function of energy.  $\delta_{d_3}$ , however, determines the magnitude at threshold. In an electron-positron collider, this scattering process and be studied via  $e^+e^- \rightarrow v\bar{v}hh$  channel. The cross section of the channel can be written as

$$\sigma = a^4 \sigma_{SM} (1 + A\delta_b + B\delta_{d_3} + C\delta_b\delta_{d_3} + D\delta_b^2 + E\delta_{d_3}^2), \quad (2.29)$$

where  $\sigma_{SM}$  is the cross section predicted by the SM. With suitable kinematic cuts, high-energy behaviour can be disentailed from the physics at threshold, allowing the extraction of  $\delta_b$ ,  $\delta_{d_3}$  and hence the coupling strength  $g_{VHH}$  and  $g_{HHH}$ . Suitable observables are variables that increases with increasing centre-of-mass energy. Two examples of such variables are the invariant mass of the two Higgses system,  $m_{hh}$ , and the sum of their transverse momenta,  $H_T$ . figure ?? shows that the  $m_{hh}$  and  $H_T$  distributions are sensitive to the values of  $\delta_b$  and  $\delta_{d_3}$ . The figure shows the result of a general level study performed in [13].

In the equation 2.29,  $a$ , which is proportional to  $g_{VHH}$ , only enters as an overall factor. However,  $a$  also appears in the definition of  $\delta_b$  and  $\delta_{d_3}$ . To extract  $\delta_b$  and  $\delta_{d_3}$ ,  $a$  should be a constant ideally. For a multi-TeV electron-positron collider, the cross section for single Higgs production is far greater than that of the double Higgs. Figure 2.3 shows the comparison of the cross section as a function of the centre-of-mass energy. Therefore,  $g_{VHH}$  and  $a$  will be measured precisely before investigating the double Higgs production.



**Figure 2.2:** Normalized differential cross sections  $d\sigma/dm_{hh}$  and  $d\sigma/dH_T$  for  $e^+e^- \rightarrow \nu\bar{\nu}hh$  at the Compact Linear Collider, with  $\sqrt{s} = 3$  TeV after the identification cuts, for several values of  $\delta_b$  and  $\delta_{d_3}$ . Plot is taken from [13].

For the purpose of measuring  $g_{WWH}$  and  $g_{HHH}$  via double Higgs production,  $\alpha$  in the equation 2.29 can be treated as a constant.



**Figure 2.3:** Cross section as a function of centre-of-mass energy for the Higgs production processes at an electron-positron collider for a Higgs mass of 126GeV. The values shown correspond to unpolarised beams and do not include the effect of beamstrahlung. Plot is taken from [15].

# Chapter 3

## Detector and Physics at Future Linear Colliders

*“ILC will be built next year”*

— Mysterious person

Since the discovery of a particle consistent with being the SM Higgs boson in LHC at 2012 [1, 16], our understanding of Standard Model has improved greatly. Yet limited by the underlying QCD interaction from proton-anti-proton collision, one has great difficulty to measure the properties of the Higgs precisely. Next generation electron-positron linear collider could hopefully make precision measurements of the Higgs sector and the Top quark sector [?].

### 3.1 ILC

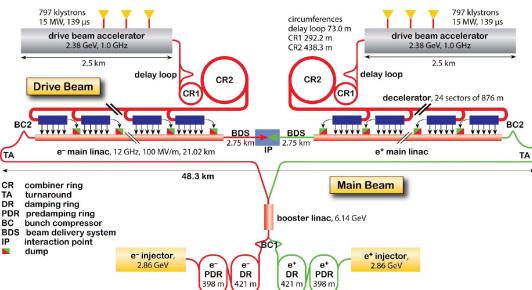
Two leading candidates for next generation electron-positron linear collider are the International Linear Collider (ILC) [17], and the Compact Linear Collider (CLIC) [18]. The ILC is a high-luminosity electron-positron linear collider with centre-of-mass energy from 200 GeV up to 1 TeV. The machine would be build at different stages. The first stage would have a centre-of-mass energy of 250/350 GeV. The second stage would be 500 GeV with a possible upgrade to 1 TeV. Thirty years of development leads to the technical design report in 2013 [19]. A layout of the collider complex is shown in figure 3.1.



**Figure 3.1:** A layout of the International Linear Collider complex, taken from [19].

## 3.2 CLIC

The other potential next generation electron-positron linear collider, the Compact Linear Collider (CLIC), has a higher reach of the centre-of-mass energy up to 3 TeV. The CLIC is designed as a staged machine. The first stage, with centre-of-mass energy 380 GeV, is a compromise of precision measurement between both top quark and Higgs physics. The final stage 3 TeV is motivated by the physics reach of detecting new physics, and measurement rare decays of Higgs. The second stage is around 1.4 TeV, which bridges between the first stage and the final stage. A layout of the CLIC complex is shown in figure 3.2.



**Figure 3.2:** A layout of the Compact Linear Collider at 3 TeV, taken from [20].

## 3.3 CLIC vs ILC

Due to the similarities of the two two linear collider programs, the development with CLIC detector concepts start with ILC detector concepts. CLIC\_ILD and CLIC\_SiD are developed based on ILD and SiD. The two main differences are the high centre-of-mass energy and the 0.5 ns between bunch crossings at CLIC. More incoherent pairs and more hadronic two-photon events are produced at high energy. Particles produced in

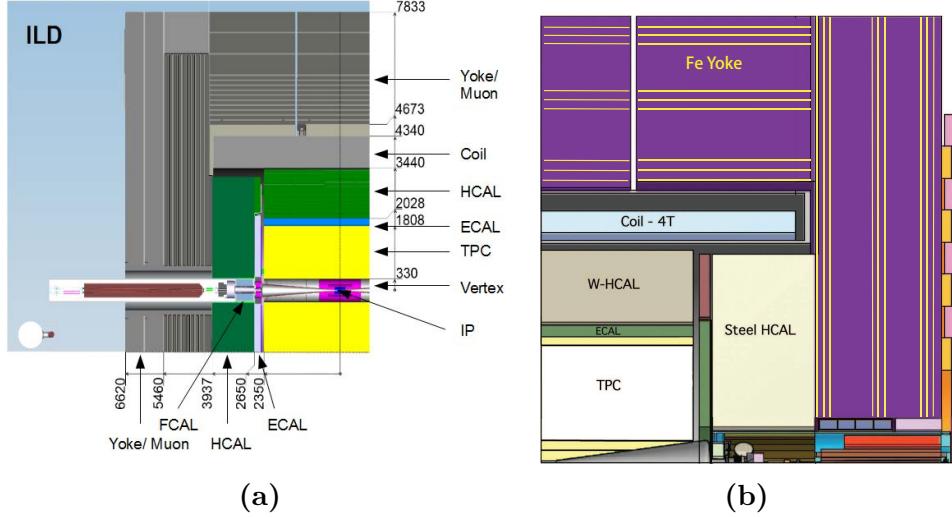
the forward region via t-channel are more important due to a stronger boost. These differences leads to a modification in the detector design and the reconstruction software for the CLIC. A comparison of CLIC\_ILD and ILD longitudinal cross sections can be seen in figure 3.3. A comparison of key parameters of the ILD and CLIC\_ILD detector concepts is shown in table 3.1.

Concept	ILD	CLIC_ILD
Tracker	TPC/Silicon	TPC/Silicon
Solenoid Field (T)	3.5	4
Solenoid Field Bore (m)	3.3	3.4
Solenoid Length (m)	8.0	8.3
VTX Inner Radius (mm)	16	31
ECAL $r_{\min}$ (m)	1.8	1.8
ECAL $\Delta r$ (mm)	172	172
HCAL Absorber B / E	Fe	Fe / W
HCAL Interaction Length	5.5	7.5
Overall Height (m)	14.0	14.0
Overall Length (m)	13.2	12.8

**Table 3.1:** A comparison of key parameters of the ILD and CLIC\_ILD detector concepts. ECAL  $r_{\min}$  is the smallest distance from the calorimeter to the main detector axis. HCAL Absorber B / E indicates the absorber material for the barrel (B) and the endcap (E). The table is adapted from [18].

## 3.4 Physics at future linear colliders

The physics program for the CLIC and the ILC, which is a driving force for the detector design, have some common goals. ILC has a reach of centre-of-mass from 200 GeV to 1 TeV, whilst CLIC can reach from 350 GeV to 3 TeV. Both machines are capable of precision higgs coupling measurements, top mass and coupling measurements, search for new physics such as supersymmetry particles. The ILC can also operate at low energy to be a Z and a H factory for ultra precision Z mass and H mass measurement. CLIC, however, has the advantage of higher energy reach, which allows measurements of rare events, such as higgs triple self-couplings and quartic couplings. The discovery potential for the CLIC is also greater.



**Figure 3.3:** Figure 3.3a and Figure 3.3a shows longitudinal cross section of top quadrant of the ILD and the CLIC-ILD detector concepts, taken from [19] and [18] respectively. From interaction point (IP) outwards, there is a tracking system comprising a large time projection chamber (TPC) augmented with silicon tungsten layer, highly granular electromagnetic calorimeters (ECAL) and hadronic calorimeters (HCAL), muon chambers, forward calorimeters (FCAL), magnetic coils and iron yokes. Numbers are in units of mm.

### 3.5 Impact of physics requirements on the detector design

#### 3.5.1 Jet energy resolution requirements on the detector design

The physics goal of jet energy resolution at the ILC and the CLIC is to separate W and Z hadronic decays via reconstruction their di-jet masses [18, 19]. This translates to a requirement of 3.5-5% of the energy resolution. This level of precision is unlikely to be achieved with a traditional calorimetry design. A traditional energy flow to calorimetry measures jet energies as a sum of the energy in the calorimeters. The jet energy resolution is parameterised by

$$\frac{\sigma_E}{E} = \frac{\alpha}{\sqrt{E(\text{GeV})}} \oplus \beta \quad (3.1)$$

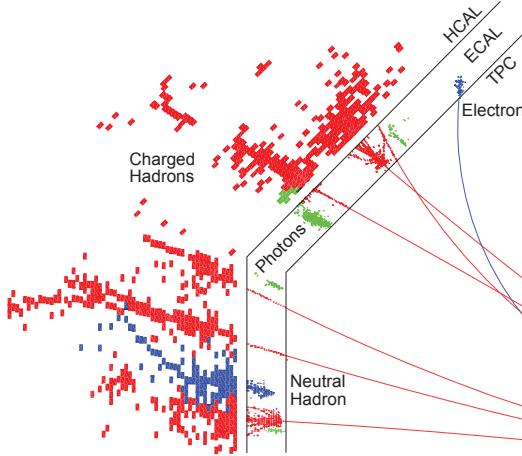
The stochastic term  $a$  is typically greater than 60% and  $b$  is order of a few percents. For the jet energy resolution of 3.5%,  $a$  should be less than 30% which is unlikely to be achieved by a traditional calorimeter. On the contrary, particle flow approach has demonstrated its ability to reach the goal [21, 22].

In a typical jet, using measurements on the particle composition from the LEP [23, 24], about 62% of the jet energy is from charged particles, 27% from photons, 10% from long-lived neutral hadrons, and 1.5% from neutrinos. In a traditional approach to calorimetry, about 72% jet energy is measured in the electromagnetic (ECAL) and the hadronic (HCAL) calorimeters combined. The jet energy resolution is thus limited by the energy resolution of the hadronic calorimeters, which typically is  $\gtrsim 55\%/\sqrt{E(\text{GeV})}$ .

The particle flow approach to calorimetry improves the jet energy resolution by fully reconstructing the four momenta of all visible particles in the detector. The jet energy is the sum of the individual particles' energies, where the energy of the charge particles are measured in the tracking detectors, and the energy of neutral particles are measured in calorimeters. In this manner, the hadronic calorimeter only measures about 10% of the energy, which would greatly improve the overall energy measurement. Assuming 30% of the jet energy, photon energy, is measured with  $\sigma_E/E = 15\%/\sqrt{E(\text{GeV})}$ , and 10% of the jet energy, which are hadrons, measured with  $\sigma_E/E = 55\%/\sqrt{E(\text{GeV})}$ , a jet energy of  $\sigma_E/E = 19\%/\sqrt{E(\text{GeV})}$  can be obtained. This satisfies the jet energy resolution for separating W and Z hadronic decays. In reality, this level of performance is unattainable due to incorrect association of energy deposits to particles. At jet energy beyond tens of GeVs, the “confusion” rather than the intrinsic detector performance limits the particle flow performance. This has stringent requirements in the ECAL and the HCAL design.

The particle flow calorimetry requires to fully reconstruct particles and associate calorimeter hits to tracks. This is demanding for the software design and the detector design. The software details of the PandoraPFA, which is a successful particle flow implementation, are described in section 4.4. The detector needs to be highly granular for the excellent spatial resolution to be able to correctly associate calorimeter hits to the inner detector tracks. This forms the main motivation in the calorimeter design.

Figure 3.4 is a typical topology of a 250 GeV jet, simulated with CLIC\_ILD detector concept. The particles with the calorimeter hits and tracks are labelled with different colours. Clusters of calorimeter hits in the highly granular ECAL and HCAL are associated with tracks from the inner tracking detector, TPC. Photons are identified using the characteristics longitudinal and transverse electromagnetic shower profiles.



**Figure 3.4:** A typical topology of a 250 GeV jet, simulated with CLIC-ILD detector concept, taken from [22].

Hadronic showers are separated from electromagnetic showers due to the small transverse spread of the electromagnetic shower. Therefore, the inner tracking detector should be highly efficient and has very little material. For the calorimeter, the ECAL and HCAL, both should be highly granular. The material of the calorimeter should be dense and has a large ratio of interaction length to radiation length.

### 3.5.2 Other requirements on the detector design

Other physics requirement for the detectors for the ILC and the CLIC are summarised from [18, 19].

The requirement of tracking momentum resolution is driven by the Higgs boson mass resolution via Higgsstrahlung process,  $e^-e^+ \rightarrow ZH$ . Higgs mass can be reconstructed precisely as the recoil mass against the Z momenta, which is obtained via  $Z \rightarrow \mu^+\mu^-$ . For the ILC operating at  $\sqrt{s} = 250$  GeV, the momentum resolution needs to be  $\sigma_{p_T}/p_T^2 \lesssim 5 \cdot 10^{-5} \text{ GeV}^{-1}$ . For the CLIC at high  $\sqrt{s}$ , the momentum resolution needs to be  $\sigma_{p_T}/p_T^2 \lesssim 2 \cdot 10^{-5} \text{ GeV}^{-1}$ .

The performance requirement of the vertex detector is determined by efficient b-quark and c-quark tagging. The ability to identify secondary vertices and tracks, which are not originated from the interaction point, is the prerequisite for the flavour tagging. The

impact parameter resolution can be written as

$$\sigma_{d_0}^2 = a^2 + \frac{b^2}{p^2 \sin^2(\theta)} \quad (3.2)$$

where  $a$  is related to the point resolution and  $b$  is related to multiple scattering. The requirements for both the ILC and the CLIC detectors are  $a \lesssim 5\mu\text{m}$  and  $b \lesssim 15\mu\text{m}\text{GeV}$

The lepton identification are should be over 95% for effective lepton tagging. The forward converge of the detector should be down to a very low angle. This is more critical for the CLIC as particles are boosted at high  $\sqrt{s}$ .

## 3.6 International Large Detector

The International Large Detector, ILD, is a detector concept at the International Linear Collider, ILC. The ILD detector concept has been optimised in the view of the particle flow techniques. Particle flow approach to event reconstruction has shown to deliver the best possible jet reconstruction with proof-of-principle implementation such as PandoraPFA [chapter 5](#). Each individual particles are reconstructed with the particle flow approach. For charged particles, calorimeter hits are associated with the tracks. The measurement of charged particle relies on the excellent tracking system resolution. Neutral particle reconstruction require fine spatial resolution of the calorimeters. These form the requirements for the detector designs and optimisations.

The particle flow paradigm requires topological information for individual particle reconstruction. The sub-detector systems need to have the spatial resolution to separate charged particles from neutral particles. The result is a highly granular calorimeters with a central tracking system with excellent momentum resolution. Longitudinal cross section of top quadrant of the ILD detector concept, taken from [\[?\]](#), is shown in figure [3.3a](#) From interaction point (IP) outwards, there is a tracking system compromising a large time projection chamber (TPC) augmented with silicon tungsten layer, highly granular electromagnetic calorimeters (ECAL) and hadronic calorimeters (HCAL), muon chambers, forward calorimeters (FCAL), magnetic coils and iron yokes. Numbers are in units of mm.

This section will describe the sub-systems of the ILD detector concept in the ILD technical design report [\[\]](#), the ILD\_o1\_v05 option in Mokka simulation. This detector

concept has been used in studies described in subsequent chapters. The ILD detector concept has been optimised and documented in previous documents, such as the letter of intent [1]. The CLIC\_ILD detector concept for the CLIC in the conceptual design report [18] is a modified version of the ILD, adapted to the CLIC colliding environment. The differences between ILD and CLIC\_ILD can be seen in figure 3.3 and table 3.1, and are addressed in the discussion below.

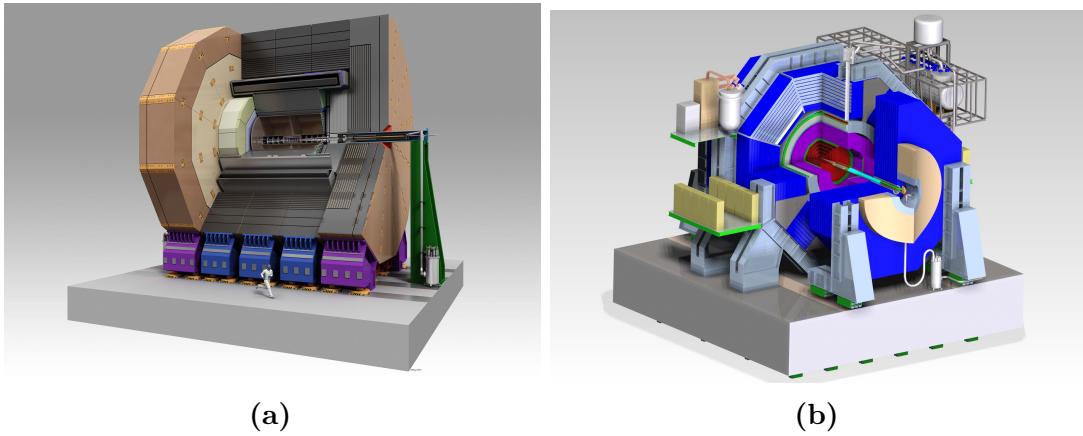
### 3.6.1 ILD vs SiD

Two detectors concepts have been designed for the ILC to deliver the physics program. Precision tests for Standard Model requires an excellent jet energy resolution and di-jet mass reconstruction. Particle Flow Algorithms based event reconstruction meets the requirement and motivates the detector designs. For the best performance of the PFA, high granular calorimeter systems and highly efficient tracking systems are designed. The requirement to separate W and Z bosons in di-jet final states requires a jet energy resolution below 3.5%. The momentum resolution of  $5 \times 10^{-5}$  GeV is motivated by the Higgs boson recoil reconstruction in the Higgs-strahlung.

Motivation for two detectors is to have multiple independent measurements within one collider for cross-checking, complementary measurements and competition between collaborations. Two detectors are both general purpose detectors. Silicon Detector, SiD, is a compact detector with a large magnetic field of 5 T. It uses silicon tracking modules. The International Large Detector, ILD, is a larger detector with a time projection chamber as the main tracking unit. Both detectors have high granular calorimeters optimised for the particle flow. A view of both detector concepts can be seen in figure 3.5

## 3.7 Overview of ILD sub-detectors

The ILD detector concept is designed as a general purpose detector. Closest to the interaction points are the precision vertex detector and a tracking system. The tracking system consists of silicon tracking with a time projection chamber. Surrounding the tracking system is a high granular calorimeter system. The outer solenoid provides a magnetic field of 3.5 T. The most outer iron return yoke acts as a muon calorimeter.



**Figure 3.5:** Figure 3.5a and figure 3.5b show the view of the International Large Detector and the Silicon Detector for the International Linear Collider, taken from [19].

### 3.7.1 Vertex Detector

The pixel-vertex detector(VTX) needs to be close to the interaction point to reconstruct secondary vertex. As the TPC is the main tracking detector, the VTX mainly measures the impact parameter of tracks. The structure is three double layers with a barrel geometry. Double layer lowers the material budget and improves the impact parameter measurements. The first double layer is half length of the other two to avoid the high occupancy region of direct low omentulum hits from the incoherent pair background.

For the CLIC, the same structure is used. The first layer is moved outwards due to a larger high occupancy region with higher centre-of-mass energy. The detector is also required to provide time stamping at nanoseconds level.

### 3.7.2 Tracking Detectors

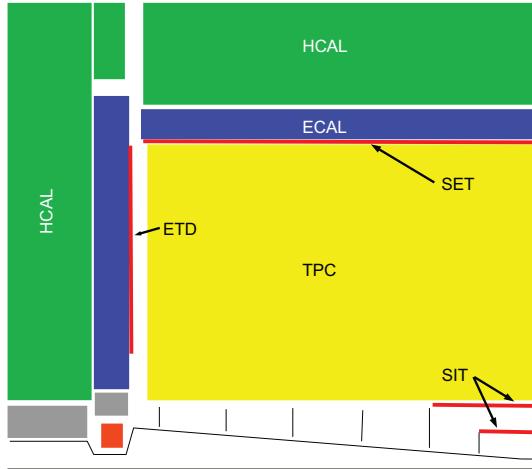
The hybrid tracking system consists of a large volume time projection chamber (TPC), a Silicon Inner Tracker (SIT), a Silicon External Tracker (SET) in the barrel region, a end cap tracking component (ETD) behind the endplate of the TPC, and a silicon forward tracker (ETD) in the forward region. The SIT, SET, and ETD are made up two single-sided strip layers tilted by a small angle. The ETD is a system of two silicon-pixel disks and five silicon-strip disks. The silicon envelope tracking system and the TPC are shown in figure 3.6.

The main part of the tracking system, the TPC, can measure a large number of three dimensional spatial points. Continuous tracking allows precise reconstruction of non-pointing tracks. The TPC is optimised for point resolution and minimum material, as required for the best calorimeter and particle flow performance.

The barrel silicon trackers improve the the overall momentum resolution. They provide additional high precision space points and additional redundancy between the TPC, the VTX, and the calorimeters. The ETD provide the low angle coverage which is not covered by the TPC.

For the CLIC\_ILD, the hybrid structure is used. The outer silicon tracking system is more important at the CLIC to achieve a high momentum resolution at high centre-of-mass energy, as it is challenging using a TPC to sperate two tracks in high energy jets and to identify events in the collection of 312 bunch crossings in 156 ns.

TODO angular acceptance



**Figure 3.6:** A top quadrant view of the ILD silicon envelope system, SIT, SET, ETD, and ECAL as included in MOKKA full simulation, adapted from the figure in [19].

### 3.7.3 Electromagnetic Calorimeter

The Silicon-Tungsten sampling electromagnetic calorimeters in the ILD consist of a nearly cylindrical barrel and two end cap systems, optimised for particle flow. The ECAL measures photon energies and separates photons from other particles. The fine granular ECAL also sits inside the HCAL, which hosts the first part of the hadronic showers and greatly helps to separate hadronic showers.

The particle flow paradigm has a large impact on the ECAL design with many requirements. In addition for the ECAL to measure and separate photons, it also needs to reconstruct detail shower profiles to separate electromagnetic showers from hadronic showers, as approximately 50% of hadronic showers starts in the ECAL. These requirements can be fulfilled with an excellent three dimensional granular ECAL.

From test beam data and simulation studies, a sampling calorimeter with longitudinal and transverse segmentation below one Molière radius and below one radiation length at the front the calorimeter is needed. The most compact design is realised with Tungsten as absorber material and silicon pad diodes as active material. A cross section of the ECAL is shown in figure 3.7. Tungsten is a dense material with a large ratio of interaction length to radiation length. This helps to separate electromagnetic showers from hadronic showers by making electromagnetic showers transversely narrow. Silicon pad size of 5.1 by 5.1 mm cover large areas. They are simple and reliable to operate. The choice of thin silicon layers offers a great spatial resolution at a cost of the energy resolution in favour of the particle flow.

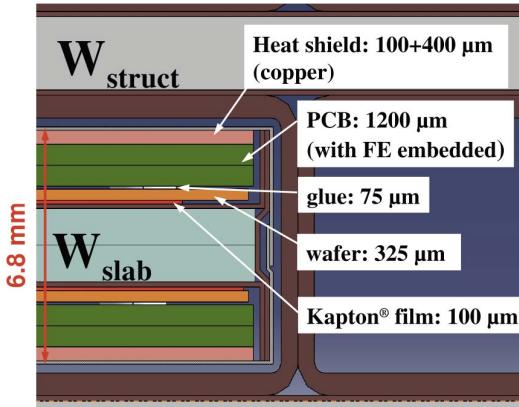
The longitudinal segregation is a compromise between the cost and the performance. The total 30 layers, which is about 20 cm, provides about 24 radiation lengths. The first 20 layers use 2.1 mm thick absorber plates, which is twice finer sampling than the last 10 layers with 4.2mm thick absorber plates. The test beam data with electron shows the energy resolution of the ECAL concept to be  $16.6/\sqrt{E(\text{GeV})} \oplus 1.1\%$ , which is compatible with the values assumed for the full ILD detector simulation.

For the CLIC\\_ILD, the same ECAL from the ILD is assumed, as the requirements of a CLIC detector are satisfied. The increased centre-of-mass energy results in extra energy leakage. But only a small fraction of particles are affected and the leakage is controlled by the HCAL.

### 3.7.4 Hadronic Calorimeter

The requirements of sampling hadronic calorimeter is, again, driven by the need of the particle flow. The need of three dimensional granularity in transverse and logically direction is satisfied by a sampling calorimeter.

The principle role of the HCAL is to separate neutral hadron showers from other particles, and to measure neutral hadron energies. The neutral hadron contribution of the jet energy is around 10% on average. A moderate fine granular HCAL is a good



**Figure 3.7:** A cross section through electromagnetic calorimeter layers, taken from [19].

balance between cost and performance. The chosen layout is 48 longitudinal layers with 3 by 3cm scintillator tiles, using an analogue read out system. The layout of a technological prototype, the "EUDET prototype" is shown in figure 3.8a [25].

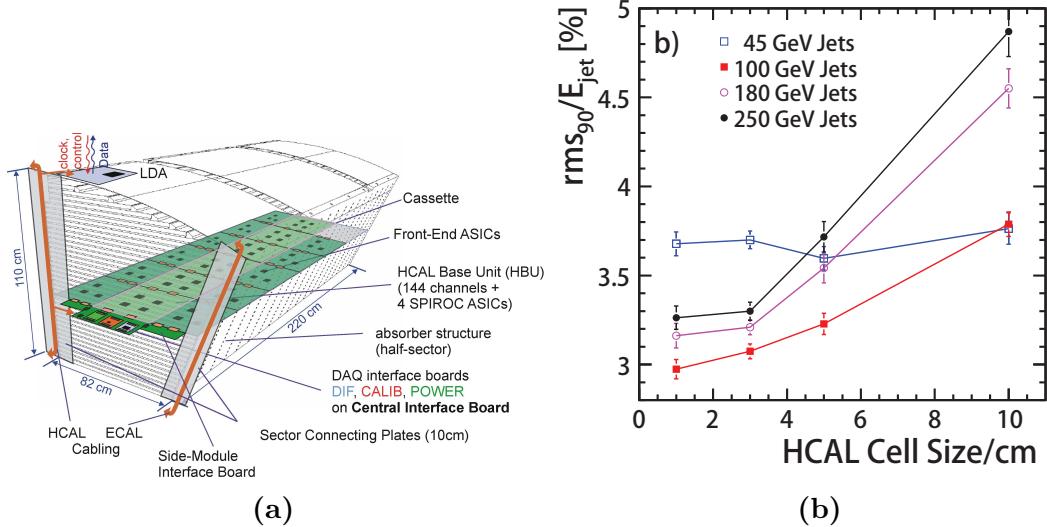
The longitudinal system provide about 6 radiation lengths including the ECAL, which is sufficient to contain the hadronic showers. The transverse cell sizes has been optimised for the best jet energy resolution. It is found that no substantial gain below 3 cm and performance degradation above 3 cm. Hence 3 cm cell size is chosen for the HCAL. The jet energy resolution as a function of HCAL scintillator cell size with different jet energies is shown in figure 3.8b.

The absorber material, stainless steel is chosen for mechanical and calorimetric reasons. Steel allows a self-supporting structure without auxiliary supports. Also steel has a moderate ratio of interaction length to radiation length.

For the CLIC\_ILD, extra layers are added to contain the hadronic shower at high energy. The increased thickener is justified by the simulation studies, where the jet energy resolution degrades quickly for a thinner HCAL.

### 3.7.5 Solenoid

A large superconducting solenoid outside the calorimeters produces a nominal 3.5 T magnetic field. For the CLIC\_ILD, the magnetic field is increased to 4 T for a better performance at a high centre-of-mass energy.



**Figure 3.8:** Figure 3.8a shows the schematic view of a CALICE AHCAL technological prototype module. Figure 3.8b shows the jet energy resolution as a function of the hadronic calorimeter scintillator cell sizes, with different energies. Both figures are taken from [19].

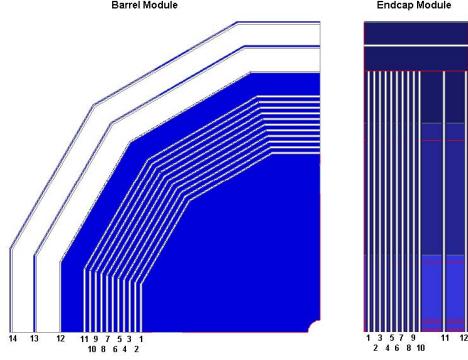
### 3.7.6 Yoke and Muon system

An iron yoke instrumented with scintillator strips active layers returns the magnetic flux, and acts as a muon detector and tail catcher calorimeter at the same time. The layout is shown in figure 3.9. The agreed maximum magnetic field at 15 m radial distance from the detector is 50 Gauss to ensure safety [26]. A highly efficient muon detector is provided by the 3 by 3 cm scintillator strips. As a tail catcher calorimetry, the first layer of the muon detector, catches the energy leakage from the HCAL and the ECAL. It has been shown a 10% improvement of single particle energy resolution is possible with the tail catcher [27].

For the CLIC\\_ILD, due to the different magnetic field strength, the iron yoke thickness differs to that of the ILD.

### 3.7.7 Very Forward Calorimeters

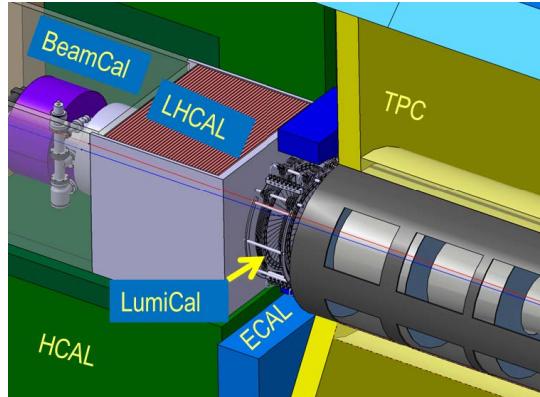
The forward region detectors provide luminosity measurements and forward coverage of calorimeters. A system of precision and radiation resistant calorimeters are required. The luminosity calorimeter counts Bhabha scattering to measure the luminosity to precision



**Figure 3.9:** Sensitive Layers of the ILD muon system, taken from [19].

of  $10^{-3}$  at 500 GeV centre-of-mass energy. The beam calorimeter (BeamCAL) extend the forward coverage, which are hit by many beamstrahlung pairs after each bunch crossing. BeamCAL estimates a bunch-by-bunch luminosity. An additional hadron calorimeter, LHCAL, at the forward region extends the angular coverage of the HCAL to that of the LumiCAL. Electron tagging is possible with the very forward calorimeters [28], which aids event reconstruction at high centre-of-mass energy.

The CLIC\_ILD adopted a similar very forward calorimetry system as that of the ILD. Modifications to the design due to CLIC 3 TeV centre-mass-of energy can be found in [18].



**Figure 3.10:** The forward calorimeters of the ILD, taken from [19]. LumiCal, BeamCal, and LHCAL are the luminosity calorimeter, beam calorimeter, and forward hadronic calorimeter.

### 3.7.8 ILD vs CLIC\_I LD

Beam Calorimeter acceptance is defined as  $|\cos(\theta_Z)|$  is between 0.01 and 0.04 rad and length in z direction is between 3181 and 3441 mm. Luminosity Calorimeter acceptance is defined as  $|\cos(\theta_Z)|$  is between 0.038 and 0.11 rad and length in z direction is between 2539 and 2714 mm.

ggHad

incoheraent pair beamsstrulung bunch crossing



# Chapter 4

## Simulation and Reconstruction

*“How to open a pandora box?”*

— A wise Chinese

In previous chapters, overviews of the theory and the future linear collider experiments have been described. Since the work presented in this document is for future collider, simulation and monte carlo method is used throughout the document. In this chapter, the simulation and event reconstruction chain are discussed, with emphasis on the PandoraPFA event reconstruction, which provides the background for the photon reconstruction algorithms in chapter 5.

Simulation and reconstruction of events for the future Linear Colliders, ILC and CLIC, share common software framework. The noticeable difference will be discussed.

### 4.1 Monte Carlo event generation

For the simulated study, Monte Carlo (MC) event generation is the first step. Most events, electron-positron interaction, are generated with WHIZARD software, [29, 30], with no polarisation of the electron and positrons. Some simple events are generated with HEPEVT. PYTHIA [31] is used to describe parton showering, hadronisation and fragmentation. The parameters for PYTHIA are tuned to OPAL data from the LEP [32]. TAUOLA [33] debrides the tau lepton decay with correct spin correlations of the decay products. The Initial State Radiation (ISR) is simulated in WHIZARD with the ISR

photons being collinear with the beam direction. The Final State Radiation (FSR) is simulated with default parameters in PYTHIA.

For the CLIC simulated samples, the luminosity spectrum, which is generated with GUINEAPIG [34], is simulated in WHIZARD. The  $\gamma\gamma \rightarrow \text{hadrons}$  background events are hadronised with PYTHIA, and superimposed on the physics process simulations to save computational resources.

TODO  $\gamma\gamma \rightarrow \text{hadrons}$  simulation -5 to 25ns etc

## 4.2 Event Simulation

The event simulation software is GEANT4 [35], and the detector geometry description is provided by MOKKA [36]. QGSP\_BERT physics list is used to describe the hadronic showers decay in the detector.

## 4.3 Event Reconstruction

Reconstruction software runs in Marlin framework [37], as a part of the iLCSoft. Event reconstruction contains following steps: digitisation of simulated calorimeter hits, reconstruction of tracks in the tracking system using pattern recognition algorithms, and particle flow objects (PFOs)reconstruction with PandoraPFA [21, 22]. Reconstruction does not include calorimeters hits in the forward calorimeters, due to computational reasons (see section 7.3.3).

For the CLIC detector simulations, suppression of  $\gamma\gamma \rightarrow \text{hadrons}$  background is performed.

Details of the reconstruction can be found in [17, 18]. Particle flow reconstruction via PandoraPFA will be discussed in details, which provides the background for the photon reconstructions in PandoraPFA in chapter 5.

## 4.4 PandoraPFA

Tradition calorimetric approach is unable to meet the mass and energy requirements for the future linear collider. The particle flow approach with PandoraPFA has a proof-of-principle demonstration of its capability to reach required resolution. The particle flow approach also put stringent requirements on the detector design, which is described in chapter 3. By associating calorimeter hits to the tracks, around 60% of the jet energy from charged particles is measured by the tracker, which has a much better resolution than the calorimeter. Small cell sizes of the calorimeters are required to identify hits from different particles. The traditional sum of calorimeter cell energies is replaced by particle flow reconstruction algorithms, a complex pattern recognition problem. The PandoraPFA algorithm has been developed and used in the ILC and CLIC simulation studies.

Started with the ILD detector concept, PandoraPFA has been adapted to the CLIC condition and shows its ability to deliver required energy resolutions [18]. Recent the code base of the PandoraPFA has been restructured. The core base codes for basic object and memory managements are factorised in the Pandora C++ Software Development Kit [38]. There are over 60 linear collider specific reconstruction algorithms, each aims to address a particular topological in the reconstruction.

In the subsequent paragraphs, the main steps in the PandoraPFA reconstructions are described. The details of the reconstruction can be found in citeThomson:2009rp,Marshall:2012ry,Marshall

Inputs of PandoraPFA are digitised calorimeter hits and reconstructed tracks. The output are reconstructed particles with four-momenta, Particle Flow Objects (PFOs).

### 4.4.1 Track selection

Tracks from the tracking system are selected based on their topological properties, how likely they are from physical processes, and whether they are consistent with the tracker resolution. Tracks passed the selection are used for the subsequent reconstruction. Special topologies of tracks are identified, such as when a neutral particle decays or converts into a pair of charged tracks, leaving a “V0” shape tracks. This is identified by searching for a pair of tracks originated from a single point. Another topologies include “kinks”, when a charged particle decays to a single charged particles with neutral particles, and “prongs”, when a charged particles decays to multiple charged particles. This information

are stored and passed on to the subsequent reconstruction, along side with the helical track fit (using last 50 reconstructed hits) and the track projection to the front of the ECAL.

#### 4.4.2 Calorimeter selection

The input information of a calorimeter hit is the position, the layer in the calorimeter and the energy response from the calorimeter.

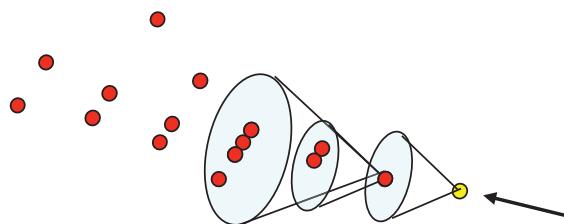
Calorimeter hits are selected based on a series of criterion. The selected hits need to have energies above the threshold, using the conversion of a minimum ionising particle (MIP) equivalent, and using directly the converted energy. Similar to tracks, only calorimeter passed the selection are used in later steps.

Geometry information and likelihood of the hit originated from a minimum ionising particle (MIP) are calculated.

Isolated hits, often originated from low energy neutrons in a hadronic shower, are difficult to associate to the correct hadronic shower. They are identified and not used in the clustering. But their energy is added in the very last particle flow object (PFO) creation step.

#### 4.4.3 Cone Clusters Algorithm

Before discussing the rest of the PandoraPFA reconstruction, it is necessary to introduce the cone based clustering algorithm, which is widely used in the calorimeter in PandoraPFA. The clustering algorithm produces basic working objects, Clusters.



**Figure 4.1:** Illustration of the cone clustering algorithm, taken from [39]

There are two main types of clustering algorithms: cone based and sequential combination (see section ??). The main clustering scheme PandoraPFA is cone clustering, for grouping calorimeter hits. Illustrated in figure 4.1, cone clustering has a specified opening angle of the seed hit. Because the direction of particle flows is largely unchanged from the originated particle, whether it is a electromagnetic shower, QCD radiation or hadronisation, these cone clusters have similar direction and energy to the originated particle. Therefore it is applicable to use cone based clustering algorithms for building clusters.

The seed for the cone clustering is typically the projection of a energetic track to the front of the ECAL. A high energy calorimeter hit can also be used as a seed. A cone with a specified opening angle and depth will be formed around the seed. The four-momentum of calorimeter hits sum to the cone's four-momentum.

TODO Build from inner to outer, then every layer outer in inner see Mark's paper

#### 4.4.4 Particle Identification

Dedicated particle identification algorithms aim to identify muons and photons before associating calorimeter hits to tracks. The details of the photon reconstruction algorithms and photon related algorithms are described in chapter 5. By removing the hits from muons and photons, the reconstruction of charged particles is improved as it reduce the pattern recognition problem. Identified muons and photons do not participate in the clustering and re-clustering stages, but re-entre the construction at the fragment removal stage (see section 4.4.9).

#### 4.4.5 Clustering

The cone clustering algorithm described in section 4.4.3 is used to group calorimeter hits from innermost to outmost psuedo-layer. The output Clusters are further processed, merged or split based on their topological properties.

#### 4.4.6 Topological cluster association

Initial clustering scheme is aggressive at splitting clusters. Small clusters are merged based on clear topological signatures. These merging signatures include combining track

segments, connecting tack segments with gaps, connecting track segment to a hadronic shower, and merging clusters when they are within close proximity.

#### 4.4.7 Track-cluster association

Clusters are associated to tracks, according to the proximity of the first layer of the cluster and the track projection to the front of the ECAL.

#### 4.4.8 Re-clustering

The cluster association scheme work well for low energy (less than 50 GeV) jet. For a high energy jet, particles and the subsequent hadronic showers are more boosted and more likely to overlap each other. Therefore, it is important to re-cluster based on the compatibility of the cluster energy and the associated track momentum. A cluster may be split into two. Two clusters maybe be re-clustered based on the track-cluster association. The re-clustering algorithm is applied iteratively to find a more correct clustering of calorimeter hits.

#### 4.4.9 Fragment removal

The late stage of the reconstruction will focus on merging low energy clusters, especially non-photon neutral clusters. These neutral clusters are likely to be fragments of charged clusters, instead of being a physical particle. The merging criterion are mostly based on the proximity and the energy comparison.

One algorithm will attempt to split up photon clusters, where each is originated from two close by photons. Photon related algorithms are described in details in chapter 5.

#### 4.4.10 Particle Flow Object Creation

Particle Flow Objects (PFOs) are created at the last step. Tracks are associated to the clusters based on the proximity. Simple but effective particle identification for electrons, muons are applied. Photon identifications have been applied at various stages of the reconstruction.

PFOs are the output of the PandoraPFA reconstruction. The four-momentum of these PFOs are used heavily for the downstream analysis. The electron, muon and photon identification are also used in physics analysis, such as one described in chapter 6.

## 4.5 MC truth linker

Hits contribution to MC particle. Main contributed MC particle from energy weighted hits Main contributed Jets etc. from hits

## 4.6 CLIC specific simulation and reconstruction

## 4.7 Luminosity spectrum

## 4.8 Suppression of $\gamma\gamma \rightarrow \text{hadrons}$ backgrounds

For the CLIC, significant  $\gamma\gamma \rightarrow \text{hadrons}$  background is present. It is crucial to remove the beam induced background as they don't represent the underlying physics process.

Two Marlin process has been developed to suppress these background, a track selector and a PFO selector [22].

The track selector aims to remove poor quality and fake tracks. It places simple quality cut and a simple time of arrival cut. If the arrival time of the track at the front of the ECAL, using the helical fit, differs more than 50 ns from using a straight line fit, the track will be rejected.

The PFO selector utilise the high spatial resolution from the high granular calorimeter. PFOs from  $\gamma\gamma \rightarrow \text{hadrons}$  often have low  $p_T$  and have a range of time. PFOs from physics processes have a range of  $p_T$ , and have time close to the brunch crossing time. These two distinctive features allow  $\gamma\gamma \rightarrow \text{hadrons}$  background to be separated. The optimal suppression uses different  $p_T$  and time cuts for the central part of the detector, and for the forward part of the detector, and uses different cuts for photons, neutral PFOs and charged PFOs. Three configurations of these cuts are developed, namely “loose”, “normal”, and “tight” selections. As the name suggested, “loose” selection corresponds

to a looser cut of  $p_T$  and time. The optimal configuration depends on the  $\sqrt{s}$  of the collision, and the physics process to study.

The background suppression is used in analysis described in chapter 7 TODO add efficiency table pandora TODO add timing pT cut and table of energies

## 4.9 CLIC simulated particle masses

## 4.10 Reconstruction Processors

Automated analysis is the only way to deal with the vast amount of data generated in the high energy physics. In the last chapter we described the automated reconstruction tools in details. This chapter is dedicated to the common automated analysis tools and techniques, which will be used in the analysis described in subsequent chapters.

For the linear collider, thanks to the high granular calorimeter, the starting point for analysis would be individual Particle Flow Objects, as well as individual tracks. Each of the PFOs encodes four-momentum and position information. For tracks, they would have momentum and position information.

However, sometimes it is interesting to group PFOs and tracks into jets, where a jet is the result of hadronisation process from high energy particles like quarks or gluons.

## 4.11 Jet algorithm

A jet is typically a visually obvious structure in a event display. The momentum and the direction of a jet tend to resemble the originated particle. Despite the relative easiness of identifying jets visually, it presents a challenge for a pattern recognition program to identify jets effectively and efficiently.

Early work on jet finding started in 1977 [50], where later development can be found in reviews [51–53].

There are two large families of jet finding algorithm, cone based algorithms, and sequential combination algorithms. Cone based algorithm is briefly discussed in section ?? in the context of the PandoraPFA user case.

Sequential combination algorithms typically calculate a pair-wise distance metric. Pairs with the smallest metric will be combined. The metric will be calculated and updated after a combination. This procedure will be repeated until some stopping criterion are satisfied.

The chosen jet algorithm implementation is FastJet C++ software package [54, 55], providing a wide range of jet finding algorithms. The implementation in Marlin software package is called MarlinFastJet. The symbols in the subsequent discussion follow the convention in [54]

**$k_t$  algorithm** longitudinally-invariant  $k_t$  algorithm [56, 57] is one of the common sequential combination algorithms for  $\bar{p}p$  collider experiment. In the inclusive variant, the symmetrical pair-wise distance metric between particle  $i$  and  $j$ , and the beam distance, are defined as

$$d_{ij} = d_{ji} = \min(p_{Ti}^2, p_{Tj}^2) \frac{\Delta R_{ij}^2}{R^2}, \quad (4.1)$$

$$d_{iB} = p_{Ti}^2, \quad (4.2)$$

where  $p_{Ti}$  is the transverse momentum of particle  $i$  with respect to the beam ( $z$ ) direction, and  $\Delta R_{ij}^2$  is the measurement of angular separation of particle  $i$  and  $j$ , defined as  $\Delta R_{ij}^2 = (y_i - y_j)^2 + (\phi_i - \phi_j)^2$ , where  $y_i = \frac{1}{2} \ln \frac{E_i + p_{zi}}{E_i - p_{zi}}$  and  $\phi_i$  are particle  $i$ 's rapidity and azimuthal angle.  $R$  is a free parameter controlling the jet radius.

If  $d_{ij} < d_{iB}$ , particle  $i$  and  $j$  are merged, with the four-momentum of particle  $i$  updated as the sum of the two particles. Otherwise, particle  $i$  is set to be a final jet, and deleted from the particle list. The above procedure is repeated until no particle left.

The exclusive variant is similar. First difference is that when  $d_{iB} < d_{ij}$ , the particle  $i$  is discarded and part of the beam jet. The second difference is that when both  $d_{ij}$  and  $d_{iB}$  are above some threshold,  $d_{cut}$ , the clustering will stop. In practise, exclusive mode allows a specified number of jets to be found, which will automatically choose the  $d_{cut}$ . The inclusive mode would find as many jets as the algorithm allows.

### 4.11.1 Durham algorithm

Durham algorithm [?], also known as  $e^+e^- k_t$  algorithm, is commonly used  $e^+e^-$  collider experiment. It has a single distance metric:

$$d_{ij} = 2 \min(E_i^2, E_j^2)(1 - \cos(\theta_{ij})), \quad (4.3)$$

where  $E_i$  is the energy of particle  $i$ .  $\theta_{ij}$  is the polar angle difference between particle  $i$  and  $j$ . Durham algorithm can only be run at exclusive mode, which means that the clustering will stop when  $d_{ij}$  is above some threshold,  $d_{cut}$ .

Comparing to  $k_t$  algorithm, it uses energy instead of  $p_T$  in the distance metric, and it did not have a beam jet. This is because that for the  $e^+e^-$  collider in the past, the beam induced background was not severe and collisions energy is known,  $\sqrt{s}$ .

**Jet algorithm for the CLIC** Although CLIC is a  $e^+e^-$  collider, the significant beam-induced background adds a large amount of energy from  $\gamma\gamma \rightarrow \text{hadrons}$  process (see section ??). Therefore, traditional  $e^+e^-$  jet algorithms, like Durham algorithm (see section ??), is not suitable for the CLIC environment. Studies has shown that jet algorithms for  $\bar{p}p$  collider have better performance [18, 58].

A more recent attempt at marrying merits from both Durham and  $k_t$  algorithms has resulted in Valencia jet algorithm [59]. It had shown promising improvement comparing to  $k_t$  algorithm, which is used in the parallel  $HH \rightarrow b\bar{b}b\bar{b}$  sub-channel analysis.

### 4.11.2 $y$ parameter

A commonly used variable for number of jets is the  $y$  parameter.  $y$  parameter describes the transition of exclusive jet algorithm going from  $N$  clustered jets to  $N+1$  clustered jets. For example,  $y_{23}$  would be the  $d_{cut}$  value for a exclusive jet algorithm, above which the jet algorithm returns 2 jets, below which the jet algorithm returns 3 jets (see section ?? for jet algorithm). Numerically  $y$  parameter is often much smaller than one. A typically way to convert the small number to a human acceptable range is to take the minus logarithm of the number.

### 4.11.3 LCFIPlus

The LCFIPlus software package is based on the LCFIVertex package, which was used in the simulation studies for ILC Letter of Intent [61, 62] and CLIC Concept Design Report [18]. Current software is built in mind of a future  $e^+e^-$  collider. Although the software is modular and can be used in any order, here it will be described in the order used in a physics analysis.

The input are PFOs. The vertex finding algorithms perform vertex fitting and identify primary and secondary vertex. There is a “V0” particle rejection step, which is neutral particles decaying or converting into a pair of charged tracks. The topology is similar to the decay of b or c hadrons. Hence it is important to remove the V0 particles to improve the heavy quark flavour tagging (see section 4.4.1 for a similar V0 rejection).

Once the primary and secondary vertices are found, PFOs are clustered in to jets. This jet clustering scheme ensures that the secondary vertices and the muons identified from semileptonic decay fall in the same jet. Therefore, it is consistent with the hadronic decay. Jet algorithms used are Durham and Durham modified algorithms(see section ??).

The next step is to refine vertices finding to improve the b jet identification from c jet. Since the existence of two close by vertices is strongly correlated to a b jet, the vertices refining step will reconstruct as many secondary vertices correctly as possible.

The last step is to gather the information about vertices and jets, and deploy a multivariate analysis. The multivariate classifier used, Boosted Decision Tree, is implemented in TMVA software package [63], which is discussed later in section ?? . A series of flavour sensitive variables are calculated, and the classification is divided into four subset: jet with zero, one, or two properly reconstructed vertices, or a single-track pseudovertex. For each subset, a jet can either be classified to a b jet, a c jet, or a light flavour quark jet (u, d or s). The multiclass classifier’s response is normalised across different subset, and they will be referred in the subsequent physics analysis as the tag value. See section ?? for a discussion on multiclass classifier.

The samples for training the multiclass classifier are  $e^+e^- \rightarrow Z\bar{v}v$  at  $\sqrt{s} = 1.4$  TeV, where Z decays to  $b\bar{b}$ ,  $c\bar{c}$ , or  $u\bar{u}/d\bar{d}/s\bar{s}$ .

The flavour tagging is performed after the initial jet reconstruction, and all the PFOs in the reconstructed jets are the input to the LCFIPlus flavour tagging processor.

Therefore, the classifier in the LCFIPlus processor is trained for a specific PFO collection and a specific jet reconstruction algorithm. In this analysis, the classifier is trained with the optimal jet reinstruction choice, discussed in section 7.4.1. The output of the processor for a jet is three values, corresponding to the likelihood of the jet being a b jet, a c jet, or a light flavour quark jet. The selection efficiency of b-jets and c-jets with training samples is shown in figure 7.5.

## 4.12 MVA

Multivariate analysis has become increasingly common in high energy physics. MVA can be viewed as an advanced tool for regression or classification. Comparing to the traditional cut based method, modern machine learning technique offers much improvement in data analysis. Software package for MVA used throughout this document is TMVA [63].

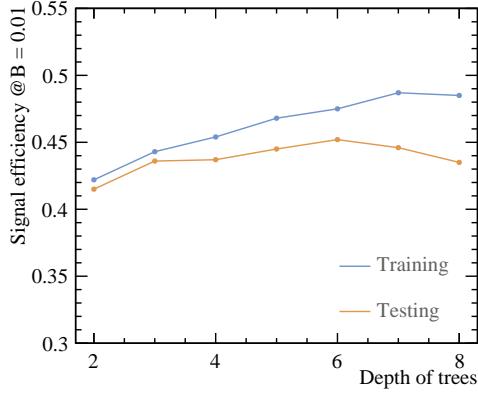
A typical machine learning MVA classification involves two classes, also known as signal and background. A machine learning model, also known as a classifier in TMVA, needs to be trained with training data. The model requires a set of discriminative variables, which separate the signal from background. The trained model will be applied onto the testing data for signal extraction. Response of the model could be a classification of signal or background, or could a response in a continuous spectrum, where the user decides the value to separate signal from background.

Strictly, there should be three statistically independent samples for the MVA. One sample is for the training. Another sample for the validation, including optimisation and checking for overfitting. The last sample is for testing. However, due to technical reason (TMVA only natively supports two samples), sometimes the same sample is used for the validation and the testing, which is acceptable with large statistics.

This classification scheme can be easily extended to multiple classes, implemented in TMVA with multiclass class. The multiclass class is used in the tau decay mode classification in section ?? and in the flavour tagging classifier in section 7.5.

### 4.12.1 Optimisation and overfitting

The optimisation of the model refers to selecting the optimal free parameters of the model. One could build a complex model which fits the training samples very well, but



**Figure 4.2:** Example of MVA overtraining

it would not be optimal for another testing sample. A simple model is less prone to statistical fluctuation of samples, however, it might be too simple to achieve the optimal modeling. The former case is known as overfitting, or overtraining. The latter case is called underfitting, or undertraining.

The compromise is clear. The optimal model is one between overfitting and underfitting. In practice, this involves building the model with increasing complexity, and finding the point where overfitting occurs.

Figure 4.2 shows a typical overfitting plot. Overfitting is defined when the efficiency of signal selection in the training samples increases, but the efficiency in the testing sample decreases. Here the example is chosen from double Higgs analysis at  $\sqrt{s} = 3 \text{ TeV}$ , using Boosted Decision Tree model. The efficiency of signal selection is defined as the signal fraction when background fraction is 1%, report by the TMVA training process. In figure 4.2 , the depth of the tree, or the number of layers in the tree, reflects the complexity of the model. From tree depth 2 to 5, the efficiency for both testing and training samples increases. From tree depth 6 onwards, the overfitting occurs. In this particular example, one should choose a tree depth fewer than 7 to avoid overfitting.

There are methods to assign the error on the selection efficiency. Thus one can make a better choice of parameters to avoid overfitting. These methods were not implemented due to the technical capacity provided by the TMVA.

### 4.12.2 Choice of models

The model, also known as the classifier in TMVA, can be as simple as cut based, likelihood or linear regression. It can also be as complicated as non linear tree, non linear neutral network or support vector machine. Regardless of model complexity, the choice of most optimal classifier is often data driven. Also, given the free parameters in each model, the comparison between different models without individual tuning is not rigorous. Nevertheless, as researchers in the machine learning suggested, the boosted decision tree is probably the best out-of-the-box machine learning method. Neutral network could potentially be better than the boosted decision, but it requires more tuning, and it is less intuitive to interpret the model. For these reasons, boost decision tree (BDT) is often the choice of machine learning model in the high energy physics. And it is used in various physics analysis in this document.

Before describing BDT in detail, we will first visit some simple models.

#### Rectangular Cut

Probably the most intuitive model, the rectangular cut method optimise cuts to maximise some specific metric. The metric could be the signal efficiency for a particular background efficiency. Alternatively, the metric can be the significance,  $\frac{S}{\sqrt{S+B}}$ , where S and B are signal and background numbers, respectively.

Discriminative variables gives better separation power when they are gaussian-like and statistically independent. Therefore it is common to decoorelate the variables and gaussian transform them before using the rectangular cut MVA.

Because its simplicity, the cut method is often performed manually, much more often in the time pre-date the wide spread of machine learning methods. It is still commonly used for the pre-selection step before the MVA (see section 7.7), and other simple cases. Unless specified, the optimal cuts proposed in this document for various physics analysis are found using the rectangular cut method manually.

## Projective Likelihood

Projective likelihood model (PDE) is used in PandoraPFA for the photon ID due to its simplicity and low requirement on computing resources. It is discussed in details in section 5.5.1

## Boosted decision tree

Boost decision tree (BDT) is a non linear tree based model. Its rather complex nature requires a careful explanation of many concepts within the BDT.

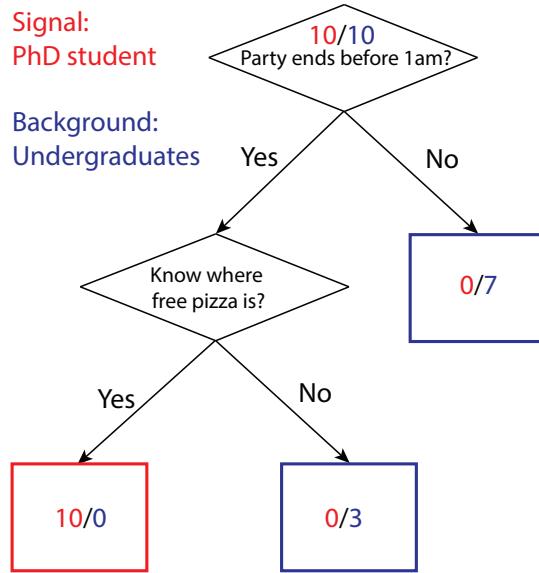
**Decision tree** Decision tree is a binary tree, where each node, the splitting point, uses a single discriminative variable to decide whether a event is signal-like (“goes down by a layer to the left”), or background-like (“goes down by a layer to the right”). At each node, samples are divided into signal-like and background-like sub-samples. The tree growing starts at the root node, and stops at certain criterion, which could be the minimum number of events in a node, the number of layers of the tree, or a minimum/maximum signal purity.

The training of the decision tree is to determine the optimal cut at the node by minimising the metric. The probability of the cut producing the signal is  $p$ . Three commonly used metrics for two-class classification are

1. Misclassification error:  $1 - \max(p, 1-p)$ ,
2. Gini index:  $2p(1-p)$ ,
3. Cross-Entropy or deviance:  $-p \log p - (1-p) \log (1-p)$ .

The using of a trained decision tree is to transverse along the tree. The event is classified as signal or background depending on whether it falls in the signal-like or background-like end node.

Figure 4.3 illustrate a simple example of a decision tree. The signal is the PhD student and the background is the undergraduate student. The depth of this imbalance binary tree is 2. A node is represented by a diamond. The signal-like node is the red rectangle and the background-like nodes are blue rectangles. The tree is constructed with two possible cut, “Party ends before 1am” and “Know where free pizza is”. The attribute of samples is listed in table 4.1. To demonstrate the choice of the first layer cut, the



**Figure 4.3:** Example of a decision tree

Number	Party ends before 1am	Know where free pizza is
PhD student	10	10
Undergrad student	3	5

**Table 4.1:** The attribute of samples for the decision tree example.

Gini index metric is used. If the first cut is “Party ends before 1am”, the probability of the cut producing the signal,  $p$ , is  $\frac{10}{13}$ . Gini index is  $2p(1-p) \simeq 0.36$ . If the first cut is “Know where free pizza is”,  $p = \frac{10}{15}$ . Gini index is  $2p(1-p) \simeq 0.44$ . Therefore, the first cut is “Party ends before 1am”.

The simple tree in figure 4.3 is grown fully as each end node contains signal or background only. To use the trained decision tree, if there is student who ends party before 1am and knows where free pizza is, then the student is classified as a PhD student.

**Improve decision tree** Decision tree has a low bias, but high variance. This means it is very easy to construct a tree that fits the training data very well, but the tree would not be optimal for the testing sample. To overcome the instability of the decision tree, many methods have been developed. The most successful one is boosting.

Boosting: it is a technique where the misclassified events receives a higher weight than the correctly classified events. Therefore, when the training is iterated, the misclassified

events would receive higher and higher weights and more likely to classify correctly. The boosting is done at every iteration, which can be few hundred or few thousand times. This will create a “forest” of many trees. The final output could be a majority vote, by traversing the event to the end node for each tree in the forest.

Bagging: also known as boot-strap, it is a method that selects a simple random sub-set of the training sample, and apply the model. In this case, every boosting iteration takes a bagged sample, rather than the whole sample.

TMVA implementation of the BDT for the output is using a likelihood estimator, depending on how often an event is classified as signal in the forest. The likelihood number is later used to select signal from background.

### 4.12.3 Optimisation of Boosted Decision Tree

Many parameters of the BDT can be tuned. The tuned parameters are described below. The optimal values are obtained by choosing the best performance without overfitting with samples at  $\sqrt{s} = 3 \text{ TeV}$ . The same values are used  $\sqrt{s} = 1.4 \text{ analysis TeV}$ .

The most important parameter is the depth of a tree, which determines how many end nodes a tree has, or the degrees of freedom of a tree. The related parameter is the number of trees. Experience shows that using many small trees yields the best result. The performance as a function of the depth is shown in figure 4.2. The chosen value is 4.

The number of trees is another important parameter. Intuitively large number of trees leads to overfitting. However, it has been shown that a large number does not lead to overfitting, using the definition above. There is a debate on how to determine the optimal number of trees. The chosen value is 4000.

The minimum number of events in a node, which is a stopping criteria for tree growing, affects the size of the tree. But it is less influential than the depth of the tree. The chosen value is 0.25% of the total events.

The boosting has two variants, adaptive boost and gradient boost. For all the BDT used in this document, adaptive boost is used.

The learning rate of the adaptive boost, which controls how fast the weight changes for events in each boosting iteration. Experience shows small learning rate with many trees work better than large learning rate with few trees. The chosen value is 0.5.

The usual choice of the metric for the optimal cuts is either Gini index or cross-entropy. Gini index metric is chosen. It makes little difference to performances, comparing to the cross-entropy metric.

Number of bins per variables for the cut is necessary to make tree growing efficient. Discrete binned variables are faster to computer than continuous variables. The parameter does not impact the performance much. However, variables should be pre-processed before going into the model. For example, the variable should be limited to a sensible range to avoid the extremes. The variable should also be transformed to obtain a more uniform distribution, if the original distribution is highly skewed. The chosen value is 40.

For the end node, it is determined as either signal-like or background-like, based on the majority for the training event in the end node. Numerically, it corresponds to 1/0. However, the end node could also use signal purity as the output, resulting in a continues spectrum of [0,1]. The chosen method is the continuous response.

#### 4.12.4 Multiple classes

The above discussion is done assuming two classes - signal and background. The argument can be easily extended to multiple classes. There are two ways for the training. "One v.s. one" is each class is trained against each other class. And the overall likelihood is normalised. The second way to train is called "one v.s. all", which is when each class is trained against all other classes.

Using a three-class example, A, B and C, "one v.s. one" scheme trains A against B, B against C, and C against A. Then the likelihood is normalised. "One v.s. all" would train A against B plus C, B against A plus C, and C against A plus B.

TMVA multiclass implementation uses "one v.s. all" scheme. Multiclass is used in falvour tagging of jets, section ??, and in the tau lepton final state separation study, section ??.

### 4.13 Event shape variables

Event shape variables are some useful global variables to describe the shape of the event, for example whether it is back-to-back, or homogenous in the solid angle.

The classical event shape thrust [42], is defined as

$$T = \max_{\hat{t}} \frac{\sum_i |\hat{t} \cdot \vec{p}_i|}{\sum_i |\vec{p}_i|} \quad (4.4)$$

where  $\vec{p}_i$  is the momentum vector of the particle  $i$ . Summation is over all particles in the event. Thrust axis,  $\hat{t}$ , is a unit vector. (Principle) Thrust value,  $T$ , is 1 for a perfect pencil-like back-to-back two-jet event, and 0.5 for a perfect spherical event. The thrust value is useful in picking out back-to-back two-jet event. Thrust axis is useful to separate each jet in a back-to-back two-jet event.

It is derived from the sphericity tensor [64], defined as

$$\mathbf{S}^{\alpha\beta} = \frac{\sum_i p_i^\alpha p_i^\beta}{\sum_i |\vec{p}_i|^2}, \quad (4.5)$$

where  $\vec{p}_i$  is the momentum vector of the particle  $i$ . Summation is over all particles in the event.  $\alpha$  and  $\beta$  refer to the x, y, z coordinate axis. Eigenvalues of tensor  $\mathbf{S}$  can be found, or in this case diagonalisation of the matrix  $\mathbf{S}$ , denoted with  $\lambda_1, \lambda_2, \lambda_3$ . The normalisation condition requires  $\lambda_1 \geq \lambda_2 \geq \lambda_3$  and  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ . Sphericity,  $S$ , is defined in terms of  $\lambda$ ,

$$S = \frac{3}{2}(\lambda_1 + \lambda_2). \quad (4.6)$$

$S$ , is 0 for a perfect pencil-like back-to-back two-jet event, and 1 for a perfect spherically symmetric event.

Sphericity tensor [64], is defined as

$$\mathbf{S}^{\alpha\beta} = \frac{\sum_i p_i^\alpha p_i^\beta}{\sum_i |\vec{p}_i|^2}, \quad (4.7)$$

where  $\vec{p}_i$  is the momentum vector of the particle  $i$ . Summation is over all particles in the event.  $\alpha$  and  $\beta$  refer to the x, y, z coordinate axis. Eigenvalues of tensor  $\mathbf{S}$  can be found, or in this case diagonalisation of the matrix  $\mathbf{S}$ , denoted with  $\lambda_1, \lambda_2, \lambda_3$ . The normalisation condition requires  $\lambda_1 \geq \lambda_2 \geq \lambda_3$  and  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ . Sphericity,  $S$ , is defined in terms of  $\lambda$ ,

$$S = \frac{3}{2}(\lambda_1 + \lambda_2). \quad (4.8)$$

Event Number	True Signal	True Background
Selected Signal	$N_S$	$N_1$
Selected Background	$N_2$	$N_B$

**Table 4.2:** A toy example to demonstrate definitions of efficiency and purity.

$S$ , is 0 for a perfect pencillike back-to-back two-jet event, and 1 for a perfect spherically symmetric event.

Aplanarity is another event shape variable that distinguishes spherical symmetrical events from planar and linear events. The definition is

$$S = \frac{3}{2}(\lambda_1), \quad (4.9)$$

where  $\lambda_1$  is the largest eigenvalue in the diagonalised sphericity tensor.

## 4.14 Miscellaneous

An event in a collider experiment refers to one collision and the subsequent energy deposition in the detector. An event corresponds to a certain type of physics process.

Often we are dealing with extracting a type of events, from a large number of other events. The signal, or signal events refer to events of interests. Other events are referred to as the background, or background events.

Typical metrics of signal selection is efficiency and purity. This toy example illustrates definitions of efficiency and purity.

Signal selection efficiency is defined as  $\frac{N_S}{N_S+N_2}$ . Signal selection purity is defined as  $\frac{N_S}{N_S+N_1}$ . Significance is a quantity that is similar to purity,  $\frac{N_S}{\sqrt{N_S+N_1}}$ .

When we are describing particles, light lepton,  $l_l$ , refer to electrons,  $e^-$ , and muons,  $\mu^-$ . Light quarks,  $q_l$ , refer to up quark, u, down quark, d, and strange quark, s.

Computational intensive jobs are processed either on the Cambridge High Energy Physics grid, or the CLIC computing grid.

# Chapter 5

## Photon Reconstruction in PandoraPFA

*“Photons have mass? I didn’t even know they were Catholic.”*

— Woody Allen

Photon reconstruction is an important part of particle reconstruction. A good photon reconstruction provides a good single photon completeness and purity, as well as a good photon separation resolution. Such a good photon reconstruction is crucial for reconstructing heavy particles, for many physics processes involving these particles decaying into photons, such as  $\tau$  lepton and  $\pi^0$ .

### 5.1 Electromagnetic shower

### 5.2 Overview of photon reconstruction in PandoraPFA

PandoraPFA provides a framework for particle reconstruction [], as described in chapter 5. In the linear collider content, it has a vast library of algorithms developed through years by many people. Each algorithm addresses one topological issue in the particle reconstruction []. The essential part of the PandoraPFA is track-cluster association and reclustering to find the best track-cluster pair. Algorithms that removes trackless clusters,

such as removing muon clusters or photon clusters, would provide a clean environment for the track-cluster association, hence improving the jet energy resolution.

Photon identification in the PandoraPFA has two main mechanisms. The basic mechanism (see section 4.4.10) performs photon identification after track-cluster association and the reclustering processes. The second more sophisticated photon identification, the photon reconstruction algorithm, is performed before the track-cluster association and reclustering process. This algorithm identifies photon electromagnetic shower cores carefully in the dense jet environment.

The photon reconstruction algorithm in PandoraPFA version 1 improves jet energy resolution by correctly identifying photon electromagnetic shower cores and leaving a cleaner environment for the track-cluster association. However, the peripheral calorimeter hits to the shower cores may be left as fragments, and reconstructed as separate particles. This lowers the reconstructed photon completeness and makes the number of reconstructed photons a less useful physical quantity. Also, the algorithm in PandoraPFA version 1 leaves rooms for improvement of photon separation resolution.

This section presents a solution to the photon fragments issue. The newly introduced PandoraPFA algorithms also improves the photon separation resolution. Algorithms related to photon reconstruction, fragmental removal and photon splitting, which are written or introduced by authors, will be discussed below.

### 5.3 Photon reconstruction algorithm

The photon reconstruction algorithm refers to the more sophisticated photon identification of the two main identification mechanisms, before the track-cluster association and reclustering process (see section 4.4.4). The algorithm has the following steps: coarsely forming photon clusters, reconstructing photon candidate, photon ID test, and optional fragment removals. Reconstructing photon candidate requires further explanation of the two dimensional peak finding algorithms in section 5.4. The photon ID test involves a multi dimensional likelihood classifier, which is described in section 5.5.

### 5.3.1 Form photon clusters

This step finds large potential photon clusters. All calorimeter hits in the ECAL, which are not used in previous algorithms, are grouped into clusters using a cone based clustering algorithm. To find neutral photon clusters, which do not deposit energies in the tracking system, the cone clustering algorithm is seeded with energetic hits. The parameters for the cone clustering are generous, allowing potentially two or three photons in one cluster.

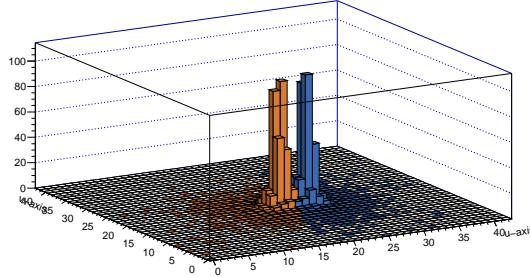
### 5.3.2 Reconstruct photon candidates

The large photon clusters are split into smaller photon candidates, using two-dimensional shower profiles. The candidates close to a track projection are deemed as non-photons. Identifying photon candidates within a large photon cluster relies on the characteristic electromagnetic showers, in particular the transverse distribution. A energetic photon or electron hits the absorber layers of the ECAL, it initiates an electromagnetic shower, where electron pair production and bremsstrahlung produce more low-energy photons and electrons. The transverse distribution is characterised by a narrow cone, widening while the shower develops.

To view the transverse shower distribution, a two-dimensional energy deposition projection is constructed in the plane perpendicular to the direction of the cluster. figure 5.1 shows the energy deposition projection of two photons candidates. U and V axis are two arbitrary orthogonal axis in the transverse plane perpendicular to the direction of photons. Z axis shows the sum of the calorimeter hit energy in GeV. The bin size corresponds to the square ECAL cell size.

By using the two-dimensional energy deposition projection, separating photons translates to separating peaks in the projection. Therefore a high performance two dimensional peak finding algorithm is the key to identify multiple photons. The peak finding algorithm will be discussed in section 5.4

The output of this step is a collection of photon candidates from a photon cluster, which will be fed to the photon ID test.



**Figure 5.1:** Two 500 GeV photons (yellow and blue), just resolved in the transverse plane perpendicular to the direction of the flight, of their energy deposition in electromagnetic calorimeter. U and V axis are two arbitrary axis perpendicular to each other in the plane. Z axis is the sum of the calorimeter hit energy in each particular bin in 2D plane in GeV.

### 5.3.3 Photon ID test

Photon ID test decides if a candidate is a photon. If a candidate is not a photon, the calorimeter hits of the candidate will be passed on to the next stage of the reconstruction. The photon ID test is a multidimensional likelihood classifier. The classifier is trained with discriminating variables, which exploit features electromagnetic showers. The classifier will be discussed in section ??

### 5.3.4 Photon Fragment removal

The optional photon fragment removal aims to merge small photon fragment to main photons. Since this step shares the same logic as the algorithm in section ??, only differing in the cut-off values for merging metrics, this step be discussed in section ??.

This step marks the end of the photon reconstruction algorithm. The output are a collection of reconstructed photons, separated from non-photon calorimeter hits.

## 5.4 Two dimensional peak finding algorithm for photon candidate

As discussed in section 5.3.2, separating photon candidates from a cluster is same as identifying peaks in a two dimensional histogram. An example of two photons is shown

in the figure 5.1. The basic algorithm treats all clusters as potential photon clusters. Since charged hadrons would deposit tracks in the tracking system, extra care is taken when a cluster is close to the projection of the track in the front of the ECAL. The basic peak finding algorithm has two main functions: identifying peaks, and assigning bins to peaks.

A two dimensional histogram is constructed using a plane orthogonal to the direction of the flight of the photon cluster. The U and V axis in figure 5.1 are two orthogonal axes in the plane. The width of the histogram is determined by the size of the ECAL square cell. The calorimeter hits of the photon cluster are then projected onto the two dimensional histogram. The height of the bin is the sum of the calorimeter hit energies in the bin.

A local peak is defined as a bin where its height is above all eight neighbouring bins. After all peak bins are found, non-peak bins are associated to one peak bin, by choosing the peak bin that minimise the metric

$$\frac{d}{\sqrt{E_{\text{peak}}}} \quad (5.1)$$

where  $d$  is the Euclidean distance between a non-peak bin and a peak bin on the histogram, and  $E_{\text{peak}}$  is the height of the peak bin, which is the energy. Alternative metrics provided in the algorithm include  $d$ ,  $\frac{d}{E_{\text{peak}}}$ , and  $\frac{d}{E_{\text{peak}}^2}$ . The default metric is chosen due to a good balance between distance and energy of the peak.

#### 5.4.1 Candidate close to track projection

If a cluster or a photon candidate is close to the projection of the track in the front of the ECAL, it is likely that the cluster or the candidate is a charged hadron. Misidentifying a charged hadron as a photon leads to significant degradation in reconstruction performance. However, if a photon next to a charged hadron is carefully reconstructed, the overall reconstruction is improved. Hence this step aims to carefully identifies photon candidate next to charged hadrons, by using track information and features of the electromagnetic shower. Photon induced electromagnetic shower in the ECAL typically start in the first few layers. As the shower develops, the direction of the shower core does not change much.

If a peak bin is within the eight neighbouring bins of the track projection onto the two dimensional plane, the peak and its associated bins are flagged as non-photons. Furthermore, the ECAL is sliced longitudinally to help identify photon candidates. For example, the default three slices will result in three ECAL fiducial spaces, each contains space from the front of the ECAL to a third, two thirds and the back of the ECAL, respectively. The peaking finding algorithm is repeated for the same cluster divided in each ECAL fiducial space. The peak is only preserved as a photon candidate if the peak exists in every fiducial space, and if its position is shifted by no more than one neighbouring bin between fiducial spaces.

### 5.4.2 Peak filtering

The performance of the two dimensional peaking finding algorithm is improved by clever programming and physics arguments. For a given two dimensional histogram, such as the one in figure 5.1, major peaks most likely correspond to physical photons, while the minor peaks more likely come from fluctuations in energy deposition. To select major peaks, every time after non-peak bins are associated with peak bins, minor peaks with fewer than three bins associated (including the peak bin) are discarded. These bins are then associated with non-discarded peaks. The algorithm also allows bins with height below a critical value to not participate in the peak finding. The default value is set such that only empty bins are not used.

### 5.4.3 Inclusive mode

The two dimensional histogram is iterated a few times during the algorithm. The time complexity is  $O(n^2)$  for a  $n \times n$  histogram (Default  $n = 41$ ). Therefore, for the purpose of speed, it is undesirable to have a very large histogram. However, since the histogram has a finite size, only energy deposition projected on the histogram would be considered for peak finding. This behaviour is suitable for photon reconstruction (section 5.3.2) and test for photon fragment removal (section 5.6). However, for photon splitting (section 5.8), there should be no calorimeter hits loss from splitting a photon. Hence inclusive mode of the peak finding algorithm is developed, and it allows energy deposition projected outside the histogram to be associated with identified peaks.

## 5.5 Likelihood classifier for photon ID

Section 5.3.3 outlines the photon ID test in the photon reconstruction algorithm. This section describes the multidimensional likelihood classifier in details, including discriminating variables. For each photon candidate, a set of kinematic variables are calculated. The classifier training typically uses simulated jet events.

### 5.5.1 Overview of Projective Likelihood

Projective likelihood model (PDE) is used in PandoraPFA for the photon ID due to its simplicity and low requirement on computing resources.

PDE implemented calculates the probability density for each discriminative variable, for signal and background. The overall signal and background likelihood are defined as products of the individual probability density. The likelihood ratio,  $R$ , is then defined as the signal likelihood over signal plus background likelihood.

To use the likelihood ratio, one way is to fit an underlying function to the probability density, which is implemented the TMVA software package. The other way is to use binned likelihood ratio,  $R$ , as the output, due to the simplicity. This is implemented in the PandoraPFA. Similarly to classifier like the rectangular cut method, PDE works better with decorrelated, gaussian like variables. The PandoraPFA implementation did not decorrelate nor transform the variables, to keep implementation fast.

### 5.5.2 Projective Likelihood in PandoraPFA

Kinematic variables to obtain the probability distribution exploit the differences between a characteristic electromagnetic shower and a hadronic shower, and the fact that a photon is more likely to be isolated from other showers and charged tracks. Two variables use the longitudinal shower distribution: the first ECAL layer of the shower, and the difference between a expected longitudinal distribution and the observed. Two variable uses the transverse shower distribution: in the transverse plane with two orthogonal axes, the r.m.s. distance of associated bins to the peak bin, and the smallest ratio of the two r.m.s. distances in each axis direction. One variable is the ratio between the photon candidate energy to the photon cluster energy. The last kinematic variable is the distance between

the photon candidate and the closest track projection. The distributions of kinematic variables are normalised to probability distribution, stored in binned histograms.

TODO explain variables explain longitudinal and transverse

Furthermore, the classifier is improved by realising the kinematic variable distributions depend on the photon energy. Thus these distributions are divided by bins of photon candidate energy. The numbers of photon and non-photon candidates in each energy bin are also different, which helps the ID test. The default energy bins edges are 0.2, 0.5, 1, 1.5, 2.5, 5, 10, 20 GeV, which covers a good range of photon energies. Candidate with energy below 0.2 GeV would not be examined in this step, as it is very unlikely to be a photon.

For a given candidate, which falls in a energy bin, the likelihood classifier output is given by

$$\text{pid} = \frac{N \prod P_i}{N \prod P_i + N' \prod P'_i} \quad (5.2)$$

where  $P_i$  and  $P'_i$  are the probability of  $i^{\text{th}}$  kinematic variable of photon and non-photon candidates.  $N$  and  $N'$  are the number photon and non-photon candidates. These are obtained during classifier training.

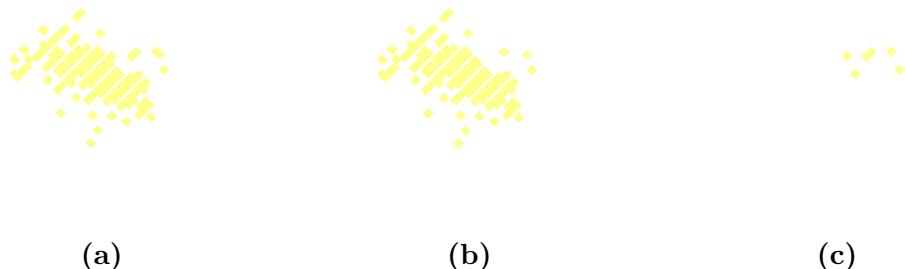
During classification, a candidate passes the photon id test if

$$\begin{cases} \text{pid} > 0.6, & \text{if } 0.2 < E < 0.5 \text{ GeV} \\ \text{pid} > 0.4, & \text{if } E \geq 0.5 \text{ GeV} \end{cases} \quad (5.3)$$

where  $E$  is the candidate energy. Two values of the  $\text{pid}$  cuts reflect the confidence of the id test with different candidate energy. The test is more cautious with low energy candidate.

## 5.6 Photon fragment removal algorithm in the ECAL

During the reconstruction, it is possible that a core of the photon electromagnetic shower is identified as a photon (the main photon). The outer part of the shower is reconstructed as a separate particle, and wrongly identified as a photon or a neural hadron (the photon/neutral fragment). figure 5.2 shows a typical creation of such a photon fragment. The fragment does not have the electromagnetic shower structure, and typically it is has much lower energy than the main photon. If a photon-fragment pair is merged, the pair should be consistent with a one-particle profile. These characteristics are used to merge fragments to main photons.



**Figure 5.2:** An event display of a typical 10 GeV photon (figure 5.2a), reconstructed into a main photon (figure 5.2b) and a photon fragment (figure 5.2c).

Photon fragment removal algorithms can exist in multiple step in the reconstruction, at the end of the photon reconstruction (see Section 5.3.4), or at the end of the reconstruction. Since these algorithms share the same base class, the latter one will be discussed. The former differs mostly in the default cut-off values for merging metrics.

A photon and a potential fragment form a pair of particles (photon-fragment pair), if they are spatially close. Kinematic and topological properties of the photon-fragment pair are examined. The pair is merged when the properties pass a set of cuts, developed by comparing true photon-fragment pairs and non photon-fragment pair. This merging test is iterated over all possible photon-fragment pairs. If multiple photon-fragment pairs pass the merging test, the pair with closest distance metric,  $d$ , will be merged.

The photon-fragment pairs is classified into photon-photon-fragment pairs and photon-neutral-hadron-fragment pairs, because they have different kinematic and topological distributions. The pairs are further classified into low energy and high energy pair,

depending on whether the fragment energy ( $E_f$ ) above 1 GeV. The cuts for merging pairs, are classified which will be explained later, are listed in table 5.1.

Low $E_f$	Photon-photon	Photon-neutral-hadron
transverse shower comparison	$d < 30, \frac{E_{p1}}{E_m+E_f} > 0.9, \frac{E_{p2}}{E_f} < 0.5, E_{p1} > E_m$	-
close proximity	-	$d < 20, d_c < 40$
low energy fragment	$d < 20, E_p < 0.4$	-
small fragment 1	$d < 30, N_{calo} < 40, d_c < 50$	$d < 50, N_{calo} < 10, d_h < 50$
small fragment 2	$d < 50, N_{calo} < 20$	-
small fragment forward region	$N_{calo} < 40, d_c < 60, E_f < 0.6,  \cos(\theta_Z)  > 0.7$	-
relative low energy fragment	$d < 40, d_h < 20, \frac{E_f}{E_m} < 0.01$	$d < 40, d_h < 15, \frac{E_f}{E_m} < 0.01$
High $E_f$	Photon-photon	Photon-neutral-hadron
transverse shower comparison	$\frac{E_{p1}}{E_m+E_f} > 0.9, E_{p2} = 0 \text{ or } (\frac{E_{p2}}{E_f} < 0.5, E_{p1} > E_m)$	$\frac{E_{p1}}{E_m+E_f} > 0.9, E_{p2} = 0 \text{ or } (\frac{E_{p2}}{E_f} < 0.5, E_{p1} > E_m)$
relative low energy fragment 1	$d < 40, d_h < 20, \frac{E_f}{E_m} < 0.02$	$d < 40, d_h < 20, \frac{E_f}{E_m} < 0.02$
relative low energy fragment 2	-	$d < 40, d_h < 20, \frac{E_f}{E_m} < 0.1, E_f > 10$
relative low energy fragment 3	-	$d < 20, d_h < 20, \frac{E_f}{E_m} < 0.2, E_f > 10$

**Table 5.1:** The cuts for merging photon-photon-fragment pairs and photon-neutral-hadron-fragment pairs for both low energy and high energy fragments.  $d$ ,  $d_c$  and  $d_h$  are the mean energy weighted intra-layer distance of the pair, the distance between centroids, the minimum distance between calorimeter hits of the pair.  $E_m$  and  $E_f$  are the main photon energy and the fragment energy.  $E_{p1}$  and  $E_{p2}$  are the two largest peaks, found by peak finding algorithm, ordered by descending energy.  $N_{calo}$  is the number of the calorimeter hits in the fragment.  $|\cos(\theta_Z)|$  is the absolute cosine of the polar angle, where beam direction is the z-axis.

Table 5.1 lists cuts for merging photon-photon-fragment pairs and photon-neutral-hadron-fragment pairs for both low energy and high energy fragments.  $d$ ,  $d_c$  and  $d_h$  are the mean energy weighted intra-layer distance between each PFO in the pair, the distance between centroids, the minimum distance between calorimeter hits of each PFO in the pair, respectively.  $E_m$  and  $E_f$  are the main photon energy and the fragment energy.  $E_{p1}$  and  $E_{p2}$  are the two largest peaks and associated calorimeter hits, found by the two dimensional peak finding algorithm (section 5.4), ordered by descending energy, using

the pair as input.  $N_{\text{calo}}$  is the number of the ECAL hits in the fragment.  $|\cos(\theta_Z)|$  is the absolute cosine of the polar angle of the main photon, where beam direction is the z-axis.

Three distance measurements have subtle difference.  $d_c$  gives the distance between centroids of each PFO in the pair, which is a quick but crude measurement.  $d_h$  is the minimum distance between calorimeter hits of each PFO in the pair. For a true photon-fragment,  $d_h$  should be close to zero as the pair should be spatially close.  $d$  is the mean energy weighted intra-layer distance between each PFO in the pair:

$$d = \frac{\sum_i^{\text{layers}} d_{l,i} E_{f,i}}{\sum_i^{\text{layers}} E_{f,i}} \quad (5.4)$$

where  $i$  indicates  $i^{\text{th}}$  pseudo-layer of the ECAL.  $d_{l,i}$  is the minimum distance between calorimeter hits of the pair in the  $i^{\text{th}}$  pseudo-layer.  $E_{f,i}$  is the energy of the fragment in the  $i^{\text{th}}$  pseudo-layer.  $d$  is a better measurement of the closeness of the pair. Similar to  $d_h$ ,  $d$  will be very small for a true photon-fragment pair.

One logic for merging is when the fragment is small with low energy and is close to the main photon. The other logic is when the pair looks like one photon in two-dimensional energy deposition projection (see section 5.3.2 and figure 5.1). Comparing low  $E_f$  and high  $E_f$  cut, the cuts are similar. High  $E_f$  cuts are more relaxed on the energy comparison for small fragment test. Comparing photon-photon-fragment pair and photon-neutral-hadron-pair, cuts for photon-neutral-hadron-pair are more conservative for low  $E_f$ , but more relaxed for high  $E_f$ . This reflects that the neutral hadron fragments originated from charged particles are more likely to be low energy.

Since all possible photon-fragment pairs are compared, this is a costly cooperation with  $O(n^2)$  time complexity for  $n$  particles. The speed is improved by considering only the pairs with  $d < 80\text{mm}$ . The algorithm occurs at the end of the reconstruction.

## 5.7 High energy photon fragment recovery algorithm

Section 5.6 described effective algorithms to remove photon fragments that are peripheral to the main photon, or the electromagnetic shower core. An example of such fragment is shown in figure 5.2. There is another type of fragment which is the leakage effect of the ECAL. When the high energy photon shower is not fully contained in the ECAL, shower

deposits energy in the HCAL, which often forms a neutral hadron in the HCAL. Photon reconstruction, as described in section 5.3, considers only calorimeter hits in the ECAL. An example of a 500 GeV photon reconstructed into a main photon in the ECAL (yellow) and a neutral hadron fragment in the HCAL (blue) is shown in figure 5.3. For the ILD detector, this ECAL leakage effect appears when the photon energy is above 50 GeV.



**Figure 5.3:** An event display of a typical 500 GeV photon, reconstructed into a main photon in the ECAL (yellow) and a neutral hadron fragment in the HCAL (blue).

With figure 5.3 as an example, high energy fragments in the HCAL is spatially close to the main photon. A fitted cone from the main photon covers most of the fragment if extended to the HCAL. These features allow a set of cuts developed to merge high energy fragments, listed in table 5.2

This algorithm would collect the photons and neutral hadrons in the HCAL as inputs. It occurs after the first pass of topological association in the reconstruction, which connects tracks to clusters in the ECAL and the HCAL. The algorithm would iterate over all pairs of reconstructed photons and neutral hadrons in the HCAL. For each pair, a set of variables are calculated and compared to a set of cuts (table 5.2). Photon-fragment pairs passing the cuts will be merged.

Fragment in the HCAL should be spatially close to the main photon, measured by three metrics.  $d_c^l$  is the distance between centroids of the last outer layer of the main photon and the first inner layer of the fragment.  $d_{cone}^l$  is the distance between fitted cones using the last outer layer of the main photon and the first inner layer of the fragment.  $d_{cone}$  is the distance between fitted cones using the main photon and the fragment.

High energy fragment recovery	Cuts
distance comparison	$d_c^l \leq 173 \text{ mm}$ , $d_{cone}^l \leq 100 \text{ mm}$ , $d_{cone} \leq 100 \text{ mm}$
shower width comparison	$0.3 \leq \frac{w_f^l}{w_m^l} \leq 5$
projection comparison	$r_f \leq 45 \text{ mm}$
energy comparison	$\frac{E_f}{E_m} \leq 0.1$
cone comparison	$\%N_{calo,cone} \geq 0.5$

**Table 5.2:** The cuts for merging high energy photon fragment in the HCAL to the main photon in the ECAL.  $d_c^l$  is the distance between centroids of the last outer layer of the main photon and the first inner layer of the fragment.  $d_{cone}^l$  is the distance between fitted cones using the last outer layer of the main photon and the first inner layer of the fragment.  $d_{cone}$  is the distance between fitted cones using the main photon and the fragment.  $w_m^l$  and  $w_f^l$  are the r.m.s. width of the last outer layer of the main photon and the first inner layer of the fragment.  $r_f$  is the r.m.s. mean energy weighted distance of a calorimeter hit in the fragment to the direction of the main photon.  $E_m$  and  $E_f$  are the main photon energy and the fragment energy.  $\%N_{calo,cone}$  is the fraction of the calorimeter hits in the fragment in the extended fitted cone of the main photon.

The direction of the fragment should be similar to that of the main photon.  $r_f$ , the r.m.s. mean energy weighted distance of a calorimeter hit in the fragment to the direction of the main photon, has to be small for merging.

Another feature of the fragment and the main photon is that the shower width should be similar.  $w_m^l$  and  $w_f^l$  are the r.m.s. width of the last outer layer of the main photon and the first inner layer of the fragment. The ratio  $\frac{w_f^l}{w_m^l}$  needs to be in the range of 0.3 to 5. The generous upper bound is due to the HCAL is coarser than the ECAL.

When a fitted cone from the main photon is extended to the HCAL, the cone should contain a significant amount of the fragment.  $\%N_{calo,cone}$ , the fraction of the calorimeter hits in the fragment in the extended fitted cone of the main photon, has to be no less than 0.5 for the merging.

The last criteria is the fragment should have low energy relative to the main photon.  $E_m$  and  $E_f$  are the main photon energy and the fragment energy. The ratio,  $\frac{E_f}{E_m}$ , has to be less than 0.1 for the merging.

If multiple photon-fragment pairs pass the cuts with the same fragment, the pair with highest  $\%N_{calo,cone}$  will be merged.

## 5.8 Photon splitting algorithm

Algorithms described above deal with forming photons from calorimeter hits in the ECAL, merging photon fragments in the ECAL and the HCAL. Another aspect in photon reconstruction is splitting accidentally merged photons. During the particle reconstruction, it is possible that photons are accidentally merged if they are spatially close. Hence another algorithm at the end of the particle reconstruction addresses this issue and tries to split merged photons.

Merged photon is typically energetic. The merged photon should be consistent with topologies of a spatially closed photon pair. Extra care should be taken if the photon is close to a charged PFO. Many PandoraPFA algorithms deal with track clusters association and there is a greater confidence in clusters associated with tracks. These features form logics behind the algorithm.

Photon splitting	Cuts
Cuts	$E > E_{c1}$ , $E_{p2} > E_{c2}$ , $N_p < 5$
$E_{c1}$ and $E_{c2}$ values	
0 nearby charged PFO	$E_{c1} = 10$ , $E_{c2} = 1$
1 nearby charged PFO	$E_{c1} = 10$ , $E_{c2} = 5$
> 1 nearby charged PFO	$E_{c1} = 20$ , $E_{c2} = 10$

**Table 5.3:** The cuts for splitting photons, and the values for energy cut-off points.  $E$  is the photon energy.  $E_{p2}$  is energy if the second largest peak from the two dimensional peak finding.  $N_p$  is the number of peaks identified by the peak finding.  $E_{c1}$  and  $E_{c2}$  are the energy cut-off values, determined by the number of nearby charged PFOs.

The table 5.3 shows values for the splitting a photon.  $E$  is the photon energy.  $E_{p2}$  is energy if the second largest peak from the two dimensional peak finding.  $N_p$  is the number of peaks identified by the peak finding.  $E_{c1}$  and  $E_{c2}$  are the energy cut-off values, determined by the number of nearby charged PFOs. When a energetic photon is identified, and a energetic second peak can be found by the peak finding, the photon is likely from a photon pair.  $N_p$  cut is because a reconstructed photon is unlikely from more than four photons. The values of  $E_{c1}$  and  $E_{c2}$  allow more conservative approach when a photon is close to charged PFOs.

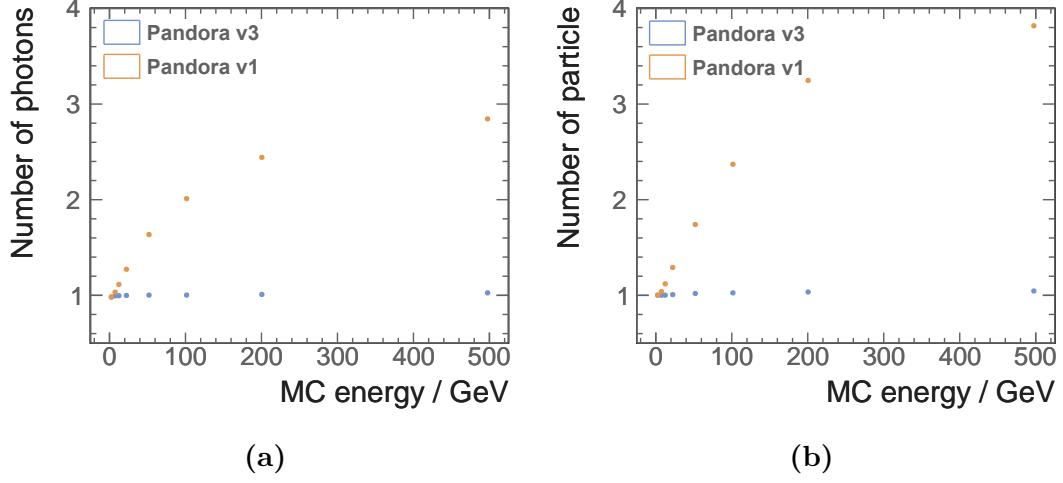
## 5.9 Photon reconstruction performance improvement

Motivations and implementations of four different algorithms for photon reconstruction, fragment removal and photon splitting have been described in the above. The main photon reconstruction algorithm in section 5.3 improves the photon completeness and the photon pair resolution, due to the improved two dimensional peak finding algorithm in section 5.4. The fragment removal algorithms in section 5.6 and section 5.7 further reduce the photon fragments in the ECAL and the HCAL. The photon splitting algorithm in section 5.8 exploits the peak finding algorithms and improves the photon separation resolution. Because of the high photon reconstruction completeness, the jet energy resolution receives a small improvement.

This section reviews the performance improvement with the introduced algorithms, using single photon, photon pair and jet samples. The performance was compared using PandoraPFA version 1 and version 3, where the photon algorithms were introduced in PandoraPFA version 2. The ILD detector model is used. The photon pair simulated events were generated with a uniform distribution in the solid angle for a range of the opening angles between the pair. The events selected such that there is no early photon conversion and the monte carlo photon deposits energies in the calorimeter. The events are further restricted to photon decaying in barrel and end cap region only, to minimise the detector effect.

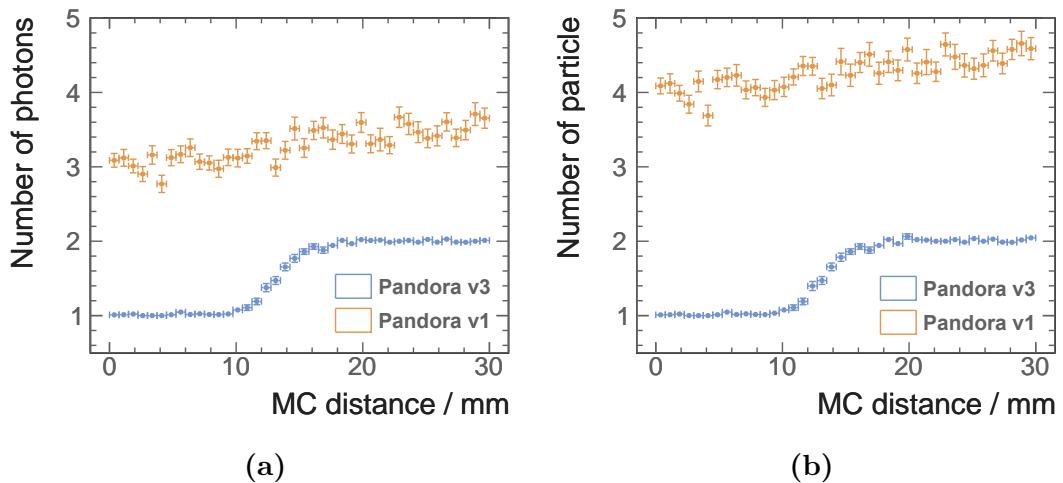
figure 5.4a shows the reduction in fragments identified as photons, using a single photon per event sample. For the blue dots, the average number of photon stays below 1.05, where the true value is 1, even at high energy. A similar trends shows in figure ??, where the extra fragments identified as neutral hadrons have taken into account. For a 100 GeV photon, the average numbers of photon and particle are reduced to 1 from 2 and 2.4. For a 500 GeV photon, the average numbers of photon and particle are reduced to 1.05 from 2.8 and 3.8.

figure 5.5 illustrates a similar reduction in the photon fragments and the neutral hadron fragments using two photons of 500 and 50 GeV per event sample. The high energy photon are more likely to create fragments. And the imbalance in the two photon energies makes it more difficult to separate correctly. The figure shows the MC distance separation from 0 to 30 mm, which corresponds to approximately 6 ECAL square cell. In both figure 5.5a and figure 5.5b, the average numbers of photon and particle are below



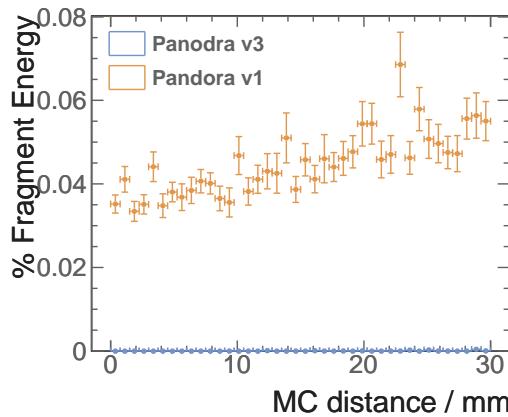
**Figure 5.4:** figure 5.4a and figure 5.4b shows the average number of reconstructed photons and reconstructed particles, as a function of their true energy using a single photon per event sample. The top orange and bottom blue dots are reconstructed with PandoraPFA version 1 and version 3. The photon reconstruction is changed in PandoraPFA version 2.

2.05 at 30 mm apart, which is significantly better than reconstruction in PandoraPFA version 1. Two photons start to be resolved at 10 mm apart, and fully resolved at 20 mm apart. The resolution is better than reconstruction in PandoraPFA version 1, which is difficult to extract due to excess fragments.



**Figure 5.5:** figure 5.5a and figure 5.5b shows the average number of reconstructed photons and reconstructed particles, as a function of the MC distance separation in the calorimeter, using two photons of 500 and 50 GeV per event sample. The top orange and bottom blue dots are reconstructed with PandoraPFA version 1 and version 3. The photon reconstruction is changed in PandoraPFA version 2.

Another metric to reflect the improvement in photon reconstruction is the fragment energy fraction of the total energy as function of the distance separation. Shown in figure 5.6, using two photons of 500 and 50 GeV per event sample, a reduction in fragment energy can be seen clearly. With improved reconstruction, the average fragment energy fraction is below 0.1% up to 30 mm apart, whilst around 5% energy would be in fragments with reconstruction in PandoraPFA version 1.



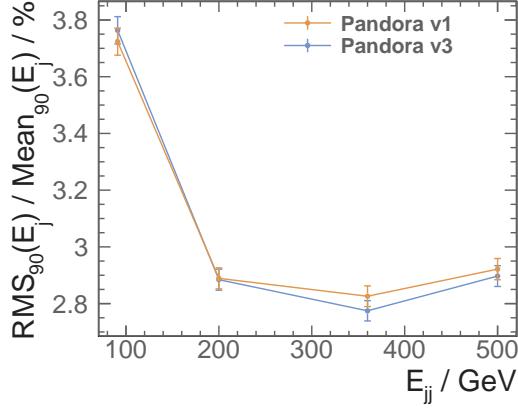
**Figure 5.6:** figure 5.6 shows the average fraction fragments energies of the total energy, as a function of the MC distance separation in the calorimeter, using two photons of 500 and 50 GeV per event sample. The top orange and bottom blue dots are reconstructed with PandoraPFA version 1 and version 3. The photon reconstruction is changed in PandoraPFA version 2.

TODO write RMS90

The improvement in completeness and resolution in photon reconstruction, as shown in single photon and double photon reconstruction, leads to a small improvement in the jet energy resolution at high energy. Jet energy resolution is defined as the root mean squared divided by the mean for the smallest width of distribution that contains 90% of entries, using  $Z' \rightarrow u/d/s$  sample. The di-jet energy is sampled at 91, 200, 360 and 500 GeV. Shown in figure 5.7, the jet energy resolutions are better at 360 and 500 GeV with improved photon reconstruction.

The improvement of the photon is also demonstrated in chapter 6, where tau lepton decay modes are classified. Excellent photon reconstruction leads to a high classification rate.

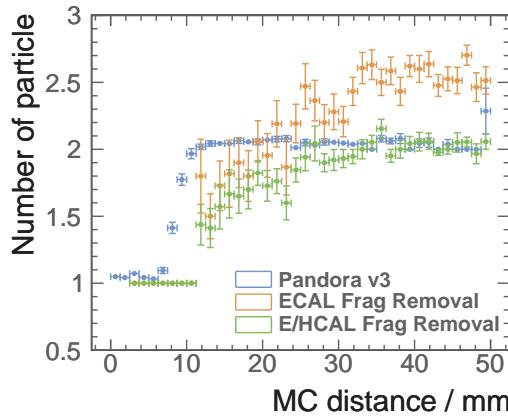
TODO write other people using it



**Figure 5.7:** figure 5.7 shows jet energy resolution as a function of the di-jet energy using  $Z' \rightarrow u/d/s$  sample. The top orange and bottom blue dots are reconstructed with PandoraPFA version 1 and version 3. The photon reconstruction is changed in PandoraPFA version 2.

## 5.10 Breakdown of photon reconstruction improvement

As stated before, photon reconstruction algorithm in section 5.3 and photon splitting algorithm in section 5.8 improves the photon completeness and the photon pair resolution. The fragment removal algorithm in section 5.6 removes fragments in the ECAL. High energy fragment removal algorithm in section 5.7 removes fragments in the HCAL. To show the incremental improvement, the average number of particle for a high energy photon pair, 500 - 500 GeV is shown in figure 5.8. With fragment removal algorithm in the ECAL, the number of fragment is reduced significantly comparing to figure 5.5b, shown as the orange dots. The high energy fragment removal algorithm further reduces the number of fragments, shown as the green dots. At 40 mm apart, with both fragment removal algorithms, there is less than 0.05 fragment per photon pair, which is similar to the best performance. The introduction of the revised photon reconstruction and photon splitting improves the photon separation resolution. Photons pair starts to be resolved at 5 mm instead of 10 mm for 500 - 500 GeV pair. Also two photons are fully resolved at 15 mm instead of 40 mm apart.



**Figure 5.8:** Figure shows the average number of photons, as a function of the MC distance separation in the calorimeter, using two photons of 500 and 500 GeV per event sample. The blue, orange, and green dots are reconstructed with PandoraPFA version 3, PandoraPFA version 1 with fragment removal in the ECAL (section ??), and PandoraPFA version 1 with fragment removal in the ECAL and the HCAL. The photon reconstruction is changed in PandoraPFA version 2.

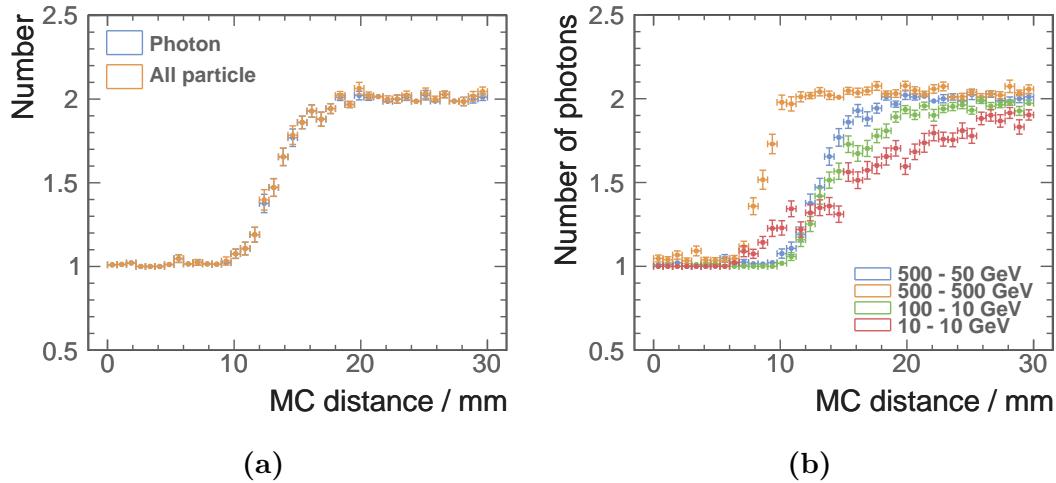
## 5.11 Photon reconstruction performance

In section 5.9, the improved performance of the photon reconstruction is demonstrated with different metrics, using single photon, double photons and jet samples. In this section, the features of the photon reconstruction will be described.

For simple samples such as two photons per event, there are very few fragments. Shown in figure 5.9a for 500 and 50 GeV photons pair sample, the average number of photons beyond 20 mm apart is 2 within errors. The average number of particle is less than 0.05 larger than the average number of photons.

The resolving power of a photon pair depends on energies of two photons. figure 5.9b is an example of average number of photon reconstructed for differen photon pairs. When the energies of two photons are similar, the resolving distance is shorter. This is because that the two photon showers have similar sizes, and the peak finding algorithm can exploit the symmetry. For example, 500 - 500 GeV photon pair and 10 - 10 GeV photon pair start to be resolved at 6 mm apart, which is about 1 ECAL cell. The asymmetrical photon pair, 500 - 50 GeV and 100 - 10 GeV pair, starts to be resolved at 10 mm apart, which is about 2 ECAL cell.

For the energetic photon, it is more difficult to remove fragments, but it is easier to identify the photon. The electromagnetic shower core is more dominant than the



**Figure 5.9:** figure 5.9a shows the average numbers of photon and particle using two photons of 500 and 50 GeV per event sample. figure 5.9b shows the average numbers of photon for four different photon pairs: 500 - 50, 500 - 500, 100 - 10 and 10 - 10 GeV.

peripherals. Therefore separating two energetic photons is easier than separating two low energy photons. This can be seen in figure 5.9b. At 20 mm apart, two photons in 500 - 500 GeV pair are fully resolved, where approximately 60% of two photons in 10 - 10 GeV pair are resolved.

This set of photon related algorithms have been incorporated into the default reconstruction chain in PandoraPFA. The CLIC simulation studies have benefited from the improved photon reconstructions in various physics process, such as  $H \rightarrow \gamma\gamma$ .

# Chapter 6

## Tau Lepton Final State Separation

*“MVA: Turn numbers into gold.”*

— TMVA

### 6.1 Introduction

Why study tau

Tau lepton has been examined closely in the past. The decay and the spin of the decay product were direct tests to the standard model. The spin of the decay product, using a Higgs decaying to tau tau channel, allows one to determine the spin of the higgs. Also, as tau is short-lived, only its decay products can be detected and reconstructed in the detector. Therefore, the ability to reconstruct and separate different tau decay modes is benchmark of detector performances.

This chapter will describe a tau final decay separation study. The processor developed for the study is used to test different detector models, as a proof-of-principle of detector optimisation using tau decay separation. Lastly, the spin of the Z was studied using one Z decaying to two tau tau channel.

## 6.2 Simulation and reconstruction

$e^-e^+ \rightarrow \tau^-\tau^+$  channel is used for the tau decay mode separation study. Generator software WHIZARD 1.95 [29] is used to generate simulated Monte Carlo (MC) samples. Hadronisation is described with PYTHIA 6.4 [31], which is tuned to the LEP results [30]. The spin effect of tau lepton decay is described by TAUOLA [33].

Final state radiation (FSR) was simulated. The initial state radiation (ISR) and the beam induced background were not simulated.

Events were simulated with the CLIC\_ILD detector concept, using software with MOKKA [36], based on the GEANT 4 package [35]. Events were reconstructed with ilcsoft version v01-17-07 [40] and PandoraPFA version v02-02-00 [38], where the photon reconstruction is described in [41].

## 6.3 Generator level cut

To study the difference between different tau decay modes, clear topological difference is required. Therefore, events were considered if the event passes a set of cuts at generator level, listed here

- the final state photons not converting to electron pair in the tracker,
- the tau leptons decaying in the barrel and the end cap regions, which are defined as polar angle between 0.3 to 0.6 rad and 0.8 to 1.57 rad, and
- the visible energy of the tau lepton decay products more than 5 GeV, where the visible energy of the tau lepton decay is defined as the energy of the tau minus the energy of the tau neutrino.

The angular requirement is due to the gap region between the barrel and the end cap of calorimeters, which degrades the PFO resolution significantly.

Around two million events were simulated for this study.

**Table 6.1:** Branching ratios of the seven major  $\tau^-$  decays, taken from [2].  $\tau^+$  decays similarly to  $\tau^-$ .

Decay final state	Branching ratio / %
$e^- \bar{\nu}_e \nu_\tau$	$17.83 \pm 0.04$
$\mu^- \bar{\nu}_\mu \nu_\tau$	$17.41 \pm 0.04$
$\pi^- \nu_\tau$	$10.83 \pm 0.06$
$\rho(\pi^-\pi^0)_{770} \nu_\tau$	$25.52 \pm 0.09$
$a_1(\pi^-\pi^0\pi^0)_{1260} \nu_\tau$	$9.30 \pm 0.11$
$a_1(\pi^-\pi^-\pi^+)_{1260} \nu_\tau$	$8.99 \pm 0.06$
$\pi^- \pi^- \pi^+ \pi^0 \nu_\tau$	$2.70 \pm 0.08$

## 6.4 Decay modes

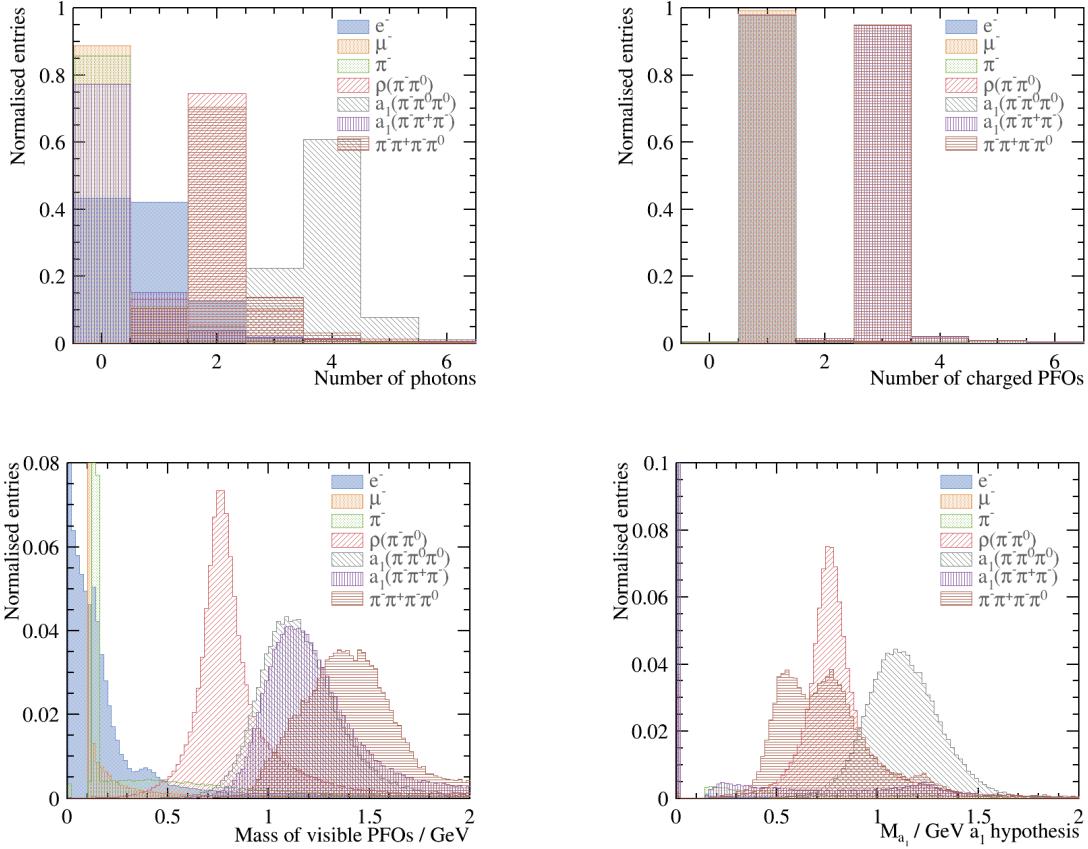
Seven major decay final states of the tau lepton shown in table 6.1 were studied, covering 92.58 % of all tau decays [2]. Decay modes not listed in the table have branching fractions lower than 1% each. These final states can be classified into three categories: leptonic decays ( $e^- \bar{\nu}_e \nu_\tau$  and  $\mu^- \bar{\nu}_\mu \nu_\tau$ ), one-prong with photons ( $\pi^- \nu_\tau$ ,  $\rho(\pi^-\pi^0)_{770} \nu_\tau$  and  $a_1(\pi^-\pi^0\pi^0)_{1260} \nu_\tau$ ), and three-prong with photons ( $a_1(\pi^-\pi^-\pi^+)_{1260} \nu_\tau$  and  $\pi^- \pi^- \pi^+ \pi^0 \nu_\tau$ ).

The studied channel,  $e^- e^+ \rightarrow \tau^- \tau^+$ , contains two  $\tau$  decaying in opposite directions. To select decay products of one  $\tau$ , the fiducial detector space was divided into two halves. Event shape variable thrust is used to separate two halves. The classical event shape thrust [42], is defined as

$$T = \max_{\hat{t}} \frac{\sum_i |\hat{t} \cdot \vec{p}_i|}{\sum_i |\vec{p}_i|} \quad (6.1)$$

where  $\vec{p}_i$  is the momentum vector of the particle  $i$ . Summation is over all particles in the event. Thrust axis,  $\hat{t}$ , is a unit vector. (Principle) Thrust value,  $T$ , is 1 for a perfect pencillike back-to-back two-jet event, and 0.5 for a perfect spherical event. The sign of dot product between thrust axis and PFO momentum determines which half the PFO falls into.

## 6.5 Discriminative variables



**Figure 6.1:** Normalised distribution for selected discriminative variables for seven final states,  $e^- \bar{\nu}_e \nu_\tau$ ,  $\mu^- \bar{\nu}_\mu \nu_\tau$ ,  $\pi^- \nu_\tau$ ,  $\rho(\pi^-\pi^0)_{770} \nu_\tau$ ,  $a_1(\pi^-\pi^0\pi^0)_{1260} \nu_\tau$ ,  $a_1(\pi^-\pi^+\pi^-)_{1260} \nu_\tau$  and  $\pi^-\pi^+\pi^0 \nu_\tau$ , separated using truth information, for  $\sqrt{s} = 100$  GeV for nominal CLIC\\_ILD detector model. The top left, top right, bottom left and bottom right plots are the normalised entries against the number of photons, number of charged PFOs, invariant mass of visible PFOs, and the invariant mass of  $a_1(\pi^-\pi^0\pi^0)_{1260}$  for hypothesis test, respectively. There is a clear distinction between different final states in each plot.

Some variables with most discriminative power are shown in figure ???. In total 29 variables used in the multivariate analysis. The reason for the large number of variables is due to training seven decay modes at once, which will be discussed later.

Here is a full list of all variables used in the multivariate analysis. Energy of the  $\tau$  is assumed to be the same as the energy of  $e^\pm$  colliding beam, which is half of the  $\sqrt{s} = 100$  GeV energy. Recoil momenta were calculated assuming the  $e^- e^+$  collision happened

at the centre of mass energy. Both assumptions are largely valid when there is no ISR contribution.

- $\frac{E_{ECAL,HCal}}{E_{tot}}$ , **charged**: Sum of energy deposited in ECal and HCal, divided by the energy of charged particles
- $\frac{E_{ECAL,HCal}}{E_{tot}}$ , **all**: Sum of energy deposited in ECal and HCal, divided by the energy of all particles
- $m_{vis}$ : Invariant mass of visible particles in GeV
- $\frac{E_{vis}}{E_{\tau^-}}$ : Sum of energy of all particles, divided by the energy of  $\tau^-$
- $\frac{E_{charged}}{E_{\tau^-}}$ : Sum of energy of charged particles, divided by the energy of  $\tau^-$
- $\frac{E_{\mu^-}}{E_{\tau^-}}$ : Sum of energy of muons, divided by the energy of  $\tau^-$
- $\frac{E_e^-}{E_{\tau^-}}$ : Sum of energy of electrons, divided by the energy of  $\tau^-$
- $\frac{E_\gamma}{E_{\tau^-}}$ : Sum of energy of photons, divided by the energy of  $\tau^-$
- $\frac{E_{\pi^-}}{E_{\tau^-}}$ : Sum of energy of charged pions, divided by the energy of  $\tau^-$
- $N_{charged}$ : Number of charged particles
- $N_{\mu^-}$ : Number of muons
- $N_e^-$ : Number of electrons
- $N_\gamma$ : Number of photons
- $N_{\pi^-}$ : Number of charged pions
- $m_\gamma$ : Invariant mass of photons in GeV
- $m_{charged}$ : Invariant mass of charged particles in GeV
- $m_{neutral}$ : Invariant mass of neutral particles in GeV
- $m_{\pi^-}$ : Invariant mass of charged pions in GeV
- $m_{\pi^0}, \rho(\pi^-\pi^0)_{770}$  hypothesis: Fitted invariant mass of  $\pi^0$  for  $\rho(\pi^-\pi^0)_{770}$  hypothesis test

- $m_{\rho(\pi^-\pi^0)_{770}}, \rho(\pi^-\pi^0)_{770}$  hypothesis: Fitted invariant mass of  $\rho(\pi^-\pi^0)_{770}$  for  $\rho(\pi^-\pi^0)_{770}$  hypothesis test
- $m_{\pi^0}, a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis: First fitted invariant mass of  $\pi^0$ , for  $a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis test, ordered by closeness to the true  $\pi^0$  mass
- $m_{\pi^0}, a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis: Second fitted invariant mass of  $\pi^0$ , for  $a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis test, ordered by closeness to the true  $\pi^0$  mass
- $m_{a_1(\pi^-\pi^0\pi^0)_{1260}}, a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis: Second fitted invariant mass of  $a_1(\pi^-\pi^0\pi^0)_{1260}$ , for  $a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis test
- $\bar{E}_{\text{cell}}$ : Average energy deposited in a calorimeter cell in GeV
- $d_{\text{trans,shower}}$ : Transverse shower width for electromagnetic shower profile, averaged for all clusters in the ECal
- $l_{\text{long,shower}}$ : Longitudinal start layer for electromagnetic shower profile, averaged for all clusters in the ECal
- $\Delta l_{\text{long,shower}}$ : Longitudinal discrepancy for electromagnetic shower profile, averaged for all clusters in the ECal
- %MIP: Fraction of calorimeter hits registered as minimum ionised particles, averaged for all clusters in the ECal
- $\frac{E}{p}$ : Energy divided by momentum, averaged for all clusters in the ECal

Number of photons is an important variable for separating decay modes. This information is only available due to the excellent photon reconstruction. Shown in figure ??, the majority of  $\mu^- \bar{\nu}_\mu \nu_\tau$ ,  $\pi^- \nu_\tau$  and  $a_1(\pi^-\pi^-\pi^+)_{1260} \nu_\tau$  final states have zero photon reconstructed. The  $e^- \bar{\nu}_e \nu_\tau$  final state event have one photon reconstructed instead of zero, due to the FSR effect.  $\rho(\pi^-\pi^0)_{770} \nu_\tau$  and  $\pi^- \pi^- \pi^+ \pi^0 \nu_\tau$  have nearly 80% events with two reconstructed photons, whilst  $a_1(\pi^-\pi^-\pi^+)_{1260} \nu_\tau$  have over 60% events with four reconstructed photons. The loss in efficiency is due to the increasing difficulty to separate nearby photons.

The number of charged PFOs can clearly separate the leptonic and 1-prong final states, from the 3-prong final states, shown in figure ???. The efficiency of leptonic final states are over 98%.

The invariant mass of the visible PFOs shows clear differences between different final states.  $\rho(\pi^-\pi^0)_{770}\nu_\tau$ ,  $a_1(\pi^-\pi^0\pi^0)_{1260}\nu_\tau$  and  $a_1(\pi^-\pi^-\pi^+)_{1260}\nu_\tau$  distribution show clear resonance at  $\rho$  and  $a_1(1260)$ .  $e^-\bar{\nu}_e\nu_\tau$ ,  $\mu^-\bar{\nu}_\mu\nu_\tau$  and  $\pi^-\nu_\tau$  distribution show much smaller invariant mass and  $\pi^-\pi^-\pi^+\pi^0\nu_\tau$  shows a large invariant mass than  $a_1(1260)$ . The  $e^-\bar{\nu}_e\nu_\tau$  final state has a long tail of invariant mass due to the extra photons from the FSR.

$\frac{E_{ECAL,HCal}}{E_{tot}}$ , charged and  $\frac{E_{ECAL,HCal}}{E_{tot}}$ , all are both very effective at picking out leptonic decay modes.

For final states containing  $\rho$  and  $a_1(1260)$  resonance, it is useful to use minimisation test for right pairing of the resonance. We will use  $a_1(\pi^-\pi^0\pi^0)_{1260}$  as an example.  $\rho(\pi^-\pi^0)_{770}$  is very similar.

The minimisation of  $a_1(\pi^-\pi^0\pi^0)_{1260}$  hypothesis states

$$\chi^2_{a_1(1260)} = \left( \frac{m_{a_1(1260),\text{fit}} - m_{a_1(1260)}}{\sigma_{a_1(1260)}} \right)^2 + \left( \frac{m_{\pi^0,\text{fit}} - m_{\pi^0}}{\sigma_{\pi^0}} \right)^2 + \left( \frac{m_{\pi^{0*},\text{fit}} - m_{\pi^0}}{\sigma_{\pi^0}} \right)^2, \quad (6.2)$$

where  $m_{\pi^0,\text{fit}}$  and  $m_{\pi^{0*},\text{fit}}$  are the invariant masses of all possible two photons combinations,  $\sigma_{a_1(1260)}$  and  $\sigma_{\pi^0}$  are the half width of the invariant mass distribution of reconstructed  $a_1(1260)$  and  $\pi^0$  using the truth information, and  $m_{a_1(1260)}$  and  $m_\pi$  are the masses of  $a_1(1260)$  and  $\pi^0$ , taken from [2]. If there are only two or three photons, the  $\chi^2_{a_1(1260)}$  expression will be reduced and not including  $m_{\pi^{0*},\text{fit}}$  term, assuming two photons are merged in the reconstruction. If there are fewer than two photons, the  $\chi^2_{a_1(1260)}$  expression would only contain  $m_{a_1(1260),\text{fit}}$  term.

For the  $\rho(\pi^-\pi^0)_{770}\nu_\tau$  final state, a similar  $\chi^2_{\rho(770)}$  test for  $\rho(770)$  hypothesis is used to extract  $m_{\rho(770),\text{fit}}$  and  $m_{\pi^0,\text{fit}}$  variables.  $\chi^2_{\rho(770)}$  is similar to  $\chi^2_{a_1(1260)}$  with  $\rho(770)$  replacing  $a_1(1260)$  and only one  $m_{\pi^0,\text{fit}}$  term.

figure ?? shows the  $m_{a_1(1260),\text{fit}}$  where  $\rho(\pi^-\pi^0)_{770}\nu_\tau$ ,  $a_1(\pi^-\pi^0\pi^0)_{1260}\nu_\tau$  and  $\pi^-\pi^-\pi^+\pi^0\nu_\tau$  final states contribute to the  $a_1(1260)$  resonance, although only  $a_1(\pi^-\pi^0\pi^0)_{1260}\nu_\tau$  final has a real  $a_1(1260)$  resonance. This is due to the structure of the  $\chi^2_{a_1(1260)}$  minimisation function allowing final states with more than two photons and one  $\pi^\pm$  to contribute.

Last six variables in the list help to differentiate an electron final state to that of a charged pion. A charged pion that starts showering early in the calorimeter could have a

similar topology to an electromagnetic shower. Nevertheless, a good separation between the two can be achieved with the help of these variables.

## 6.6 Multivariate Analysis

For the multivariate analysis, the multiclass class of the TMVA package [43] was used to perform a multiclass classification, which trains the seven final states simultaneously. The multiclass class is an extension of the standard two-class signal-background classifier.

There are two ways for the training. "One v.s. one" is each class is trained against each other class. And the overall likelihood is normalised. The second way to train is called "one v.s. all", which is when each class is trained against all other classes.

Using a three-class example, A, B and C, "one v.s. one" scheme trains A against B, B against C, and C against A. Then the likelihood is normalised. "One v.s. all" would train A against B plus C, B against A plus C, and C against A plus B.

TMVA multiclass implementation uses "one v.s. all" scheme. For each final state, the multiclass classifier will train the final state as the signal against all other final states as the background. This process is repeated for each final state. The classifier output for a single event is a normalised response for each final state, where the sum is one. The response of each final state of a event can be treated as the likelihood. The event is classified into a particular final state if the final state has the highest classifier output response. The advantage of using the multiclass is that the correlation between different final states are accounted for and the classifier output are correctly adjusted for multiple final states, hence one event can only be classified into one final state. The issue with the multiclass is that discriminative variables for each final state need enter the training stage, resulting in a large number of variables.

Half of the randomly selected samples were used in the training process and the other half were used for testing.

The TMVA multiclass classifier used is boosted decision tree with gradient boosting (BDTG), as it was found to give for the best performance. The MVA classifier is trained and optimised to give the best overall separation across all final states. MVA will be discuss further in section ??

## 6.7 Result

Reco ↓ True →	$e^- \bar{\nu}_e$	$\mu^- \bar{\nu}_\mu$	$\pi^-$	$\rho(\pi^-\pi^0)$	$a_1(\pi^-\pi^0\pi^0)$	$a_1(\pi^-\pi^-\pi^+)$	$\pi^- \pi^- \pi^+ \pi^0$
$e^- \bar{\nu}_e$	<b>99.8</b>	-	0.9	1.1	0.8	-	-
$\mu^- \bar{\nu}_\mu$	-	<b>99.5</b>	0.5	-	-	-	-
$\pi^-$	-	0.3	<b>93.2</b>	0.9	-	0.4	-
$\rho(\pi^-\pi^0)$	-	-	4.1	<b>93.0</b>	10.5	0.6	2.8
$a_1(\pi^-\pi^0\pi^0)$	-	-	-	4.3	<b>88.2</b>	-	1.0
$a_1(\pi^-\pi^-\pi^+)$	-	-	1.0	0.3	-	<b>96.6</b>	6.9
$\pi^- \pi^- \pi^+ \pi^0$	-	-	-	0.4	0.4	2.4	<b>89.3</b>

**Table 6.2:** The percentage of reconstructed decay modes corresponds to underlying true decay modes, with  $\sqrt{s} = 100$  GeV for nominal CLIC\\_ILD detector model. Bold numbers show the correctly reconstructed percentages. Numbers less than 0.25% are not shown. Statistical uncertainties are less than 0.25%. Final states include  $\nu_\tau$ , which is not shown.

The reconstruction efficiencies for the seven final state of the tau decaying with c.o.m. energy of 100 GeV for the nominal CLIC\\_ILD detector are shown in table 6.2. The perfect reconstruction would result in only terms in the diagonal.

The unprecedented high classification rate has been achieved. The improvement of photon reconstruction described in section ?? improved the ability to separate 1-prong final state. Most notably, figure ?? shows number of photons have a high correct reconstruction efficiency.

For leptonic decay, the selection efficiency is above 99.5% as the tracking system have much better resolution than the calorimeter.

The  $\mu^- \bar{\nu}_\mu$  final state has very clear topology, as muon deposits energy in the muon chamber. Therefore, there is little confusion with other final states.

$e^- \bar{\nu}_e$  final state is well separated, due to the specialised variables aimed to differentiate early hadronic shower to electromagnetic shower. However, there is still about 1% confusion in one prong final state.

For one prong final states,  $\pi^-$ ,  $\rho(\pi^-\pi^0)$ , and  $a_1(\pi^-\pi^0\pi^0)$ , the confusion is mainly due to the imperfect separation of nearby photons, originated from  $\pi^0$ .

Similarly the confusion between 3-prong final state,  $a_1(\pi^-\pi^-\pi^+)$ , and  $\pi^- \pi^- \pi^+ \pi^0$  is caused by the inability to resolve photon pairs.

## 6.8 Electromagnetic calorimeter optimisaiton

As discussed above, the tau decay mode separation is an benchmark test of detector performance. The ability to resolve photon pairs is crucial to separate different 1-prong states, and different 3-prong state. One of the main feature of calorimeter design affecting the photon resolution is the size of electromagnetic calorimeter (ECal) cell for the high granular calorimeter. The finer ECal cell size is, the better resolution of reconstructing individual photons.

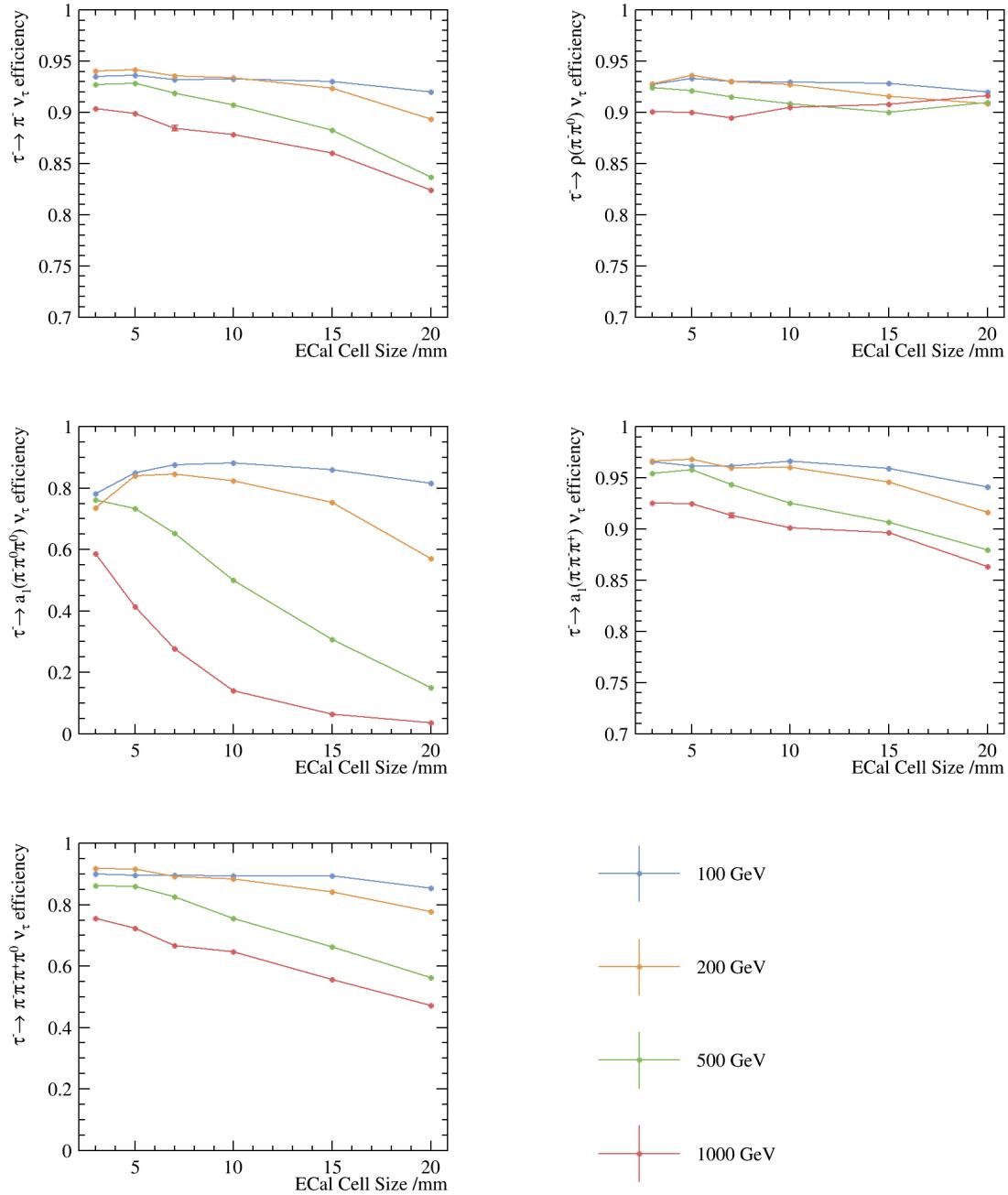
The classification is being tested with the impact of different  $\sqrt{s}$  and different ECal square cell sizes. Around two million events were simulated at each  $\sqrt{s} = 100, 200, 500$  and  $1000\text{ GeV}$ , with each different ECal square cell sizes of 3, 5, 7, 10, 15 and 20 mm. Events were simulated and reconstructed in the same way as described above, with same selection applied. MVA classifier was trained individually for each  $\sqrt{s}$  and each ECal square cell size, with same set of discriminative variables.

To access the impact different ECal square size on detector performance, in particular ECal performance, the correct reconstruction efficiency for 1-prong and 3-prong final states is used as metric. The higher  $\sqrt{s}$  of the collision would degrade the performance, as photons are more boosted and more difficult to resolve.

The leptonic decay correct reconstruction efficiency is not used as a metric as they are similar across different ECal cell sizes. This is because the  $e^\pm$  and  $\mu^\pm$  identifications mostly rely on the tracking system, which was not varied in this study. The energy deposited in the calorimeter are used for the association to the tracks but it has a small impact on the lepton identification.

figure 6.2 shows that as the ECal cell sizes increase, the reconstruction efficiencies generally decrease. Larger cell sizes have lower spatial resolutions, making the separating of nearby photons more difficult.

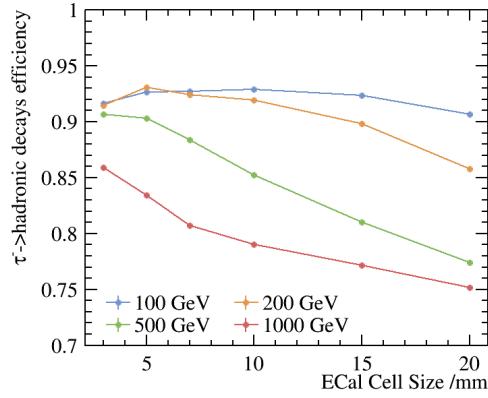
For the  $a_1(\pi^-\pi^0\pi^0)_{1260}\nu_\tau$  final state, the selection efficiency for  $500\text{ GeV}$  rises from ECal cell sizes 15 mm to 20 mm and the one for  $1000\text{ GeV}$  rises from 7 , to 20 mm actually goes up as cell size increases. This is because when the algorithm can not reconstruct



**Figure 6.2:** The selection efficiencies for various final states against the ECal cell size for different c.o.m. energies with the nominal CLIC\_IID detector model are shown. The top left, top right, middle left, middle right and bottom left plots are for the  $\pi^- v_\tau$ ,  $\rho(\pi^-\pi^0)_{770} v_\tau$ ,  $a_1(\pi^-\pi^0\pi^0)_{1260} v_\tau$ ,  $a_1(\pi^-\pi^-\pi^+\pi^0)_{1260} v_\tau$  and  $\pi^-\pi^+\pi^-\pi^0 v_\tau$  final states respectively. From the top to the bottom, blue, orange, green and red lines are representing the  $\sqrt{s} = 100, 200, 500$  and  $1000$  GeV respectively.

four photons in the  $a_1(\pi^-\pi^0\pi^0)_{1260}\nu_\tau$  final state, and the event topology would be very similar to the  $\rho(\pi^-\pi^0)_{770}\nu_\tau$  final states.

For the  $\sqrt{s} = 100$  and  $200$  GeV, the selection efficiency of the 5 mm ECal cell size is better than that of the 3 mm. One possible explanation is that the and the PandoraPFA have been optimised for the nominal ILD detector with the 5 mm ECal cell size, which shares the same ECal structure with the nominal CLIC\_ILD detector.



**Figure 6.3:** The  $\tau$  hadronic decay efficiency against the ECal cell size for different  $\sqrt{s} = e$  TeV energies with the nominal CLIC\_ILD detector model are shown. The blue, orange, green and red lines are representing the  $\sqrt{s} = 100, 200, 500$  and  $1000$  GeV respectively.

To effectively compare the overall separation power of all hadronic final states across  $\sqrt{s}$  and ECal square cell sizes, we constructed a single parameter function, the  $\tau$  hadronic decay final state efficiency function,

$$\varepsilon_{\text{had}} = \frac{\sum_i B r_i \varepsilon_i}{\sum_i B r_i}, \quad (6.3)$$

where  $B r_i$  is the branching fraction of a hadronic final state after the generator level cut.  $\varepsilon_i$  is the correct reconstruction efficiency of the final state, and the  $i$  is summing over five hadronic decay final state of  $\tau$ . Leptonic decays,  $e^- \bar{\nu}_e \nu_\tau$  and  $\mu^- \bar{\nu}_\mu \nu_\tau$ , were not included, because the variation of the leptonic decay selection efficiency is small.

In the figure 6.3,  $\tau$  hadronic decay final state efficiency,  $\varepsilon_{\text{had}}$ , against the ECal cell size with different  $\sqrt{s}$  is shown.  $\varepsilon_{\text{had}}$  decreases when cell sizes increases and when  $\sqrt{s}$  increases.  $\varepsilon_{\text{had}}$  of the 5 mm ECal cell size is better than that of the 3 mm for  $\sqrt{s} = 100$

and 200 GeV lines possibly due the optimisation of the software for the nominal ILD 5 mm ECal square cell size.

The  $\varepsilon_{\text{had}}$  is above 90% for the ECal cell size from 3 to 20 mm for the  $\sqrt{s} = 100$  GeV. For  $\sqrt{s} = 200$  GeV, the  $\varepsilon_{\text{had}}$  decreases from over 90% to 86% for the ECal square cell size from 3 to 20 mm. The degradation of the  $\varepsilon_{\text{had}}$  is more significant for the  $\sqrt{s} = 500$  and 1000 GeV, where the  $\varepsilon_{\text{had}}$  drops from over 90% to 77% and from 86% to 75% respectively, over the same range of ECal square cell size.

For  $\sqrt{s} = 100$  and 200 GeV, up to 15 mm cell sizes of ECal will give a good performance for  $\tau$  hadronic decay modes separation, and the  $\varepsilon_{\text{had}}$  is above 90%. For  $\sqrt{s} = 500$  and 1000 GeV, it is preferential to have a small ECal cell size for a good  $\tau$  hadronic decay modes separation. There is about 15% degradation of  $\varepsilon_{\text{had}}$  for ECal square cell size from 3 to 20 mm.

## 6.9



# Chapter 7

## Double Higgs Bosons Production Analysis

*“Two is better than one”*

— Sir Steve Orange, 1785–1854

Since the discovery of Higgs boson in the LHC in 2012 [1, 16], it is crucial to understand the properties of the Higgs boson and test if it is a Standard Model Higgs. The Higgs mechanism and the Higgs boson in the Standard Model have been explained in the chapter 2. A few theories of Higgs beyond the Standard Model could be tested via the double higgs production in an electro-positron collider (section 2.9). A generator level studies has shown that the precision reached by a multi-TeV linear collider, such as the Compact Linear Collider (CLIC), is much better than it by the Large Hadron Collider (LHC) even with  $3000\text{fb}^{-1}$  of data [13].

The first challenge for double Higgs bosons production analysis is that events are rare. The cross section is very small comparing to other background physics processes, making it difficult to select signal events. The second challenges is at high centre-of-mass, events are boosted and many particles are in the forward region of the detector, where the reconstruction performance is worse than the barrel region.

In this chapter, a full CLIC\_ILD detector simulation studies has been performed for the double higgs production,  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ . First event generation and simulation will be briefly discussed. An overview of the analysis, including lepton finding and jet reconstruction, is presented. This is followed by a multivariate analysis with results on

the selection efficiency. Lastly, a uncertainty on extracting Higgs triple self couplings,  $g_{HHH}$ , and quartic coupling,  $g_{WWHH}$ , is presented.

## 7.1 Analysis Straggly Overview

The double higgs production,  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ , can occur via processes in figure 7.1. Figure 7.1a is sensitive to Higgs triple self coupling  $g_{HHH}$ . Figure 7.1b is sensitive to quartic coupling  $g_{WWHH}$ . Figure 7.1c and figure 7.1d are the irreducible background processes for the study of  $g_{HHH}$  and  $g_{WWHH}$ .

The  $e^-e^+ \rightarrow ZHH$ , where  $Z$  decays to  $\nu \bar{\nu}$ , has the same  $HH\nu\bar{\nu}$  final state as in figure 7.1. The  $ZHH$  channel can be used to study the Higgs triple self couplings, and has been studied at ILC for  $\sqrt{s} = 500$  GeV [19]. However, its contribution to the  $HH\nu\bar{\nu}$  final state is small comparing to Feynman diagrams in figure 7.1, for the relevant CLIC energies  $\sqrt{s} = 1.4$  TeV and 3 TeV. Shown in figure 2.3, the cross section for  $e^-e^+ \rightarrow ZHH$  is one order of magnitude smaller than  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  via processes in figure 7.1. Therefore, the effect of  $e^-e^+ \rightarrow ZHH$  present in  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  at  $\sqrt{s} = 1.4$  TeV and 3 TeV is negligible.

The double higgs production,  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  is divided into sub-channel  $HH \rightarrow b\bar{b}W^+W^-$  and  $HH \rightarrow b\bar{b}b\bar{b}$  to allow a more detailed examination, providing a cross checks between two sub-channels and an improvement of signal selection when combined. In this chapter, the W hadronic decay mode of the  $HH \rightarrow b\bar{b}W^+W^-$  channel is studied, and called the signal channel. It is chosen because the hadronic decay has the largest cross section and does not produce neutrinos, which allows each W to be reconstructed using a di-jet. However, hadronic decay final state of the  $HH \rightarrow b\bar{b}W^+W^-$  has a very low cross section. The signal selection is challenging and aggressive background rejection methods are deployed.

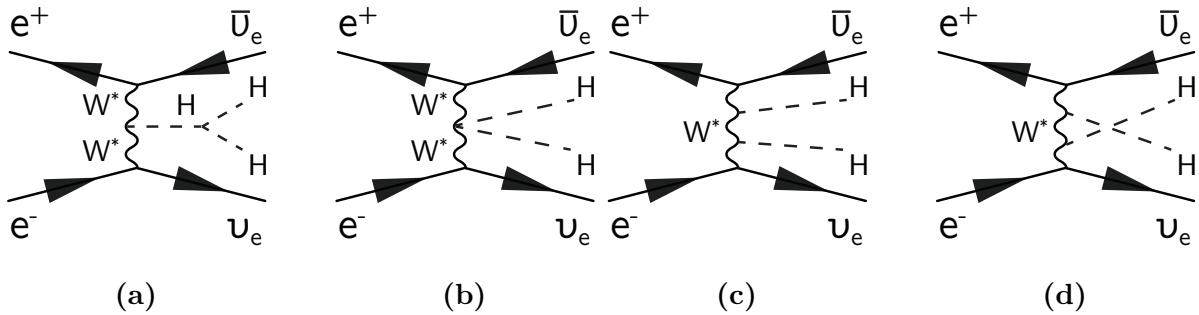
The semi-leptonic final states of the sub-channel  $HH \rightarrow b\bar{b}W^+W^-$  is studied as well. However, extra neutrinos in the final states present greater difficulty to reconstruct the two Higgs bosons, as some momenta of one Higgs boson is missing. This channel would be discussed briefly as it is adapted from the hadronic decay analysis.

$HH \rightarrow b\bar{b}b\bar{b}$  sub-channel is studied independently by collaborators. However there are collaboration between two studies, which would converge on couplings extractions.

The signal channel final state,  $HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e \rightarrow b\bar{b}qqq\bar{q}\nu_e\bar{\nu}_e$ , is a six quark final state with missing momentum. The high number of quarks require efficient jet

reconstruction and jet pairing algorithms to select the signal final states. The two b quarks in the final states allows b jet tagging information to be useful. Since the final state does not contain leptons, event-level lepton finding, typically for energetic isolated leptons, would improve the signal selection efficiency. The topology of the signal events determine the analysis strategy.

Proof-of-principle generator study was performed at CLIC using CLIC\_ILD detector model for  $\sqrt{s} = 1.4$  TeV and 3 TeV [18]. A full CLIC\_ILD detector simulation is presented in this chapter. Firstly suitable signal and background channels are identified. To help selecting the signal, light lepton and tau lepton are found and b-jet tagging is used. An event is grouped into a number of jets depending on the final state, followed by pre-selections cuts and multivariate analysis. The analysis for  $\sqrt{s} = 1.4$  TeV with  $HH \rightarrow b\bar{b}W^+W^-$  hadronic decay will be presented first. The difference in the analysis and results for  $\sqrt{s} = 3$  TeV and for the  $HH \rightarrow b\bar{b}W^+W^-$  semi-leptonic decay will be highlighted afterwards. The simultaneous extraction of couplings, extracting Higgs triple self couplings,  $g_{HHH}$ , and quartic coupling,  $g_{WWHH}$ , are presented in the final section.



**Figure 7.1:** Figures show Feynman diagrams of leading-order  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  processes at CLIC, without considering  $e^-e^+ \rightarrow ZHH$ . Figure 7.1a is sensitive to the Higgs triple self coupling  $g_{HHH}$ . Figure 7.1b is sensitive to quartic coupling  $g_{WWHH}$ .

## 7.2 Monte Carlo Sample Generation

Selected background samples are considered in the analysis and listed in Table 7.1. The signal channel is  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-$  where both W decay hadronically.

Background processes with many quarks with missing energies would be challenging to veto. Two examples are  $e^-e^+ \rightarrow q\bar{q}q\bar{q}\nu\bar{\nu}$  and  $e^-\gamma(BS) \rightarrow \nu q\bar{q}q\bar{q}$ . Single Higgs boson production would be difficult to remove as well, such as  $e^-e^+ \rightarrow q_l\bar{q}_l H\nu\bar{\nu}$ .

Some background channels are not considered because either have different different event topologies, or they have very small cross sections. For example  $e^\pm \gamma() \rightarrow q\bar{q}H\ell$  is ignored as the cross section is very small even at  $\sqrt{s} = 3$  TeV (0.07 fb with photon from EPA, 0.6 fb with photon from BS). Other background channels are not simulated due to computational limitations. For example, six-quark final states were not simulated due to constraint of the simulating software.

Electron-photon interactions are considered in this analysis. They are important as the interaction becomes significant with high  $\sqrt{s}$ . Processes initiated by photons included, where photons are produced due to the high electric field generated by the colliding beams. Processes involving real photons from beamstrahlung (BS) and “quasi-real” photons are generated separately. For the “quasi-real” photon initiated processes, the Equivalent Photon Approximation (EPA) has been used.

For processes involving Higgs production explicitly, simulated Higgs mass is 126 GeV. Multi-quark final state background samples could, in principle, contain higgs production. Therefore, they are generated with a Higgs mass of 14 TeV. This will produce negligible double higgs production cross section. The cross section for sub-channels of the signal, such as  $HH \rightarrow b\bar{b}W^+W^-$ , are corrected manually according to [44], as the internal value of the event generation software (PYTHIA) is inaccurate.

The simulation and reconstruction chain are described in chapter 5. For most background processes, events are simulated when invariant mass of quarks are above 50 GeV. For electron-photon interaction with  $qqqq\nu$  final state, events are simulated when invariant mass of quarks are above 120 GeV. These limits are necessary to generate a large amount of background samples in a feasible time, without losing much signal samples.

Finally, the main beam induced background  $\gamma\gamma \rightarrow \text{hadrons}$  is simulated and overlayed to all samples according to the integration time of each subdetector. Details can be found in section ??.

## 7.3 Lepton identification

The reconstruction is done via Marlin in iLCSoft v01-16. Separate software packages (processors) exist for identification of electrons, muons, and taus, and for jet reconstruction.

Channel	$\sigma(\sqrt{s} = 1.4 \text{ TeV}) / \text{fb}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$	0.149
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-$ , hadronic	0.018
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	0.047
$e^-e^+ \rightarrow HH \rightarrow \text{others}$	0.085
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	0.86
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	0.36
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	0.31
$e^-e^+ \rightarrow qqqq$	1245.1
$e^-e^+ \rightarrow qqqq\ell\ell$	62.1*
$e^-e^+ \rightarrow qqqq\ell\nu$	110.4*
$e^-e^+ \rightarrow qqqq\nu\bar{\nu}$	23.2*
$e^-e^+ \rightarrow qq$	4009.5
$e^-e^+ \rightarrow qq\ell\nu$	4309.7
$e^-e^+ \rightarrow qq\ell\ell$	2725.8
$e^-e^+ \rightarrow qq\nu\bar{\nu}$	787.7
$e^-\gamma(\text{BS}) \rightarrow e^-qqqq$	1160.7
$e^+\gamma(\text{BS}) \rightarrow e^+qqqq$	1156.3
$e^-\gamma(\text{EPA}) \rightarrow e^-qqqq$	287.1
$e^+\gamma(\text{EPA}) \rightarrow e^+qqqq$	286.9
$e^-\gamma(\text{BS}) \rightarrow \nu qqqq$	79.8†
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu} qqqq$	79.3†
$e^-\gamma(\text{EPA}) \rightarrow \nu qqqq$	17.4†
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu} qqqq$	17.3†
$e^-\gamma(\text{BS}) \rightarrow qqH\nu$	15.8*
$e^+\gamma(\text{BS}) \rightarrow qqH\nu$	15.7*
$e^-\gamma(\text{EPA}) \rightarrow qqH\nu$	3.39*
$e^+\gamma(\text{EPA}) \rightarrow qqH\nu$	3.39*
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow qqqq$	21406.2*
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow qqqq$	4018.7*
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow qqqq$	4034.8*
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow qqqq$	753.0*

**Table 7.1:** List of signal and background samples with the corresponding cross sections at  $\sqrt{s} = 1.4 \text{ TeV}$ .  $q$  can be  $u$ ,  $d$ ,  $s$ ,  $b$  or  $t$ . Unless specified,  $q$ ,  $\ell$  and  $\nu$  represent particles and its corresponding anti-particles.  $\gamma$  (BS) represents a real photon from beamstrahlung (BS).  $\gamma$  (EPA) represents a “quasi-real” photon, simulated with the Equivalent Photon Approximation. For processes involving Higgs production explicitly, simulated Higgs mass is 126 GeV. Otherwise, Higgs mass is set to 14 TeV. For processes labeled with \* and †, the generator level cut requires invariant mass of quarks greater than 50 and 120 GeV, respectively.

New processors have been developed and existing processors have been optimised for signal selection and background rejection. The latest functioning flavour tagging processor exist in iLCSoft v01-16. Thus newer versions of iLCSoft can not be used in this analysis.

For the signal channel,  $\text{HH} \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}qqqq$ , there is no lepton in the final state, whilst many background final states contain primary leptons, such as  $qqqq\ell\nu$ . Hence effective vetoing events with leptons would improve the signal selection efficiency. Leptons from channels like  $qqqq\ell\nu$  are primary, where their tracks are very close to the interaction point. These leptons are typically energetic, and often isolated from other particles. Whilst electrons and muons are stable enough to deposit energies in calorimeters, tau lepton is very short lived and decays in the tracker. Therefore only tau decay products can be reconstructed. These characteristics provide the basis for the isolated lepton finding.

### 7.3.1 Electron and muon identification

#### IsolatedLeptonFinderProcessor

IsolatedLeptonFinderProcessor in Marlin package is used, modified, and optimised. This processor identifies high energy electrons and muons that are away from other particles, “isolated”. The optimal parameters below are chosen in collaboration and tested using the signal channel and the  $e^-e^+ \rightarrow qqqq\ell\nu$  channel.

Electrons induced electromagnetic showers are mostly contained in the ECAL, while muons would deposit some energies in the ECAL. The inner detector tracks of electrons and muons from the electron-positron interaction are primary tracks, which are very close to the interaction point. The isolation criteria requires the lepton to be far away from other high energy particles. These form the logics for the IsolatedLeptonFinderProcessor. The value for the selection cut is listed in table 7.2.  $E_{\text{ECAL}}$  is the energy deposited in the ECAL.  $E_{\text{cone}}$  is the total energy of PFOs within a cone of an opening angle of  $\cos^{-1}(0.995)$  around the lepton.  $d_0$ ,  $z_0$ , and  $r_0$  are the Euclidean distance of the track starting point to the interaction point in the x-y plane, the in z direction, and in the x-y-z three dimensional space. The performance of the processor is shown in table 7.6.

IsolatedLeptonFinderProcessor	Selection
High Energy (GeV)	$E > 15$
Energy deposition $e^\pm$	$\frac{E_{ECAL}}{E} > 0.9$
Energy deposition $\mu^\pm$	$0.25 > \frac{E_{ECAL}}{E} > 0.05$
Primary Track (mm)	$d_0 < 0.02, z_0 < 0.03, r_0 < 0.04$
Isolation (GeV)	$E_{cone}^2 \leq 5.7 \times E - 50$

**Table 7.2:** Optimised selection criterion for IsolatedLeptonFinderProcessor

### BonoLeptonFinderProcessor

The analysis would benefit from aggressive lepton veto due to the low signal cross section. Since the IsolatedLeptonFinderProcessor is conservative, an aggressive light lepton selection processor, BonoLeptonFinderProcessor, is developed. It utilises calorimetric information provided by PandoraPFA.

The processor uses two sets of cuts and the logic is similar to ones in the IsolatedLeptonFinderProcessor. The first set of cuts uses the particle ID information from PandoraPFA. The cuts require a PandoraPFA electron or muon with high  $p_T$  and is from a primary track. The lepton should either have very high transverse momentum, or it is isolated. The second set of cuts uses ECAL energy fraction to determine light leptons. The rest of the logic is very similar to the first set. The isolation criterion is stricter than it in the first set to reduce fake rate.

Table 7.3 lists values for the selection cut. Variables are defined similarly as those in section 7.3.1.  $p_T$  is the transverse momentum.  $E_{cone1}$  and  $E_{cone2}$  are the total energy of PFOs within a cone of an opening angle of  $\cos^{-1}(0.995)$  and  $\cos^{-1}(0.99)$  respectively around the lepton. The performance of the processor is shown in table 7.6.

### Comparison: IsolatedLeptonFinderProcessor v.s. BonoLeptonFinderProcessor

Two processors share similar criterion for light lepton identification. The main difference is that the BonoLeptonFinderProcessor uses the particle identification from PandoraPFA, use extra calorimetric information to determine the particle ID than simple ECAL energy fraction. BonoLeptonFinderProcessor also allows high  $p_T$  light lepton to be identified

BonoLeptonFinderProcessor	Selection
High Energy (GeV)	$E > 10$
Energy deposition $e^\pm$	PandoraPFA reconstructed & $\frac{E_{ECAL}}{E} > 0.95$
Energy deposition $\mu^\pm$	PandoraPFA reconstructed
Primary Track (mm)	$r_0 < 0.015$
a) High Transverse Momentum (GeV)	$p_T > 40$
b) Isolation (GeV)	$E \geq 23 \times \sqrt{E_{cone1}} + 5$
High Energy (GeV)	$E > 10$
Energy deposition $e^\pm$	$\frac{E_{ECAL}}{E} > 0.95$
Energy deposition $\mu^\pm$	$0.2 > \frac{E_{ECAL}}{E} > 0.05$
Primary Track (mm)	$r_0 < 0.5$
a) High Transverse Momentum (GeV)	$p_T > 40$
b) Isolation (GeV)	$E \geq 28 \times \sqrt{E_{cone2}} + 30$

**Table 7.3:** Optimised selection criterion for IsolatedLeptonFinderProcessor

in a potential non-isolated environment, which leads to the more aggressive nature of the BonoLeptonFinderProcessor. The performance of two processors on the signal and selected background samples is shown in table 7.6.

### 7.3.2 Tau identification

#### TauFinderProcessor

With a decay length of  $87\mu\text{m}$ , tau leptons decay before reaching the detector and can only be identified through the reconstruction of their decay products. The leptonic decay of tau can be identified using the two isolated lepton finder processors in section 7.3.1 and section 7.3.1. Therefore tau identification will focus on the hadronic decay.

The basic logic of a tau finder is similar to a cone clustering algorithm (section 4.4.3). A high energy track is selected as a seed and a small cone is formed around it. The PFOs inside the cone are required to be consistent with a tau hadronic decay: no more than 3 charged particles, invariant mass close to tau mass and few PFOs in the cone. The cone then forms the tau candidate. Like the lepton in the isolated light lepton finder, the tau candidate is required to be isolated from other particles. To reduce fake rate,

low momentum and very forward particles do not participate in the tau finding, as they more likely come from  $\gamma\gamma \rightarrow \text{hadrons}$  background.

TauFinderProcessor [45], an existing processor Marlin package, has been tuned in collaboration and tested. The selection criterion are listed in table 7.4. Variables are defined similarly as those in previous sections.  $\theta_Z$  is the polar angle w.r.t. the beam axis.  $N_{X^+} > 3$  and  $N_{\text{cone}}$  are number of charged particles and number of PFOs respectively in the tau candidate.  $m_{\text{cone}}$  is the invariant mass of the sum of the PFOs, tau candidate.  $E_{\text{cone}}$  are the total energy of PFOs within a cone of an opening angle between 0.03 and 0.33 rad around the tau seed. The performance of the processor is shown in table 7.6.

TauFinderProcessor	Selection
Veto $\gamma\gamma \rightarrow \text{hadrons}$ (GeV)	$p_T < 1,  \cos(\theta_Z)  > 1.1$
Seed particle (GeV)	$p_T > 10$
Tau candidate cone opening angle (rad)	0.03
Tau candidate rejection	$N_{X^+} > 3, N_{\text{cone}} > 10, m_{\text{cone}} > 2$
Isolation (GeV)	$E_{\text{cone}} < 3$

**Table 7.4:** Optimised selection criterion for TauFinderProcessor

### BonoTauFinderProcessor

For the similar reason of developing a more aggressive light lepton finder, a new more aggressive tau lepton selection processor, BonoTauFinderProcessor, is developed. BonoTauFinderProcessor utilises calorimetric information provided by PandoraPFA and has a cleaner code structure.

Similar to TauFinderProcessor, this processor identifies a high momentum particle as a tau seed. Particles are iteratively added to the search cone according to the size of the opening angle to the seed. After each particle addition, a temporary search cone is considered as the tau candidate, and tested against tau hadronic decay hypothesis and isolation conditions. The tau candidate only needs to pass one of the isolation conditions. The iterative particle addition would stop when the cone opening angle is bigger than a threshold. If multiple temporary search cones of a same seed passing the selection, the cone with smallest opening angle is chosen to form the final tau candidate. To reduce fake tau decay products from  $\gamma\gamma \rightarrow \text{hadrons}$  background, low energy particles do not participate in the tau finding.

Table 7.4 lists the selection criterion for BonoTauFinderProcessor. Variables are defined similarly as those in previous sections.  $\theta_S$  is the opening angle of the search cone. `cone1` and `cone2` are defined as a cone of an opening angle of  $\cos^{-1}(0.95)$ , and  $\cos^{-1}(0.99)$  respectively around the tau seed.  $r_0$  is referring to tau seed particle. The performance of the processor is shown in table 7.6.

TauFinderProcessor	Selection
Veto $\gamma\gamma \rightarrow \text{hadrons}$ (GeV)	$E < 1$
Seed particle (GeV)	$p_T > 5$
Maximum search cone opening angle (rad)	$\theta_S \leq \cos^{-1}(0.999)$
Tau candidate rejection	$N_{X^+} \neq 1, 3, m_{\text{PFO}} > 3$
Isolation (GeV) 1	$N_{\text{cone1}} = 0, p_{T\text{cone}} \geq 10$
Isolation (GeV) 2	$N_{X^+} = 1, N_{\text{cone1}} = 1, r_0 > 0.01$
Isolation (GeV) 3	$N_{X^+} = 3, N_{\text{cone1}} = 1, p_{T\text{cone}} \geq 10, \theta_S < \cos^{-1}(0.9995)$
Isolation (GeV) 4	$N_{X^+} = 1, N_{\text{cone2}} = 0, r_0 > 0.01, p_{T\text{cone}} \geq 10$
Isolation (GeV) 5	$N_{X^+} = 3, N_{\text{cone2}} = 0, p_{T\text{cone}} \geq 10, \theta_S < \cos^{-1}(0.9995)$

**Table 7.5:** Optimised selection criterion for BonoTauFinderProcessor

### Comparison: TauFinderProcessor v.s. BonoTauFinderProcessor

Two processors share similar size of search cone and isolation cone. The main difference is that the BonoTauFinderProcessor has an interactive approach to build up tau candidate, which allows a dynamic tau search cone size. The BonoTauFinderProcessor has looser cut on minimum  $p_T$  and invariant, but stricter isolation criterion. This leads to a more aggressive tau finder. The performance of processors is shown in table 7.6.

### 7.3.3 Very forward electron identification

At high  $\sqrt{s}$ , particles are boosted and it is important to extract information in the forward calorimeters to aid signal selection. Certain background channels, for example photon-electron interactions, can have energetic electrons in the forward calorimeters,

the LumiCAL and the BeamCAL. These extracting information from these forward calorimeters is different to information in the barrel and the end cap. As forward calorimeters have the angular low angular acceptance, most particles in these forward detector would be very forward particles from beam induced background. However, previous study [28] has shown that sufficiently high energy electrons can be efficiently identified in the BeamCAL and LumiCAL.

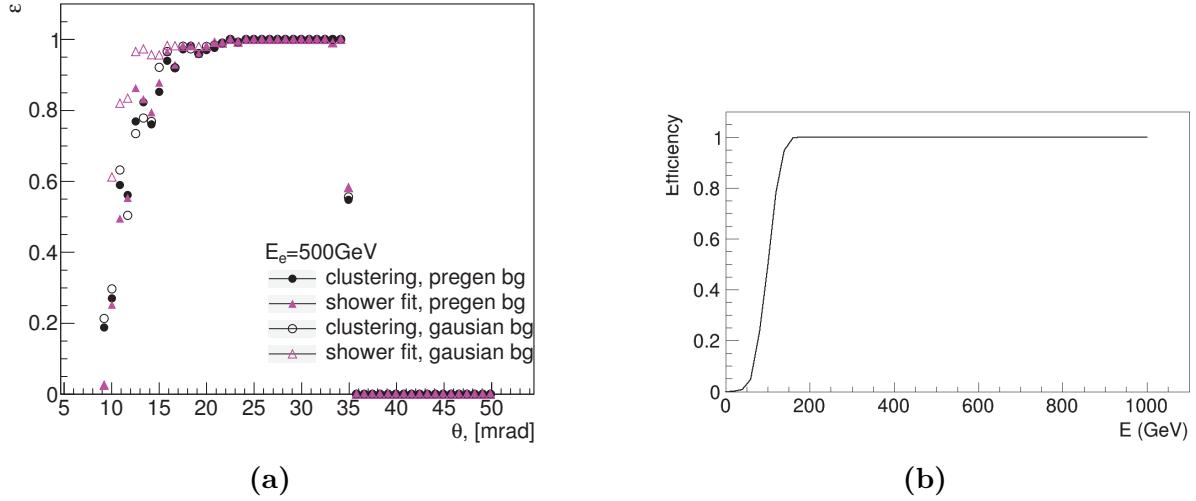
For the CLIC\_ILD detector concept, energies deposited in the LumiCAL and the BeamCAL are not simulated. This is because the thousands of beam induced background particles per bunch corssing () (see section ??) requires expensive computational resources. The current simulation software could not handle the simulation in a feasible time. Therefore, the adopted approach is to parameterise the background particle energy deposition and leads to electron tagging algorithms.

For the BeamCAL, [46] describes an electron tagging algorithm developed using  $\sqrt{s} = 3 \text{ TeV}$  collision environment, by comparing the simulated electron and background energy distributions. An electron is tagged if the energy is significantly larger than the expected background energy distributions.

Events are overlaid with background energy deposition integrated over 40 bunch crossing. A C++ code library has been developed and used in this analysis. The tagging efficiency for electrons with energy 500 to 1500 GeV are bin in histograms at interval of 100 GeV. There is no tagging for electrons with energy below 500 GeV or about 1500 GeV. An indicative performance plot of 500 GeV electron tagging efficiency as a function of the polar angle is shown in figure 7.2a.

The input of the BeamCAL electron tagging algorithm is the four momenta of the MC electron. Since the algorithm assumes collision at  $\sqrt{s} = 3 \text{ TeV}$ , for the  $\sqrt{s} = 1.4 \text{ TeV}$  user case, the momenta of the MC electron is scaled down by a factor of  $\frac{3}{1.4}$  to use the algorithm.

For the LumiCAL, the  $H \rightarrow \mu\mu$  analysis in [48, 49] has developed an algorithm for electron tagging in the LumiCAL, with similar logic as the algorithm for the BeamCAL. Figure 7.2 shows the LumiCAL electron tagging efficiency as a function of the electron energy, for polar angle  $\theta = 50 \text{ mrad}$ , where events are overlaid with background energy deposition integrated over 100 bunch crossings. As figure 7.2 is the only available performance plot for the LumiCAL electron tagging, the LumiCAL electron tagging in this analysis is based on the plot.



**Figure 7.2:** Figure 7.2a shows BeamCAL 500 GeV electron tagging efficiency as a function of the polar angle with different methods to model backgrounds and fittings, taken from [46]. Figure 7.2b shows the LumiCAL electron tagging efficiency as a function of the electron energy, for polar angle  $\theta = 50$  mrad, taken from [47].

Assuming that the LumiCAL electron tagging efficiency is described as in figure 7.2, for all polar angles and for  $\sqrt{s} = 1.4$  TeV and 3 TeV, LumiCAL electrons tagging efficiency,  $\varepsilon$ , is parameterised as

$$\varepsilon = \begin{cases} 0, & \text{if } E < 50 \text{ GeV} \\ 0.99 \times \frac{\text{erf}(E-100)+1}{2}, & \text{otherwise} \end{cases} \quad (7.1)$$

where  $E$  is the energy of the electron or the photon.  $\text{erf}$  is the error function. For each MC electron in the LumiCAL, a random number between 0 and 1 is generated. If the random number is less than  $\varepsilon$ , then the MC electron is tagged.

Due to lack of tracking ability in the forward region (see section ?? and section ?? for detector coverage), electrons and photons would have the same electromagnetic shower profile (see section ??), with the ECAL spatial resolution. Therefore, MC photons and electrons appear indistinguishable to the BeamCAL and LumiCAL and both are tagged by the above algorithms.

### 7.3.4 Lepton identification performance

The performance of all lepton finding processors on the signal and selected background samples is shown in table 7.6. The percentages are the fraction of events after rejecting

Efficiency (1.4 TeV)	Signal	$e^+e^- \rightarrow qqqq\ell\nu$
IsolatedLeptonFinderProcessor	99.3%	50.3%
BonoLeptonFinderProcessor	99.1%	39.9%
TauFinderProcessor	97.5%	52.3%
BonoTauFinderProcessor	89.7%	38.5%
ForwardFinderProcessor	98.9%	95.1%
Combined	86.6%	16.8%

**Table 7.6:** Isolated lepton finder processors performance on the signal and selected background samples at  $\sqrt{s} = 1.4$  TeV.

Selection / Efficiency (1.4 TeV)	Signal	$e^-\gamma(BS) \rightarrow e^-qqqq$
Combined light lepton finder	87.6%	67.5%
ForwardFinderProcessor	98.9%	53.6%
Combined	86.6%	30.8%

**Table 7.7:** Very forward electron and photon finder performance on the signal and selected background samples at  $\sqrt{s} = 1.4$  TeV.

events with lepton identified. BonoLeptonFinderProcessor and BonoTauFinderProcessor are more aggressive at rejecting background than the processors IsolatedLeptonFinderProcessor and TauFinderProcessor. By combining the processors, 86.6% signal events remain and 16.8% of  $e^+e^- \rightarrow qqqq\ell\nu$  events survive after rejecting events with lepton identified.

The ForwardFinderProcessor is most effective at rejecting backgrounds with leptons in the forward region. Table 7.7 shows the performance for the signal and the  $e^-\gamma(BS) \rightarrow e^-qqqq$  background. 53.6% of  $e^-\gamma(BS) \rightarrow e^-qqqq$  background are left after the lepton veto.

### 7.3.5 Other lepton identification processors

Other isolated lepton selection processors available in Marlin package, including IsolatedLeptonTagging and TauJetClustering, have been tested. The results, after tuning of parameters, were unsatisfactory. They either performed poorly comparing to the

processors above, or became redundant after using the processors above. Therefore, these other lepton selection processors are not used in this analysis.

## 7.4 Jet reconstruction

For the signal channel,  $\text{HH} \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}q\bar{q}q\bar{q}$ , one Higgs boson decays to two b quarks, resulting in two jets from hadronisation. Similarly the other Higgs boson decays to two W bosons, where each W boson decays into two quarks. Therefore, the expected number of jets is six. Physical bosons, W and H, can be reconstructed from suitable jets, which allows information to be extracted about the signal channel. Therefore it is important to have efficient jet reconstruction to maximum information extraction. In this section, the optimisation of the jet reconstruction is discussed.

### 7.4.1 Jet reconstruction optimisation

A overview of the jet algorithms can be found in section ???. Longitudinal invariant,  $k_t$ , jet algorithm was chosen for the jet clustering. The free parameters for  $k_t$  algorithm is the R parameter, which controls the radius of the jet. The other parameter to optimised regularise There is also the choice of the PFO collection, which incorporate different level of time and  $p_T$  cuts, to reduce beam induce background (see section ??). Both parameters are optimised for  $\sqrt{s} = 1.4 \text{ TeV}$  and  $\sqrt{s} = 3 \text{ TeV}$ .

The metric for optimising the R parameter, and the PFO collection, is the invariant mass resolution of H and W. With the signal events, jets will be paired to give physical bosons using cheated MC truth information. Hence invariant mass resolution of H and W are obtained.

The sample for the optimisation is  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ . The signal channel, hadronic decay of  $\text{HH} \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}q\bar{q}q\bar{q}$ , is chosen using the MC truth information by examining the decay chain of MC particles. The signal event is then processed through  $k_t$  jet algorithm in 6-jet exclusive mode. The six jets are paired up using MC truth information (see section ??), to the corresponding Higgs and W boson. Four invariant mass distributions are obtained: two Higgs masses,  $m_{H_{bb}}$ ,  $m_{H_{WW^*}}$ , and two W masses  $m_W$ ,  $m_{W^*}$ .  $W^*$  indicates the off-mass-shell W boson, because when a Higgs decays into

two W bosons, one W is off the mass shell, as the Higgs mass is less than the sum two W masses (see section 2.8).

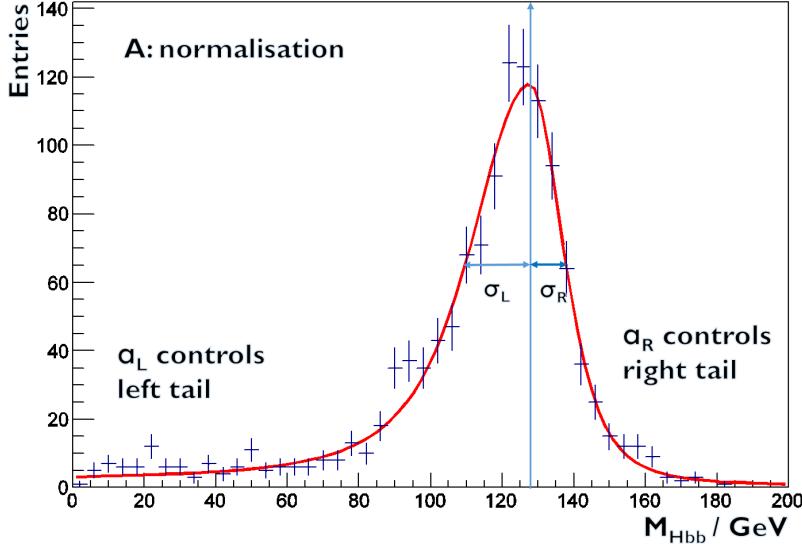
**Mass resolution fit** Three invariant mass resolutions are worth comparing for different jet reconstruction,  $m_{H_{bb}}$ ,  $m_{H_{WW^*}}$ , and  $m_W$ . The optimal jet reconstruction should produce the a sharp mass peak around the particle's simulated mass (see section ??). An example of  $m_{H_{bb}}$  invariant mass distribution is shown in figure 7.3. To quantitatively access the mass distribution resolution, a functional form is fitted. The base fitting function is a gaussian function. Although the underlying mass distribution of particles like  $m_W$  is a Breit-Wigner distribution, as the detector resolution is worse than the W width, the overall mass distribution, which is a convolution between a narrow W Breit-Wigner distribution and a wide gaussian distribution for the detector resolution, is more gaussian like. The  $m_{H_{bb}}$  mass distribution is gaussian like, but with asymmetrical width. This is due to b quarks decaying to neutrinos, leading to a loss of detectable particles and loss of momentum. Therefore there are more events with lower invariant mass and thus asymmetrical width is obtained. As only the peak region of the mass distribution is Gaussian like, tail parameters are added to the fitting function in order to fit the whole range of the mass distribution. The fitting function takes the form of

$$f(m) = A e^{-\frac{(m-\mu)^2}{g}} \begin{cases} g = 2\sigma_L + \alpha_L(m - \mu), & \text{if } m < \mu, \\ g = 2\sigma_R + \alpha_R(m - \mu), & \text{if } m \geq \mu, \end{cases} \quad (7.2)$$

where  $m$  is binned mass distribution, with 50 bins in range [0, 200] GeV.  $\mu$  is the fitted mass peak.  $\sigma_L$  and  $\sigma_R$  allow asymmetrical width of the distribution.  $\alpha_L$  and  $\alpha_R$  control the fit of tails.  $A$  is the normalisation factor.

**Optimal R and PFO collection** The mass fit is performed for  $m_{H_{bb}}$ ,  $m_{H_{WW^*}}$ , and  $m_W$  distributions. The optimal jet reconstruction should have the mass peak close to the particle's simulated mass (see section ??), and a narrow peak width. Due to the asymmetrical fit, the overall relative width is  $(\sigma_L + \sigma_R)/M$ . Smaller width indicates better mass resolution. This mass fit is repeated for reconstructions with  $R$  values of 0.5 to 1.3 at interval of 0.1, and with three PFO collections: loose, normal, and tight (see section 4.8).

Figure 7.4 shows the mass peak and the relative width as a function of  $R$  and PFO collections, for  $m_{H_{bb}}$ ,  $m_{H_{WW^*}}$ , and  $m_W$ . The mass peak value increases as  $R$  increase.

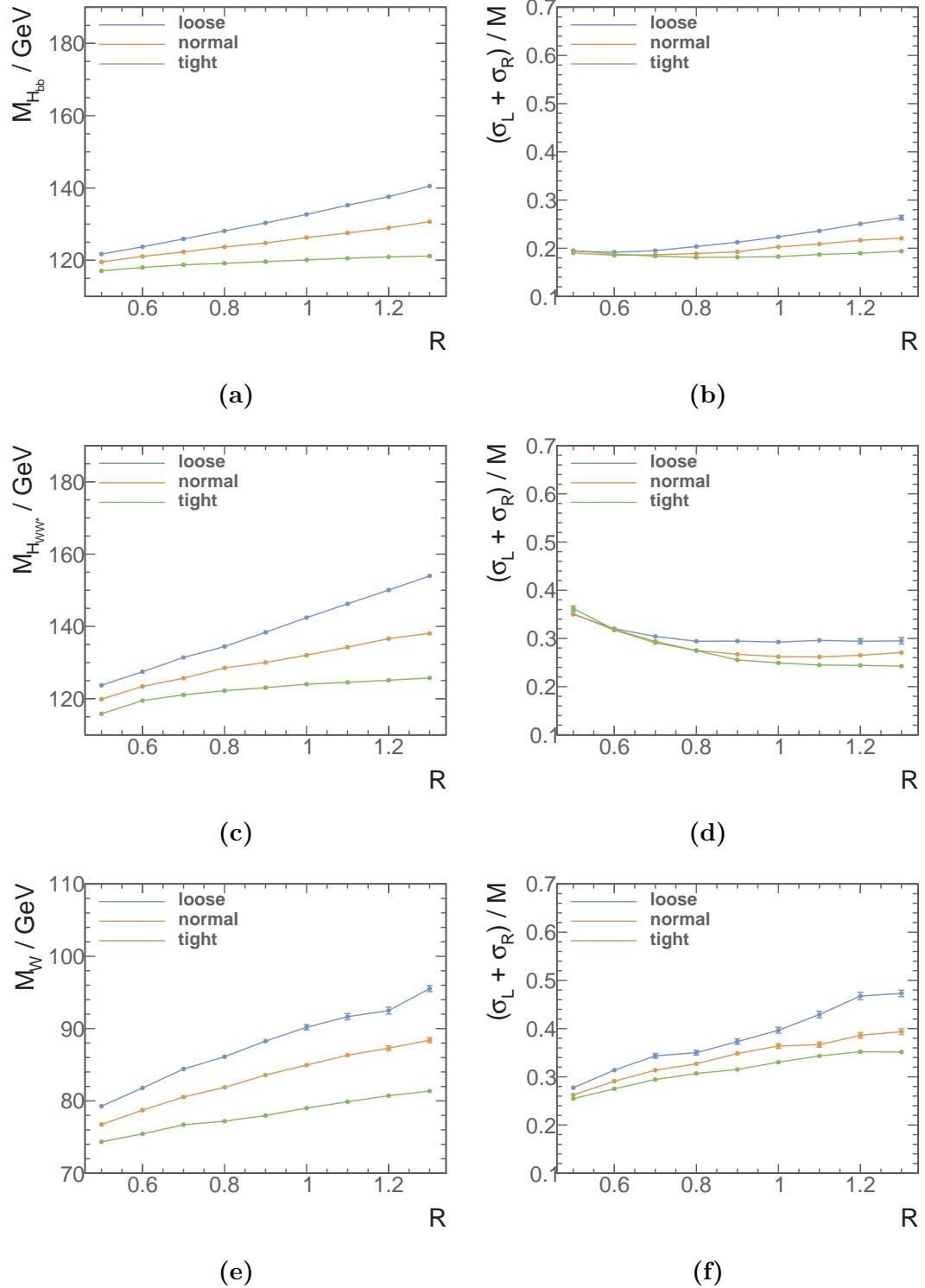


**Figure 7.3:** A typical example of  $m_{H_{bb}}$  mass distribution with the superimposed fitting function in red.

This is because more particles are included in jets with increasing  $R$ . Hence a larger invariant mass is obtained. For the relative width, values for  $H_{bb}$  increase with increasing  $R$ , but values for  $H_{WW^*}$  decrease. This is due to the compensating effect, as  $H_{WW^*}$  decays to 4 jets, preferring large jet radius, whilst  $H_{bb}$  decays to 2 jets, prefers small jet radius. For the similar reason,  $W$  prefers small jet radius and the values increase with increasing  $R$ .

The choice of PFO collection impacts number of PFOs in the event. The loose PFO selection has the most PFOs in the event, hence the largest invariant mass and worst mass resolution. This trend is consistent comparing loose to normal to tight PFO collections.

The optimal choice, normal selected PFO collection with  $R = 0.7$  give a good fitted mass for  $H_{WW^*}$  and  $W$ . The mass is slightly too low for the  $H_{bb}$ . The small  $R$  is good for  $m_{H_{bb}}$  and  $m_W$  resolution.  $m_{H_{WW^*}}$  resolution is relatively flat when  $R > 0.7$ . Hence normal selected PFO collection with  $R = 0.7$  is the optimal choice. The extracted fitted parameters of optimal jet reconstructions are summarised in table 7.8.



**Figure 7.4:** figure 7.4a, 7.4c, and 7.4e show fitted mass of  $H_{bb}$ ,  $H_{WW^*}$ , and  $W$  as a function of  $R$  parameter, for loose, normal and tight selected PFO collection, at  $\sqrt{s} = 1.4$  TeV. figure 7.4b, 7.4d, and 7.4f show realtive width of  $H_{bb}$ ,  $H_{WW^*}$ , and  $W$  as a function of  $R$  parameter, for loose, normal and tight selected PFO collection, at  $\sqrt{s} = 1.4$  TeV.

Fitted jet parameters $\sqrt{s} = 1.4 \text{ TeV}$	
$\mu_{H_{bb}}$	$122.3 \pm 0.2$
$\sigma_{L,H_{bb}}$	$15.2 \pm 0.2$
$\sigma_{R,H_{bb}}$	$7.55 \pm 0.16$
$\mu_{H_{WW^*}}$	$125.7 \pm 0.2$
$\sigma_{L,H_{WW^*}}$	$29.4 \pm 0.3$
$\sigma_{R,H_{WW^*}}$	$7.18 \pm 0.17$
$\mu_W$	$80.5 \pm 0.2$
$\sigma_{L,W}$	$16.2 \pm 0.3$
$\sigma_{R,W}$	$9.03 \pm 0.16$

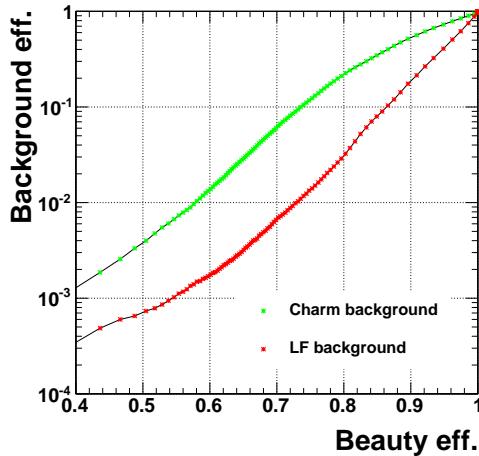
**Table 7.8:** The extracted fitted parameters of optimal jet reconstructions, normal selected PFO collection with  $R = 0.7$ , at  $\sqrt{s} = 1.4 \text{ TeV}$ .

## 7.5 Jet flavour tagging

As the signal channel,  $HH \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}q\bar{q}q\bar{q}$  has two b quarks in the final state, information can be extracted to determine the likelihood of a jet originated from a b quark. To find the likelihood of a b jet, b tag value, a software package LCFIPlus [60] is used. An overview of the flavour tagging processor is in section ??.

The flavour tagging is performed after the initial jet reconstruction, and all the PFOs in the reconstructed jets are the input to the LCFIPlus flavour tagging processor. LCFIPlus includes a multiclass classifier which needs to be trained. The samples for training the multiclass classifier are  $e^+e^- \rightarrow Z\bar{v}\nu$  at  $\sqrt{s} = 1.4 \text{ TeV}$ , where Z decays to  $b\bar{b}$ ,  $c\bar{c}$ , or  $u\bar{u}/d\bar{d}/s\bar{s}$ . The classifier in the LCFIPlus processor is trained with the optimal jet parameters, with the jet clustering step in the set to find two jets. The output of the processor for a jet is three values, corresponding to the likelihood of the jet being a b jet, a c jet, or a light flavour quark jet. The selection efficiency of b-jets and c-jets with training samples is shown in figure 7.5.

To use the LCFIPlus, all the PFOs in the initial reconstructed jet are fed into the processor. The jet clustering step in the LCFIPlus is set to find six jets. For each jet, values for the likelihood of a b jet and a c jet are obtained.



**Figure 7.5:** Performance of b-jet tagging with training samples at  $\sqrt{s} = 1.4$  TeV.

## 7.6 Jet pairing

Having optimised the jet reconstruction, and obtained the six jets from the jet clustering step in the LCFIPlus processor, the next step is to group the jets according to event topology. The hadronic decay of the  $HH \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}qqqq$  has six quarks final states, which results in six jets in the event. As the jet pairing scheme in the jet reconstruction optimisation section 7.4.1, jets are paired to with two jets for  $H \rightarrow b\bar{b}$ , two jets for hadronic decay of a  $W^*$ , two jets for of a  $W^*$ , and the two  $W$  forming a  $H$  from  $H \rightarrow b\bar{b}$ .  $W^*$  indicates the off-mass-shell  $W$  boson, because when a Higgs decays into two  $W$  bosons, one  $W$  is off the mass shell, as the Higgs mass is less than the sum two  $W$  masses (see section 2.8).

The jet pairing should reconstruct the invariant mass  $H_{bb}$ ,  $H_{WW^*}$ , and  $W$ , within the mass resolution. The theoretical possible mass peak and the mass resolution are obtained with the MC truth information, listed in the table 7.8. These numbers are used as for the jet pairing metric:

$$\chi^2 = \left( \frac{m_{ij} - \mu_{H_{bb}}}{\sigma'_{H_{bb}}} \right)^2 + \left( \frac{m_{klmn} - \mu_{H_{WW^*}}}{\sigma'_{H_{WW^*}}} \right)^2 + \left( \frac{m_{kl} - \mu_W}{\sigma'_W} \right)^2, \quad (7.3)$$

$$\sigma'_{H_{bb}} = \begin{cases} \sigma_{L,H_{bb}}, & \text{if } m_{ij} < \mu_{H_{bb}} \\ \sigma_{R,H_{bb}}, & \text{otherwise} \end{cases} \quad (7.4)$$

$$\sigma'_{H_{WW^*}} = \begin{cases} \sigma_{L,H_{WW^*}}, & \text{if } m_{klmn} < \mu_{H_{WW^*}} \\ \sigma_{R,H_{WW^*}}, & \text{otherwise} \end{cases} \quad (7.5)$$

$$\sigma'_W = \begin{cases} \sigma_{L,W}, & \text{if } m_{kl} < \mu_W \\ \sigma_{R,W}, & \text{otherwise} \end{cases} \quad (7.6)$$

where  $i$  to  $l$  indicate the one of the six jets, with all the possible combination.  $\mu$  and  $\sigma$  are the fitted invariant mass, and the fitted width from table 7.8. The asymmetrical structure of the fitting function is reflected in the jet pairing metric. The jet pairing with minimal  $\chi^2$  is chosen, with an additional requirement of at least one of two jets forming  $H_{bb}$  having a b-jet tag greater than 0.2.

## 7.7 Pre-selection

With reconstructed jets paired to the physical bosons, kinematic and topological variables can be calculated for the signal selection. A set of pre-selection cuts are places to aid the multivariate analysis. The pre-selection cuts falls into three categories: discriminative pre-selection cuts, loose cuts for the MVA, and the mutually exclusive cuts with the  $HH \rightarrow b\bar{b}q\bar{q}q\bar{q}$  sub-channel.

### 7.7.1 Discriminative pre-selection cuts

This set of pre-selection cuts are designed to discard the phase space which is dominated by the background events. The double Higgs system should have substantial invariant mass as both  $H$  are real. The two b jets in the signal final state is a signature. The final state contains neutrino. Therefore there is missing momentum in the event. These form

the logics for the pre-selection cuts. The cuts are listed in table 7.9 and shown in figure 7.6.

Figure 7.6a shows the distribution of the invariant mass of the two Higgs system, where the cut above 150 GeV is effective against samples with two quark final states. Figure 7.6c shows the distribution of the second highest b-jet tag, where the cut above 0.2 helps to reduce background events with no b-jets in final states. Figure 7.6b shows the distribution of the  $p_T$  of the two Higgs system, where the cut above 30 GeV is extremely effective against background channels with no neutrinos in the final state.

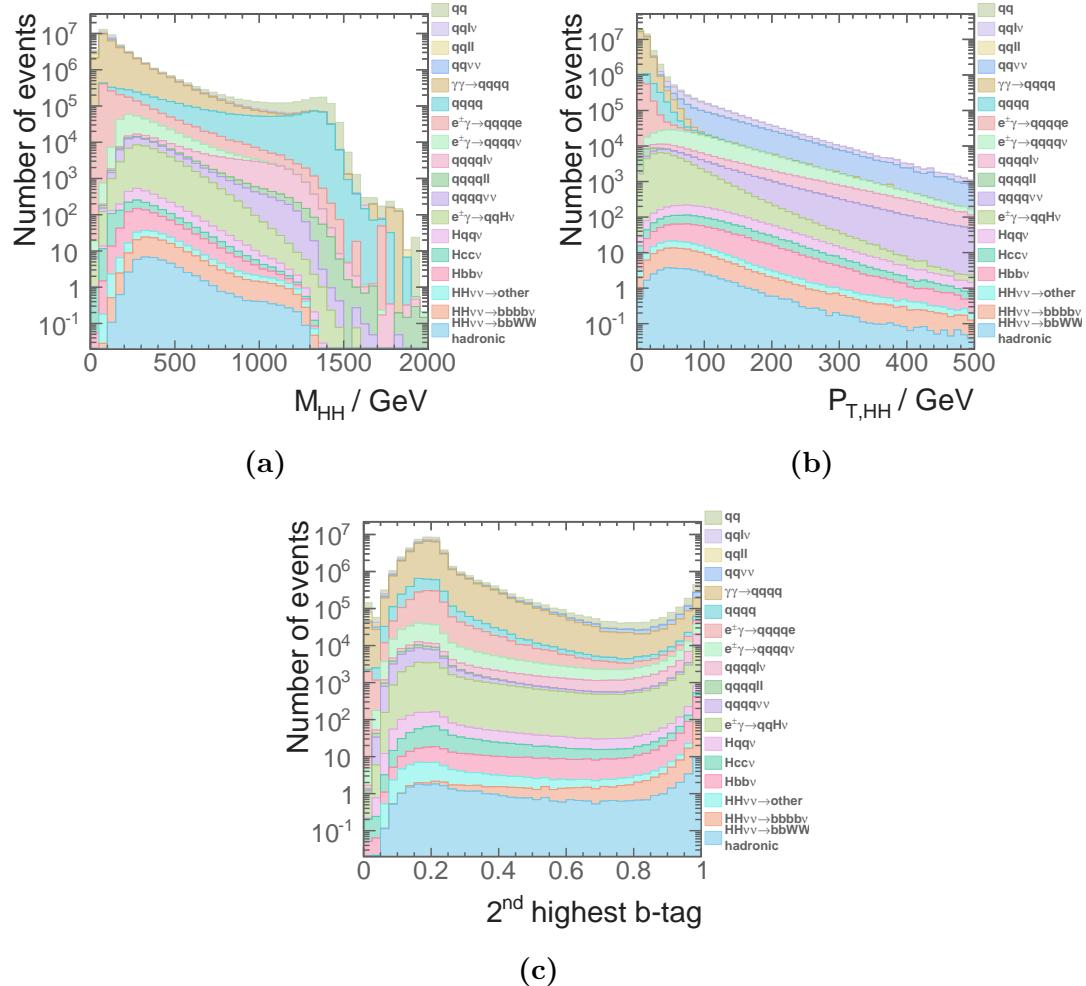
Pre-selection	$\sqrt{s} = 1.4 \text{ TeV}$
Discriminative pre-selection	$m_{HH} > 150 \text{ GeV}, B_2 > 0.2, p_{T_{HH}} > 30 \text{ GeV}$
Loose cuts for MVA	$m_{H_{bb}} < 500 \text{ GeV}, m_{H_{WW^*}} < 800 \text{ GeV},$ $m_W < 200 \text{ GeV}, m_{HH} < 1400 \text{ GeV}$
Mutually exclusive	$\Sigma B_{4\text{jets}} < 2.3, y_{34} < 3.7$

**Table 7.9:** Pre-selection cuts at  $\sqrt{s} = 1.4 \text{ TeV}$ .

The selection efficiency of the lepton veto and the pre-selection is shown in table 7.10. These pre-selection are very aggressive. The reason is that the cross sections of signal channel for is extremely small, comparing to the background. Hence only the signal events with very clear characteristic topologies would be able to pass the final MVA selection, optimised for the signal significance. Therefore, aggressive pre-selection cuts would not hurt the final signal selection. On the contrary, final signal selection efficiency would benefit, as the MVA is focused on the difficult background events, where their topologies are similar to the signal events.

### 7.7.2 Loose cuts for the MVA

A set of physics motivated loose cuts aims to reduce the range of invariant masses to increase the effectiveness of MVA. (See section ?? on MVA). The invariant masses of physical bosons are required to be within a certain range, to avoid the effect of extreme values on the MVA. The cuts are listed in table 7.9 and the performance is listed in figure 7.12.



**Figure 7.6:** Discriminative pre-selection variables for  $\sqrt{s} = 1.4 \text{ TeV}$ , after rejecting events with leptons, and jet pairing.

Channel / Efficiency $\sqrt{s} = 1.4 \text{ TeV}$	Expected number of events	Lepton ID and jet pairing	$m_{\text{HH}} > 150 \text{ GeV}$	$B_2 > 0.2$	$p_T > 30 \text{ GeV}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , hadronic	27.9	85.8%	85.6%	73.7%	66.4%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	67.6	90.8%	90.5%	90.1%	80.6%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow \text{other}$	128.0	36.2%	35.3%	27.7%	24.7%
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	1304.0	60.7%	59.8%	44.9%	42.0%
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	546.1	67.4%	57.7%	46.5%	43.4%
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	463.0	73.9%	72.6%	68.7%	64.2%
$e^-e^+ \rightarrow qqqq$	1867650.0	48.8%	46.1%	17.3%	4.7%
$e^-e^+ \rightarrow qqqq\ell\ell$	93150.0	5.0%	4.9%	1.5%	0.3%
$e^-e^+ \rightarrow qqqq\ell\nu$	165600.0	15.1%	15.1%	12.4%	11.4%
$e^-e^+ \rightarrow qqqq\nu\bar{\nu}$	34800.0	50.7%	50.0%	20.1%	18.8%
$e^-e^+ \rightarrow qq$	6014250.0	54.5%	17.5%	8.4%	2.2%
$e^-e^+ \rightarrow qq\ell\nu$	6464550.0	14.1%	5.3%	2.0%	1.6%
$e^-e^+ \rightarrow qq\ell\ell$	4088700.0	13.0%	1.1%	0.6%	0.1%
$e^-e^+ \rightarrow qq\nu\nu$	1181550.0	60.1%	12.3%	6.2%	5.8%
$e^-\gamma(\text{BS}) \rightarrow e^-qqqq$	1305787.5	23.3%	10.6%	4.4%	0.4%
$e^+\gamma(\text{BS}) \rightarrow e^+qqqq$	1300837.5	23.4%	10.5%	4.3%	0.4%
$e^-\gamma(\text{EPA}) \rightarrow e^-qqqq$	430650.0	11.1%	5.4%	2.2%	0.3%
$e^+\gamma(\text{EPA}) \rightarrow e^+qqqq$	430350.0	11.1%	5.3%	2.1%	0.3%
$e^-\gamma(\text{BS}) \rightarrow \nu qqqq$	89775.0	58.3%	56.8%	31.0%	27.7%
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu} qqqq$	89212.5	57.6%	56.1%	30.3%	27.3%
$e^-\gamma(\text{EPA}) \rightarrow \nu qqqq$	26100.0	29.6%	28.9%	15.4%	13.9%
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu} qqqq$	25950.0	29.2%	28.5%	15.0%	13.7%
$e^-\gamma(\text{BS}) \rightarrow qqH\nu$	17775	61.0%	59.8%	45.5%	34.6%
$e^+\gamma(\text{BS}) \rightarrow qqH\nu$	17662.5	61.1%	60.0%	45.6%	34.6%
$e^-\gamma(\text{EPA}) \rightarrow qqH\nu$	5085	31.8%	31.2%	23.7%	18.2%
$e^+\gamma(\text{EPA}) \rightarrow qqH\nu$	5085	31.9%	31.3%	23.8%	18.4%
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow qqqq$	2054951.5	56.3%	23.9%	9.6%	0.3%
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow qqqq$	4521037.5	33.6%	14.2%	5.7%	0.4%
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow qqqq$	4539150.0	33.7%	14.2%	5.7%	0.4%
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow qqqq$	1129500.0	21.1%	9.1%	3.7%	0.4%

**Table 7.10:** List of signal and background samples with the corresponding expected number at  $\sqrt{s} = 1.4 \text{ TeV}$ , assuming a luminosity of  $1500 \text{ fb}^{-1}$ . The selection efficiencies are presented in a “flow” fashion, as the every selection cut contains all the cuts to the left of it.

Selection efficiency	$\text{HH} \rightarrow b\bar{b}W^+W^-$ , hadronic	$\text{HH} \rightarrow b\bar{b}b\bar{b}$
$\Sigma B_{4\text{jets}} < 2.3$ and $y_{34} < 3.7$	86%	78%

**Table 7.11:** Mutually exclusive cuts at  $\sqrt{s} = 1.4 \text{ TeV}$ , shown for hadronic decay of  $\text{HH} \rightarrow b\bar{b}W^+W^-$  and  $\text{HH} \rightarrow b\bar{b}b\bar{b}$  subchannels.

### 7.7.3 Mutually exclusive cuts for $\text{HH} \rightarrow b\bar{b}W^+W^-$ and $\text{HH} \rightarrow b\bar{b}b\bar{b}$

This set of cuts is designed to divided samples, both signal and background, into two mutually exclusive sets, for the parallel analyses of two subchannels,  $\text{HH} \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}qqqq$  and  $\text{HH} \rightarrow b\bar{b}b\bar{b}$ . This eases the combining subchaneels, as correlations between subchannels do not need to be considered.

The most distinctive difference between two subchannels, is that they have different number of jets, and different number of b-jets in the final state. So variables related to number of b-jets and a number of jets are suitable for separating two subchannels.

Shown in figure 7.7, two subchannels can be clearly separated in the two dimensional parameter space. The optimal rectangular cuts were selected by scanning the two parameters, and maximising

$$\varepsilon = P(\text{subchannel}_1 | \text{selection}) \times P(\text{subchannel}_2 | \neg \text{selection}) \quad (7.7)$$

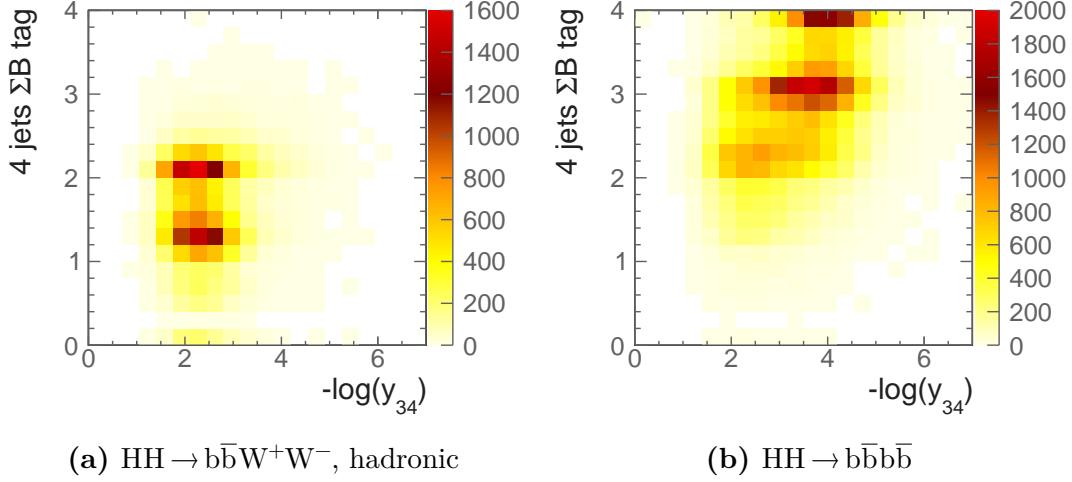
where **selection** represents the mutually exclusive cuts,  $\neg \text{selection}$  indicates the phase space not covered by the **selection**.

Variables tested includes  $\Sigma B_{4\text{jets}}$ ,  $\sum_1^3 B_{4\text{jets}}$ ,  $y_{34}$ ,  $y_{45}$ ,  $y_{56}$ ,  $y_{67}$  and other related variables. The best separation was summarised in table 7.11. The  $\Sigma B_{4\text{jets}}$  is the sum of the b tag values, when clustering an event to four jets.  $y$  parameters measures the number of jets in a event (see section ??).

The selection efficiencies after mutually exclusive cuts are listed in table 7.12. Mutually exclusive cuts veto most  $\text{HH} \rightarrow b\bar{b}b\bar{b}$  events, as designed.

Channel	Previous cuts and loose cuts	Mutually exclusive
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , hadronic	66.4%	59.7%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	80.6%	15.4%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow$ other	24.7%	20.5%
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	42.0%	39.5%
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	43.4%	31.7%
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	64.2%	25.2%
$e^-e^+ \rightarrow qqqq$	4.6%	3.4%
$e^-e^+ \rightarrow qq\bar{q}q\ell\ell$	3.3%	3.1%
$e^-e^+ \rightarrow qq\bar{q}q\ell\nu$	11.4%	9.8%
$e^-e^+ \rightarrow qq\bar{q}q\nu\bar{\nu}$	18.8%	16.6%
$e^-e^+ \rightarrow qq\bar{q}q$	2.0%	0.8%
$e^-e^+ \rightarrow qq\ell\nu$	1.6%	0.9%
$e^-e^+ \rightarrow qq\ell\ell$	0.1%	0.1%
$e^-e^+ \rightarrow qq\nu\bar{\nu}$	5.8%	4.0%
$e^-\gamma(BS) \rightarrow e^-qqqq$	0.4%	0.3%
$e^+\gamma(BS) \rightarrow e^+qqqq$	0.4%	0.4%
$e^-\gamma(EPA) \rightarrow e^-qqqq$	0.3%	0.2%
$e^+\gamma(EPA) \rightarrow e^+qqqq$	0.3%	0.3%
$e^-\gamma(BS) \rightarrow \nu qq\bar{q}q$	27.7%	25.3%
$e^+\gamma(BS) \rightarrow \bar{\nu} qq\bar{q}q$	27.3%	24.9%
$e^-\gamma(EPA) \rightarrow \nu qq\bar{q}q$	13.9%	12.6%
$e^+\gamma(EPA) \rightarrow \bar{\nu} qq\bar{q}q$	13.7%	12.3%
$e^-\gamma(BS) \rightarrow qqH\nu$	34.6%	30.6%
$e^+\gamma(BS) \rightarrow qqH\nu$	34.6%	30.6%
$e^-\gamma(EPA) \rightarrow qqH\nu$	18.2%	16.0%
$e^+\gamma(EPA) \rightarrow qqH\nu$	18.4%	16.1%
$\gamma(BS)\gamma(BS) \rightarrow qq\bar{q}q$	0.3%	0.3%
$\gamma(BS)\gamma(EPA) \rightarrow qq\bar{q}q$	0.4%	0.3%
$\gamma(EPA)\gamma(BS) \rightarrow qq\bar{q}q$	0.4%	0.3%
$\gamma(EPA)\gamma(EPA) \rightarrow qq\bar{q}q$	0.4%	0.3%

**Table 7.12:** List of signal and background samples with the selection efficiencies of loose cuts and mutually exclusive cuts at  $\sqrt{s} = 1.4$  TeV. The selection efficiencies are presented in a “flow” fashion, as the every selection cut contains all the cuts to the left of it.



**Figure 7.7:** Sum of b tag against  $y_{34}$  at  $\sqrt{s} = 1.4$  TeV, shown for hadronic decay of  $\text{HH} \rightarrow b\bar{b}W^+W^-$  and  $\text{HH} \rightarrow b\bar{b}b\bar{b}$  subchannel.

## 7.8 Discriminative variables for MVA

A series of discriminative variables were calculated to differentiate the signal and background events. These variables are fed into MVA for signal selection. The full list of variables can be found in table 7.13. The variables are grouped into categories.

Invariant mass variables are very effective at selecting signal events as no background events have double Higgs bosons in final states. Figure 7.8a and figure 7.8b show the distributions of  $m_{H_{bb}}$  and  $m_{H_{WW^*}}$  after all pre-selection cuts.

For the off-shell  $W^*$ , its energy is used as its mass distribution doesn't have a resonance. For the recoil momenta, which is calculated by assuming the collision at  $\sqrt{s}$  and a beam crossing angle 20 mrad, the pseudorapidity is used to focus on the forward region. The pseudorapidity,  $\eta$ , is defined as:

$$\eta \equiv -\ln \left[ \tan \left( \frac{\theta}{2} \right) \right], \quad (7.8)$$

where  $\theta$  is the polar angle.  $A_{12}$  measures the angle between the two constituent jets, defined as:

$$A_{12} = \pi - \cos^{-1} (\hat{\mathbf{p}}_1 \cdot \hat{\mathbf{p}}_2), \quad (7.9)$$

Category	Variable
Invariant mass	$m_{H_{bb}}, m_{H_{WW^*}}, m_W, m_{HH}$
Energy and momentum	$E_{W^*}, E_{mis}, p_{TH_{bb}}, p_{TH_{WW^*}}, p_{THH}, p_{T_{HH}}$
Angles in lab frame	$\eta_{mis}, A_{H_{bb}}, A_W, A_{HH}$
Angles in boosted frames	$\cos(\theta_{H_{bb}}^*), \cos(\theta_{H_{WW^*}}^*), \cos(\theta_W^*), \cos(\theta_{W^*}^*), \cos(\theta_{HH}^*)$
Event shape	$ S , -\ln(y_{23}), -\ln(y_{34}), -\ln(y_{45}), -\ln(y_{56})$
b and c tag	$B_{1,H_{bb}}, B_{2,H_{bb}}, B_{1,W}, B_{1,W^*}, C_{1,H_{bb}}, C_{1,W}$
Number of PFOs	$N_{H_{bb}}, N_{H_{WW^*}}, N_W, N_{W^*}$

**Table 7.13:** Variables used in MVA at  $\sqrt{s} = 1.4$  TeV

where  $\hat{p}_1$  is the normal momentum vector of jet 1.  $\cos(\theta_{12}^*)$  is the cosine of the angle between the two constituent jets in their decay rest frame. Figure 7.8c and figure ?? compares the  $A_{H_{bb}}$  and  $\cos(\theta_{H_{bb}}^*)$ . Both shows different distribution for the signal and background channels. For the signal,  $\cos(\theta_{H_{bb}}^*)$  has a flat distribution, as expected from a back-to-back decay of  $H \rightarrow b\bar{b}$ . For the background,  $\cos(\theta_{H_{bb}}^*)$  peaked at 1.

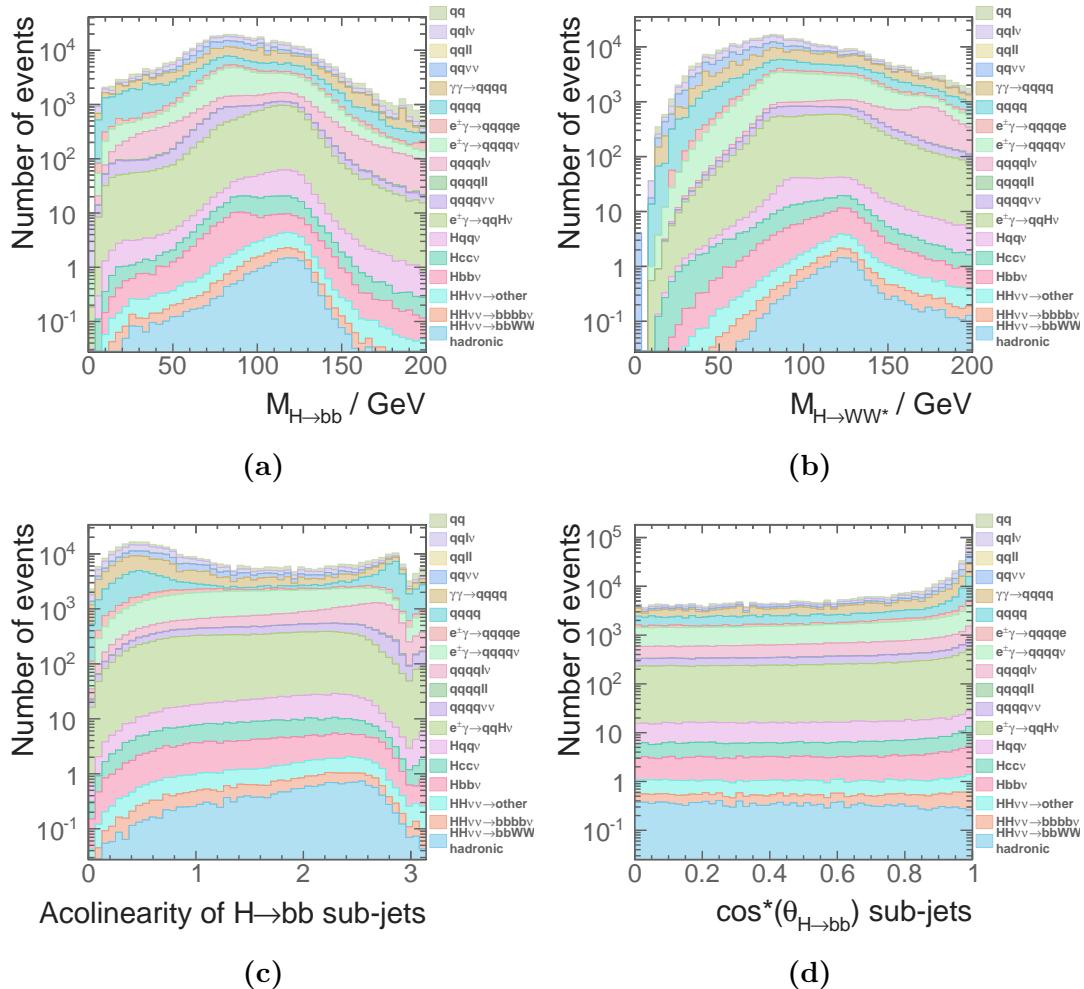
The global event shape variables includes  $y$  variables (see section 7.7.3), and the sphericity,  $S$ .  $S$  is a measurement of the spherically symmetry of the event (see section ??).

The flavour tagging variables are discussed in section 7.5.  $B_{1,H_{bb}}$  denotes the Highest b-jet tag value of two jets forming  $H_{bb}$ . The number of PFOs variables are effective against background events with fewer quarks in final states.

The optimal set of 32 variables were chosen to give the best MVA performance, whilst no strong pair-wise correlation between any two variables.

## 7.9 Multivariate analysis

After gathering the information and applying the pre-selection cuts, the signal selection is performed using multivariate analysis (MVA) with Boosted Decision Tree classifier (BDT). The parameters for boosted decision tree were optimised and checked for overtraining.



**Figure 7.8:** Stacked plots for discriminative variables for the MVA at  $\sqrt{s} = 1.4$  TeV after all pre-selection cuts.

Parameter	Value
Depth of tree	4
Number of trees	4000
The minimum number of events in a node	0.25% of the total events
Boosting	adaptive boost
learning rate of the adaptive boost	0.5
metric for the optimal cuts	Gini Index
Number of bins per variables	40
End node output	$x \in [0, 1]$

**Table 7.14:** Optimised parameters for the boosted decision tree.

A brief discussion on the MVA about the classifier, overtraining and the result can be found in see section ??.

The optimisation of the BDT follows the strategy in section ???. The optimised parameters are listed in table ???. The optimal values are obtained by choosing the best performance without overfitting with samples at  $\sqrt{s} = 3$  TeV. The same values are used  $\sqrt{s} = 1.4$  analysis TeV.

Half of the samples were used for training, and the other half used for testing and classifier optimisation.

The signal for the MVA is the hadronic decay of  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ .  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  decays to other final state does not participate in the MVA training step, as they are from the Feynman diagrams as the signal (see section 7.1), event topologies are different.  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  decaying to other final state would participate in the MVA applying stage.

## 7.10 Signal selection results

After applying the MVA, the efficiencies of the pre-selection cuts, the efficiencies of the MVA selections are listed in table 7.14, alongside with the number of events after the MVA selection. A few background channels survive after the MVA.  $e^-e^+ \rightarrow q\bar{q}H\nu\bar{\nu}$  survives as the one single Higgs plus neutrino has a very similar topology to the signal. Similar  $e^-e^+ \rightarrow qqql\nu$  can be confused as the signal when the lepton is undetected in

the forward region, or the energy is too low to be tagged.  $e^-e^+ \rightarrow q\bar{q}q\bar{q}\nu\bar{\nu}$  has a similar topologies. The electron photon and photon interactions with same final states as the above channels also survive the MVA.

Before interpreting the result, the analysis at  $\sqrt{s} = 3$  TeV and the analysis with the semi-leptonic channel of  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$  are presented.

## 7.11 $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$ hadronic decay at $\sqrt{s} = 3$ TeV analysis

The  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$  hadronic decay at  $\sqrt{s} = 3$  TeV analysis follows the same strategy as the analysis at  $\sqrt{s} = 1.4$  TeV. Brief discussion of each step and results are provided, where the differences are highlighted. Cross sections of used samples are listed in table 7.15.

The lepton finding processors are either developed or optimised with samples at  $\sqrt{s} = 1.4$  TeV, and checked against samples at  $\sqrt{s} = 3$  TeV (see section 7.3). It was found that the same set of parameters for lepton identifiers work well under  $\sqrt{s} = 1.4$  TeV and 3 TeV. The performance of the lepton processors is shown in table 7.16.

Comparing  $\sqrt{s} = 1.4$  TeV and  $\sqrt{s} = 3$  TeV, the lepton finding performance is worse for  $\sqrt{s} = 3$  TeV. This is because particles are more boosted and separation between particles is smaller. This reflects on the performance of the ForwardFinderProcessor. Whilst at  $\sqrt{s} = 1.4$  TeV, the processor only rejects 5% background and 1% signal, at  $\sqrt{s} = 3$  TeV it rejects 19% background and 4% signal, suggesting many leptons are in the forward region. These processors are complimentary and a good rejection rate with combined processors is achieved.

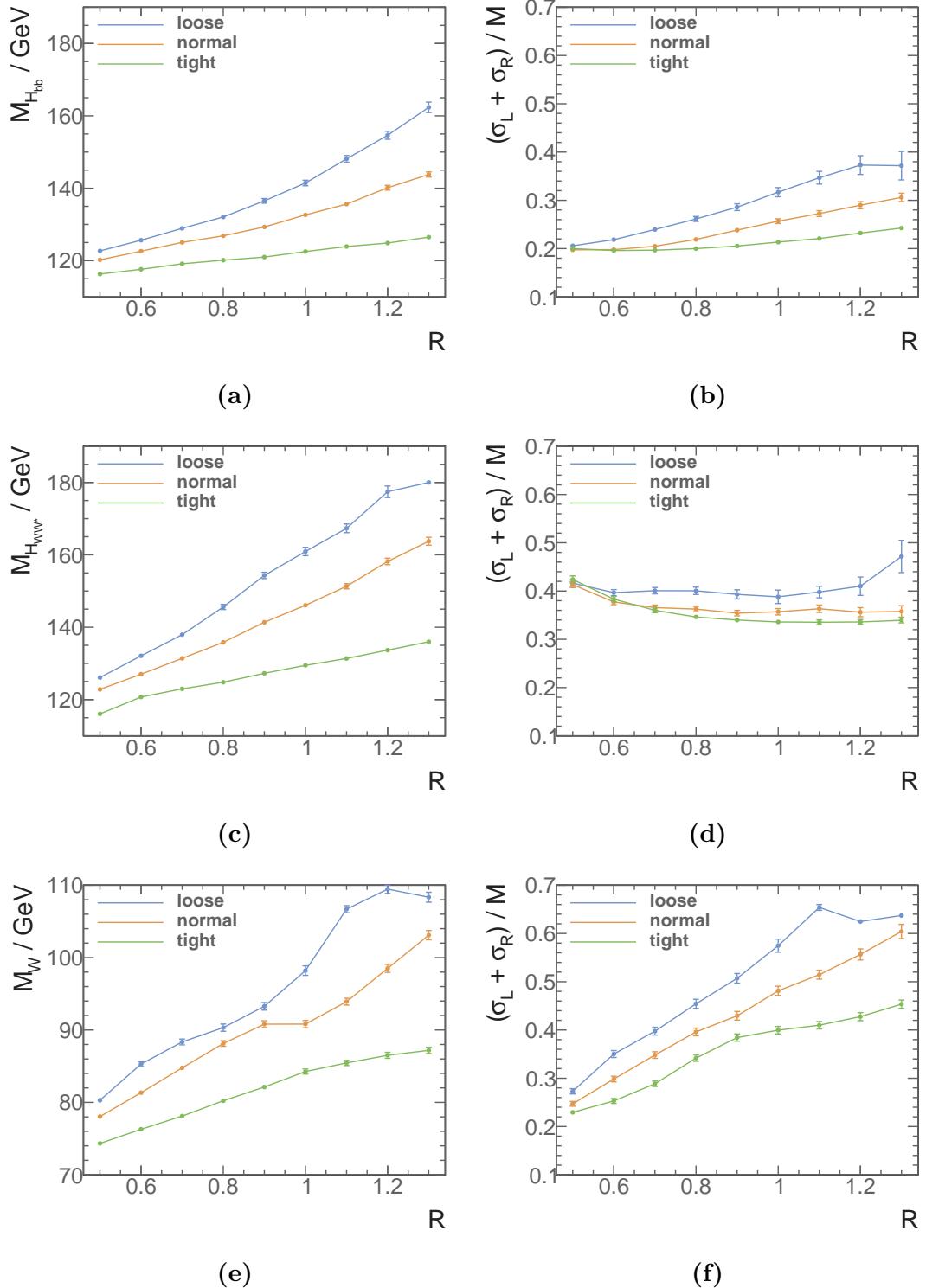
For the jet reconstruction optimisation, same strategy listed in section 7.4.1 is followed. Figure 7.9 shows fitted mass for  $H_{bb}$ ,  $H_{WW^*}$ , and  $W$ , along with the relative mass resolution. The relative resolution of  $W$  rises sharply with increasing  $R$ , hence favouring a small  $R$  and tight selected PFO collection. The optimal chosen is tight selected PFO collection with  $R = 0.7$ . The invariant mass is underestimated for the better mass resolution. The fitted values of the chosen jet reconstruction are listed in table 7.17

Channel / Efficiency $\sqrt{s} = 1.4 \text{ TeV}$	N	$\epsilon_{\text{presel}}$	$\epsilon_{\text{MVA}}$	$N_{\text{MVA}}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e, \text{ hadronic}$	27.9	59.8%	8.2%	1.29
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	67.6	15.4%	0.5%	0.05
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow \text{other}$	128.0	20.4%	1.7%	0.45
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	1304.0	39.5%	0.05%	0.29
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	546.1	31.6%	0.1%	0.16
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	463.0	24.7%	0.3%	0.37
$e^-e^+ \rightarrow qqqq$	1867650.0	3.3%	-	-
$e^-e^+ \rightarrow qqqq\ell\ell$	93150.0	0.3%	-	-
$e^-e^+ \rightarrow qqqq\ell\nu$	165600.0	9.8%	0.01%	2.06
$e^-e^+ \rightarrow qqqq\nu\bar{\nu}$	34800.0	16.5%	0.002%	0.10
$e^-e^+ \rightarrow qq$	6014250.0	0.8%	-	-
$e^-e^+ \rightarrow qq\ell\nu$	6464550.0	0.9%	-	-
$e^-e^+ \rightarrow qq\ell\ell$	4088700.0	0.08%	-	-
$e^-e^+ \rightarrow qq\nu\bar{\nu}$	1181550.0	4.0%	-	-
$e^-\gamma(\text{BS}) \rightarrow e^-qqqq$	1305787.5	0.3%	-	-
$e^+\gamma(\text{BS}) \rightarrow e^+qqqq$	1300837.5	0.4%	-	-
$e^-\gamma(\text{EPA}) \rightarrow e^-qqqq$	430650.0	0.3%	-	-
$e^+\gamma(\text{EPA}) \rightarrow e^+qqqq$	430350.0	0.3%	-	-
$e^-\gamma(\text{BS}) \rightarrow \nu qqqq$	89775.0	25.4%	0.005%	1.09
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu} qqqq$	89212.5	24.9%	0.004%	0.96
$e^-\gamma(\text{EPA}) \rightarrow \nu qqqq$	26100.0	12.6%	-	-
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu} qqqq$	25950.0	12.4%	0.008%	0.27
$e^-\gamma(\text{BS}) \rightarrow qqH\nu$	17775	30.8%	0.02%	1.00
$e^+\gamma(\text{BS}) \rightarrow qqH\nu$	17662.5	30.6%	0.02%	1.16
$e^-\gamma(\text{EPA}) \rightarrow qqH\nu$	5085	16.0%	0.04%	0.33
$e^+\gamma(\text{EPA}) \rightarrow qqH\nu$	5085	16.2%	0.08%	0.62
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow qqqq$	2054951.5	0.2%	-	-
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow qqqq$	4521037.5	0.4%	-	-
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow qqqq$	4539150.0	0.3%	-	-
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow qqqq$	1129500.0	0.3%	-	-

**Table 7.15:** List of signal and background samples with the corresponding expected number at  $\sqrt{s} = 1.4 \text{ TeV}$ , for a luminosity of  $1500 \text{ fb}^{-1}$ . The number of events, selection efficiency of pre-selection, selection efficiency of MVA after pre-selection, number of events after MVA are shown. - represents number less than 0.01.

Channel	$\sigma(\sqrt{s} = 3 \text{ TeV}) / \text{fb}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$	0.588
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-,\text{hadronic}$	0.07
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	0.19
$e^-e^+ \rightarrow HH \rightarrow \text{others}$	0.34
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	1.78
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	1.12
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	1.91
$e^-e^+ \rightarrow qqqq$	546.5*
$e^-e^+ \rightarrow qqql\ell\ell$	169.3*
$e^-e^+ \rightarrow qqql\ell\nu$	106.6*
$e^-e^+ \rightarrow qqql\nu\bar{\nu}$	71.5*
$e^-e^+ \rightarrow qq$	2948.9
$e^-e^+ \rightarrow qq\ell\nu$	5561.1
$e^-e^+ \rightarrow qq\ell\ell$	3319.6
$e^-e^+ \rightarrow qq\nu\nu$	1317.5
$e^-\gamma(\text{BS}) \rightarrow e^-qqqq$	1268.7*
$e^+\gamma(\text{BS}) \rightarrow e^+qqqq$	1267.6*
$e^-\gamma(\text{EPA}) \rightarrow e^-qqqq$	287.9*
$e^+\gamma(\text{EPA}) \rightarrow e^+qqqq$	287.8*
$e^-\gamma(\text{BS}) \rightarrow \nu qqqq$	262.5*
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu} qqqq$	262.3*
$e^-\gamma(\text{EPA}) \rightarrow \nu qqqq$	54.2*
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu} qqqq$	54.2*
$e^-\gamma(\text{BS}) \rightarrow qqH\nu$	58.6*
$e^+\gamma(\text{BS}) \rightarrow qqH\nu$	58.5*
$e^-\gamma(\text{EPA}) \rightarrow qqH\nu$	11.7*
$e^+\gamma(\text{EPA}) \rightarrow qqH\nu$	11.7*
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow qqqq$	13050.3*
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow qqqq$	2420.6*
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow qqqq$	2423.1*
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow qqqq$	402.7*

**Table 7.16:** List of signal and background samples with the corresponding cross sections at  $\sqrt{s} = 3 \text{ TeV}$ .  $q$  can be  $u, d, s, b$  or  $t$ . Unless specified,  $q, \ell$  and  $\nu$  represent particles and its corresponding anti-particles.  $\gamma$  (BS) represents a real photon from beamstrahlung (BS).  $\gamma$  (EPA) represents a “quasi-real” photon, simulated with the Equivalent Photon Approximation. For processes involving Higgs production explicitly, simulated Higgs mass is 126 GeV. Otherwise, Higgs mass is set to 14 TeV. For processes labelled with \*, the generator level cut requires invariant mass of quarks greater than 50.



**Figure 7.9:** figure 7.9a, 7.9c, and 7.9e show fitted mass of  $H_{bb}$ ,  $H_{WW^*}$ , and  $W$ , respectively, for loose, normal and tight selected PFO collection as a function of  $R$  parameter, at  $\sqrt{s} = 3 \text{ TeV}$ . figure 7.9b, 7.9d, and 7.9f show relative mass resolutions of  $H_{bb}$ ,  $H_{WW^*}$ , and  $W$ , respectively, for loose, normal and tight selected PFO collection as a function of  $R$  parameter, at  $\sqrt{s} = 3 \text{ TeV}$ .

Efficiency (3 TeV)	Signal	$e^+e^- \rightarrow q\bar{q}q\bar{q}\ell\nu$
IsolatedLeptonFinderProcessor	99.5%	66.8%
BonoLeptonFinderProcessor	99.0%	52.5%
TauFinderProcessor	97.7%	79.5%
BonoTauFinderProcessor	86.3%	60.3%
ForwardFinderProcessor	95.9%	80.7%
Combined	81.0%	23.3%

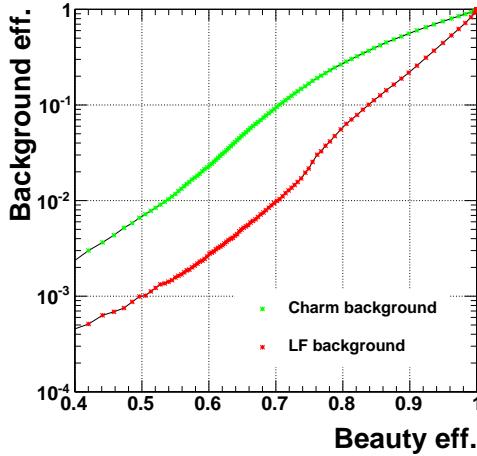
**Table 7.17:** Isolated lepton finder processors performance on the signal and selected background samples at  $\sqrt{s} = 3$  TeV.

Jet Parameters	$\sqrt{s} = 3$ TeV
$\mu_{H_{bb}}$	$119.1 \pm 0.3$
$\sigma_{L,H_{bb}}$	$15.0 \pm 0.3$
$\sigma_{R,H_{bb}}$	$8.4 \pm 0.2$
$\mu_{H_{WW^*}}$	$123.0 \pm 0.3$
$\sigma_{L,H_{WW^*}}$	$36.6 \pm 0.6$
$\sigma_{R,H_{WW^*}}$	$7.4 \pm 0.2$
$\mu_W$	$78.1 \pm 0.3$
$\sigma_{L,W}$	$13.1 \pm 0.4$
$\sigma_{R,W}$	$9.5 \pm 0.2$

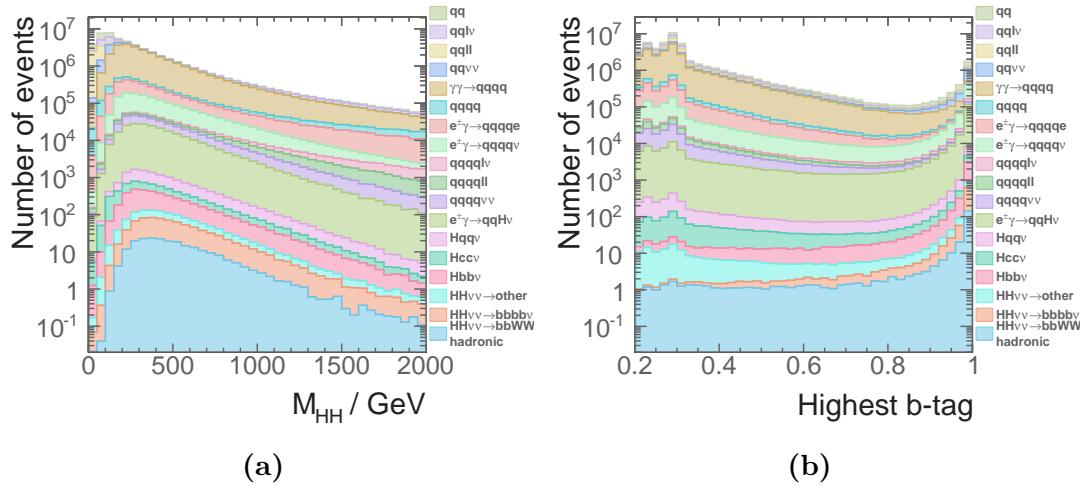
**Table 7.18:** The extracted fitted parameters of optimal jet reconstructions for tight selected PFO collection with  $R = 0.7$  at  $\sqrt{s} = 3$  TeV.

The flavour tagging processor is trained with the chosen jet parameters at  $\sqrt{s} = 3$  TeV. The performance of the flavour tagging with training samples is shown in figure 7.10. Comparing to the performance at  $\sqrt{s} = 1.4$  TeV, the flavour tagging performs worse. At high energy, particles are more collimated and more difficult to separate. Hence the performance degrades.

The pre-selection cuts at  $\sqrt{s} = 3$  TeV have changed, listed in table 7.18. Figure 7.11a shows the distribution of the invariant mass of the two Higgs system, where the cut above 150 GeV is effective against samples with two quark final states. Figure 7.11b shows the distribution of the highest b-jet tag, where the cut above 0.7 helps to reduce background events with no b-jets in final states. The cut is aggressive to compensate for the worse performance of the flavour tagging at high  $\sqrt{s}$ .



**Figure 7.10:** Performance of b-jet tagging with training samples at  $\sqrt{s} = 3$  TeV.

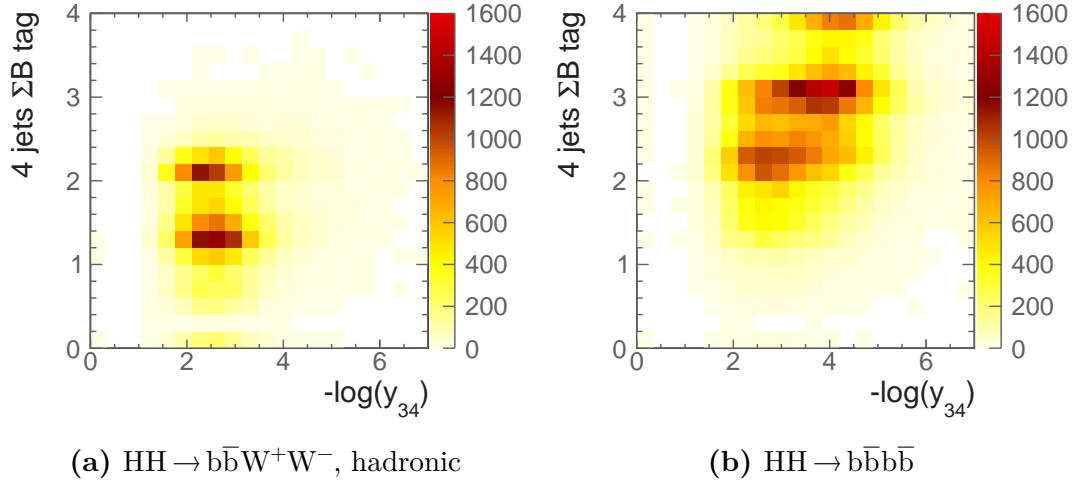


**Figure 7.11:** Discriminative pre-selection variables at  $\sqrt{s} = 3$  TeV, after rejecting events with identified leptons, and jet pairing.

Pre-selection	$\sqrt{s} = 3$ TeV
Discriminative pre-selection	$m_{HH} > 150$ GeV, $B_1 > 0.2$ , $p_{T_{HH}} > 30$ GeV
Loose cuts for MVA	$m_{H_{bb}} < 500$ GeV, $m_{H_{WW^*}} < 800$ GeV, $m_W < 200$ GeV, $m_{HH} < 3000$ GeV
Mutually exclusive	$\Sigma B_{4\text{jets}} < 2.3$ , $y_{34} < 3.6$

**Table 7.19:** Pre-selection cuts at  $\sqrt{s} = 3$  TeV.

The loose cuts for MVA at  $\sqrt{s} = 3$  TeV are largely the same as the ones at  $\sqrt{s} = 3$  TeV, apart from the cut on the invariant mass of HH due to higher  $\sqrt{s}$ . The selection efficiency



**Figure 7.12:** Sum of b tag against  $y_{34}$  at  $\sqrt{s} = 3$  TeV, shown for hadronic decay of  $\text{HH} \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}qqqq$  and  $\text{HH} \rightarrow b\bar{b}W^+W^-$  subchannel.

of the lepton veto and the pre-selection is shown in table 7.19. These pre-selection are still very aggressive. The cut on  $m_{\text{HH}}$  is effective against background with fewer number of quarks in the final states. The cut on  $B_1$  is effective against final states with no b quark.

The mutually exclusive cuts divide samples, both signal and background, into two mutually exclusive sets, for the parallel analyses of two subchannels,  $\text{HH} \rightarrow b\bar{b}W^+W^- \rightarrow b\bar{b}qqqq$  and  $\text{HH} \rightarrow b\bar{b}b\bar{b}$ . The cuts are obtained using the same strategy in section 7.7.3. The values are listed in table 7.18. The two dimensional spaces for two subchannels are shown in figure 7.12. The selection efficiencies after the loose cuts and the mutually exclusive cuts are shown in table 7.20.

The same set of variables are used at  $\sqrt{s} = 3$  TeV. The MVA BDT classifier is optimised at  $\sqrt{s} = 3$  TeV. The efficiencies of the pre-selection cuts, the efficiencies of the MVA selections are listed in table 7.21, alongside with the number of events after the MVA selection. Background channels survived the MVA are almost identical to those at  $\sqrt{s} = 1.4$  TeV. Hence see section 7.10 for the discussion.

Channel / Efficiency $\sqrt{s} = 3 \text{ TeV}$	Expected number of events	Lepton ID and jet pairing	$m_{\text{HH}} > 150 \text{ GeV}$	$B_1 > 0.7$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , hadronic	146.0	80.2%	79.9%	69.7%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	355.0	83.4%	82.9%	81.2%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow \text{other}$	675.0	36.7%	35.8%	25.2%
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	6115.4	59.5%	58.5%	40.4%
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	2249.9	64.8%	58.4%	39.3%
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	2197.7	69.7%	68.4%	64.2%
$e^-e^+ \rightarrow q\bar{q}q\bar{q}$	1093000.0	48.5%	39.7%	3.0%
$e^-e^+ \rightarrow q\bar{q}q\bar{q}\ell\bar{\ell}$	338600.0	14.7%	14.2%	0.7%
$e^-e^+ \rightarrow q\bar{q}q\bar{q}\ell\nu$	213200.0	19.7%	19.4%	10.0%
$e^-e^+ \rightarrow q\bar{q}q\bar{q}\nu\bar{\nu}$	143000.0	58.4%	57.3%	11.9%
$e^-e^+ \rightarrow q\bar{q}$	5897800.0	62.8%	13.2%	2.7%
$e^-e^+ \rightarrow q\bar{q}\ell\nu$	11121800	28.3%	11.9%	0.3%
$e^-e^+ \rightarrow q\bar{q}\ell\bar{\ell}$	6639200.0	38.3%	2.9%	0.7%
$e^-e^+ \rightarrow q\bar{q}\nu\bar{\nu}$	2635000.0	71.4%	24.1%	5.3%
$e^-\gamma(\text{BS}) \rightarrow e^-q\bar{q}q\bar{q}$	2004388.1	23.3%	21.5%	0.8%
$e^+\gamma(\text{BS}) \rightarrow e^+q\bar{q}q\bar{q}$	2002334.1	23.4%	21.6%	0.8%
$e^-\gamma(\text{EPA}) \rightarrow e^-q\bar{q}q\bar{q}$	575600.0	12.0%	11.0%	0.5%
$e^+\gamma(\text{EPA}) \rightarrow e^+q\bar{q}q\bar{q}$	575600.0	12.0%	10.9%	0.4%
$e^-\gamma(\text{BS}) \rightarrow \nu q\bar{q}q\bar{q}$	414750.0	61.7%	59.5%	20.4%
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu}q\bar{q}q\bar{q}$	414434.0	61.2%	59.1%	19.4%
$e^-\gamma(\text{EPA}) \rightarrow \nu q\bar{q}q\bar{q}$	108400.0	30.9%	29.9%	9.6%
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu}q\bar{q}q\bar{q}$	108400.0	30.7%	29.7%	9.1%
$e^-\gamma(\text{BS}) \rightarrow q\bar{q}H\nu$	92588.0	58.3%	56.2%	37.3%
$e^+\gamma(\text{BS}) \rightarrow q\bar{q}H\nu$	92430.0	58.1%	56.0%	37.1%
$e^-\gamma(\text{EPA}) \rightarrow q\bar{q}H\nu$	23400.0	30.1%	29.2%	19.4%
$e^+\gamma(\text{EPA}) \rightarrow q\bar{q}H\nu$	23400.0	29.7%	28.6%	18.8%
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow q\bar{q}q\bar{q}$	18009413.9	54.2%	49.2%	1.9%
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow q\bar{q}q\bar{q}$	3824548.1	33.5%	30.2%	1.2%
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow q\bar{q}q\bar{q}$	3828498.1	33.7%	30.3%	1.2%
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow q\bar{q}q\bar{q}$	805400.0	22.0%	19.8%	0.8%

**Table 7.20:** List of signal and background samples with the corresponding expected number at  $\sqrt{s} = 3 \text{ TeV}$ , assuming a luminosity of  $2000 \text{ fb}^{-1}$ . The selection efficiencies are presented in a “flow” fashion, as the every selection cut contains all the cuts to the left of it.

Channel	Previous cuts and loose cuts	Mutually exclusive
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , hadronic	69.5%	61.7%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	81.1%	18.8%
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow$ other	25.1%	20.0%
$e^-e^+ \rightarrow q_lq_lH\nu\bar{\nu}$	40.3%	35.9%
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	39.2%	26.2%
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	64.2%	25.9%
$e^-e^+ \rightarrow qqqq$	2.5%	1.4%
$e^-e^+ \rightarrow qqqq\ell\ell$	0.7%	0.6%
$e^-e^+ \rightarrow qqqq\ell\nu$	9.2%	7.2%
$e^-e^+ \rightarrow qqqq\nu\bar{\nu}$	11.8%	9.0%
$e^-e^+ \rightarrow qq$	2.5%	1.4%
$e^-e^+ \rightarrow qq\ell\nu$	0.3%	0.1%
$e^-e^+ \rightarrow qq\ell\ell$	0.7%	0.4%
$e^-e^+ \rightarrow qq\nu\nu$	5.3%	3.1%
$e^-\gamma(BS) \rightarrow e^-qqqq$	0.8%	0.7%
$e^+\gamma(BS) \rightarrow e^+qqqq$	0.8%	0.7%
$e^-\gamma(EPA) \rightarrow e^-qqqq$	0.4%	0.4%
$e^+\gamma(EPA) \rightarrow e^+qqqq$	0.4%	0.3%
$e^-\gamma(BS) \rightarrow \nu qqqq$	20.3%	16.8%
$e^+\gamma(BS) \rightarrow \bar{\nu} qqqq$	19.3%	15.9%
$e^-\gamma(EPA) \rightarrow \nu qqqq$	9.4%	7.8%
$e^+\gamma(EPA) \rightarrow \bar{\nu} qqqq$	8.9%	7.3%
$e^-\gamma(BS) \rightarrow qqH\nu$	37.2%	30.2%
$e^+\gamma(BS) \rightarrow qqH\nu$	37.1%	30.2%
$e^-\gamma(EPA) \rightarrow qqH\nu$	19.0%	15.7%
$e^+\gamma(EPA) \rightarrow qqH\nu$	18.4%	15.2%
$\gamma(BS)\gamma(BS) \rightarrow qqqq$	1.9%	1.7%
$\gamma(BS)\gamma(EPA) \rightarrow qqqq$	1.1%	1.0%
$\gamma(EPA)\gamma(BS) \rightarrow qqqq$	1.1%	1.0%
$\gamma(EPA)\gamma(EPA) \rightarrow qqqq$	0.7%	0.6%

**Table 7.21:** List of signal and background samples with the selection efficiencies of loose cuts and mutually exclusive cuts at  $\sqrt{s} = 3$  TeV. The selection efficiencies are presented in a “flow” fashion, as the every selection cut contains all the cuts to the left of it.

Channel / Efficiency $\sqrt{s} = 3 \text{ TeV}$	N	$\epsilon_{\text{presel}}$	$\epsilon_{\text{MVA}}$	$N_{\text{MVA}}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e, \text{ hadronic}$	146.0	61.7%	11.6%	9.89
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	355.0	18.8%	1.5%	1.05
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow \text{other}$	675.0	20.0%	3.6%	4.51
$e^-e^+ \rightarrow q_l q_l H\nu\bar{\nu}$	6115.4	36.0%	0.4%	9.42
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	2249.9	26.3%	0.5%	3.13
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	2197.7	25.8%	1.2%	6.82
$e^-e^+ \rightarrow qqqq$	1093000.0	1.4%	0.01%	1.43
$e^-e^+ \rightarrow qqqq\ell\ell$	338600.0	0.6%	-	-
$e^-e^+ \rightarrow qqqq\ell\nu$	213200.0	7.3%	0.05%	8.35
$e^-e^+ \rightarrow qqqq\nu\bar{\nu}$	143000.0	9.0%	0.05%	6.35
$e^-e^+ \rightarrow qq$	5897800.0	1.4%	-	-
$e^-e^+ \rightarrow qq\ell\nu$	11121800	0.1%	-	-
$e^-e^+ \rightarrow qq\ell\ell$	6639200.0	0.4%	-	-
$e^-e^+ \rightarrow qq\nu\bar{\nu}$	2635000.0	3.1%	-	-
$e^-\gamma(\text{BS}) \rightarrow e^-qqqq$	2004388.1	0.7%	-	-
$e^+\gamma(\text{BS}) \rightarrow e^+qqqq$	2002334.1	0.7%	-	-
$e^-\gamma(\text{EPA}) \rightarrow e^-qqqq$	575600.0	0.4%	-	-
$e^+\gamma(\text{EPA}) \rightarrow e^+qqqq$	575600.0	0.3%	-	-
$e^-\gamma(\text{BS}) \rightarrow \nu qqqq$	414750.0	16.8%	0.04%	30.7
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu} qqqq$	414434.0	15.9%	0.05%	30.3
$e^-\gamma(\text{EPA}) \rightarrow \nu qqqq$	108400.0	7.8%	0.04%	3.37
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu} qqqq$	108400.0	7.3%	0.03%	2.63
$e^-\gamma(\text{BS}) \rightarrow qqH\nu$	92588.0	30.2%	0.2%	67.5
$e^+\gamma(\text{BS}) \rightarrow qqH\nu$	92430.0	30.3%	0.2%	54.2
$e^-\gamma(\text{EPA}) \rightarrow qqH\nu$	23400.0	15.4%	0.2%	7.88
$e^+\gamma(\text{EPA}) \rightarrow qqH\nu$	23400.0	15.2%	0.3%	10.2
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow qqqq$	18009413.9	1.6%	-	-
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow qqqq$	3824548.1	1.0%	-	-
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow qqqq$	3828498.1	1.0%	-	-
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow qqqq$	805400.0	0.6%	-	-

**Table 7.22:** List of signal and background samples with the corresponding expected number at  $\sqrt{s} = 3 \text{ TeV}$ , for a luminosity of  $2000 \text{ fb}^{-1}$ . The number of events, selection efficiency of pre-selection, selection efficiency of MVA after pre-selection, number of events after MVA are shown. - represents a number less than 0.01.

## 7.12 $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$ semi-leptonic decay at $\sqrt{s} = 3$ TeV analysis

Before interpreting the results, the last discussion is on the  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu\bar{\nu}$  semi-leptonic decay at  $\sqrt{s} = 3$  TeV analysis. The  $\sqrt{s} = 1.4$  TeV analysis was performed as well. Since the selected event number is too low, there is not enough signal events to have a meaningful discussion. Hence the  $\sqrt{s} = 3$  TeV analysis is presented.

The strategy of the analysis is very similar to the hadronic decay analysis. The main difference are that there is one lepton in the final state and the final state has 4 quarks instead of 6.  $H_{bb}$  and  $W$  can not be reconstructed due to the leptonic decay of one of the  $W$ . Hence events are selected when there is one identified lepton, using the same lepton identifiers. The jet reinstruction parameters are the same as the  $\sqrt{s} = 3$  TeV hadronic decay analysis. The pre-selection cuts are the same, listed in table 7.22. There are no mutually exclusive cuts since there is no semi-leptonic analysis in the parallel analysis. The jet pairing still tries to reconstruct all the physical bosons. The same set of variables are used at  $\sqrt{s} = 3$  TeV. Parameters for the MVA BDT classifier is the same as the ones at  $\sqrt{s} = 3$  TeV. The efficiencies of the pre-selection cuts, the efficiencies of the MVA selections are listed in table 7.23, alongside with the number of events after the MVA selection. Since there are three neutrinos in the final state, reconstructing the correct event topology is more difficult. The MVA performance is worse and all most all background channels survived the MVA. Nevertheless, the dominant background are almost identical to those at  $\sqrt{s} = 1.4$  TeV. Hence see section 7.10 for the discussion.

Pre-selection	$\sqrt{s} = 3$ TeV
Discriminative pre-selection	$m_{HH} > 150$ GeV, $B_1 > 0.2$ , $p_{T_{HH}} > 30$ GeV
Loose cuts for MVA	$m_{H_{bb}} < 500$ GeV, $m_{HH} < 3000$ GeV

**Table 7.23:** Pre-selection cuts at  $\sqrt{s} = 3$  TeV.

Channel / Efficiency $\sqrt{s} = 1.4 \text{ TeV}$	N	$\epsilon_{\text{presel}}$	$\epsilon_{\text{MVA}}$	$N_{\text{MVA}}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , semi-leptonic	96.8	44.6%	21.9%	13.11
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}b\bar{b}\nu_e\bar{\nu}_e$	355.0	13.3%	10.9%	5.38
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow \text{other}$	724.2	13.1%	13.6%	12.75
$e^-e^+ \rightarrow q_l\bar{q}_lH\nu\bar{\nu}$	6115.4	7.4%	13.7%	62.63
$e^-e^+ \rightarrow c\bar{c}H\nu\bar{\nu}$	2249.9	6.3%	12.1%	17.10
$e^-e^+ \rightarrow b\bar{b}H\nu\bar{\nu}$	2197.7	15.9%	5.1%	18.03
$e^-e^+ \rightarrow qqqq$	1093000.0	0.6%	0.2%	15.04
$e^-e^+ \rightarrow qqqq\ell\ell$	338600.0	1.0%	0.06%	1.85
$e^-e^+ \rightarrow qqqq\ell\nu$	213200.0	27.6%	0.5%	270.33
$e^-e^+ \rightarrow qqqq\nu\bar{\nu}$	143000.0	1.9%	1.6%	43.78
$e^-e^+ \rightarrow qq$	5897800.0	0.4%	0.3%	60.82
$e^-e^+ \rightarrow qq\ell\nu$	11121800	0.3%	0.08%	21.24
$e^-e^+ \rightarrow qq\ell\ell$	6639200.0	0.6%	0.2%	84.14
$e^-e^+ \rightarrow qq\nu\nu$	2635000.0	0.4%	0.9%	92.55
$e^-\gamma(\text{BS}) \rightarrow e^-qqqq$	2004388.1	1.2%	-	-
$e^+\gamma(\text{BS}) \rightarrow e^+qqqq$	2002334.1	1.2%	-	-
$e^-\gamma(\text{EPA}) \rightarrow e^-qqqq$	575600.0	1.1%	-	-
$e^+\gamma(\text{EPA}) \rightarrow e^+qqqq$	575600.0	1.1%	-	-
$e^-\gamma(\text{BS}) \rightarrow \nu qqqq$	414750.0	3.7%	1.5%	226.77
$e^+\gamma(\text{BS}) \rightarrow \bar{\nu} qqqq$	414434.0	3.5%	1.6%	225.68
$e^-\gamma(\text{EPA}) \rightarrow \nu qqqq$	108400.0	11.2%	0.9%	107.90
$e^+\gamma(\text{EPA}) \rightarrow \bar{\nu} qqqq$	108400.0	10.7%	0.8%	92.75
$e^-\gamma(\text{BS}) \rightarrow qqH\nu$	92588.0	7.9%	10.7%	779.36
$e^+\gamma(\text{BS}) \rightarrow qqH\nu$	92430.0	7.9%	10.1%	741.57
$e^-\gamma(\text{EPA}) \rightarrow qqH\nu$	23400.0	22.9%	6.9%	369.52
$e^+\gamma(\text{EPA}) \rightarrow qqH\nu$	23400.0	22.7%	7.2%	381.33
$\gamma(\text{BS})\gamma(\text{BS}) \rightarrow qqqq$	18009413.9	0.4%	-	-
$\gamma(\text{BS})\gamma(\text{EPA}) \rightarrow qqqq$	3824548.1	1.0%	-	-
$\gamma(\text{EPA})\gamma(\text{BS}) \rightarrow qqqq$	3828498.1	1.0%	0.08%	28.85
$\gamma(\text{EPA})\gamma(\text{EPA}) \rightarrow qqqq$	805400.0	1.1%	-	-

**Table 7.24:** List of signal and background samples with the corresponding expected number at  $\sqrt{s} = 3 \text{ TeV}$ , for a luminosity of  $2000 \text{ fb}^{-1}$ . The number of events, selection efficiency of pre-selection, selection efficiency of MVA after pre-selection, number of events after MVA are shown. - represents no events passing the MVA.

Channel	$N_S$	$N_B$	$N_S/\sqrt{N_S + N_B}$
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , hadronic, $\sqrt{s} = 1.4 \text{ TeV}$	1.79	8.41	0.56
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , hadronic, $\sqrt{s} = 3 \text{ TeV}$	15.45	242.28	0.96
$e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e\bar{\nu}_e$ , semi-leptonic, $\sqrt{s} = 3 \text{ TeV}$	31.24	3612.39	0.52

**Table 7.25:** Number of signal and background events, and significance after MVA.

## 7.13 Result interpretation

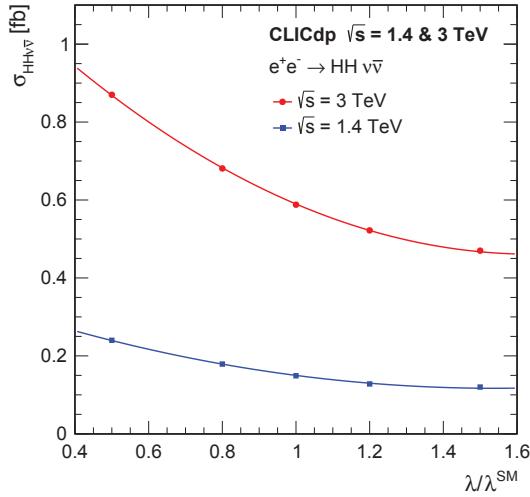
The results for  $\sqrt{s} = 1.4 \text{ TeV}$  and  $\sqrt{s} = 3 \text{ TeV}$  analysis are summarised in table 7.24.  $N_S$  is all  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  events passed the MVA. Combining the  $\sqrt{s} = 3 \text{ TeV}$  results, the expected precisions on the cross sections, which is roughly  $\sqrt{N_S + N_B}/N_S$ , at  $\sqrt{s} = 1.4 \text{ TeV}$  and  $3 \text{ TeV}$  are:

$$\frac{\Delta [\sigma(HH\nu_e\bar{\nu}_e)]}{\sigma(HH\nu_e\bar{\nu}_e)} = \begin{cases} 179\%, & \text{at } \sqrt{s} = 1.4 \text{ TeV}, \\ 92\%, & \text{at } \sqrt{s} = 3 \text{ TeV}. \end{cases} \quad (7.10)$$

As stated in the beginning of the chapter, the double Higgs production cross section is sensitive to the Higgs triple self coupling  $\lambda$ . The relative uncertainty on the coupling can be related to the uncertainty on the coupling via:

$$\frac{\Delta \lambda}{\lambda} \approx \kappa \cdot \frac{\Delta [\sigma(HH\nu_e\bar{\nu}_e)]}{\sigma(HH\nu_e\bar{\nu}_e)}. \quad (7.11)$$

$\kappa$  can be extracted by varying the  $\lambda$  and parameterise the cross section change at a general level. Figure 7.13 shows the cross section as a function of the coupling at generator level, at  $\sqrt{s} = 1.4 \text{ TeV}$  and  $\sqrt{s} = 3 \text{ TeV}$ . The negative gradient indicates that the dependence on  $\lambda$  experiences the destructive interference with other SM Feynman diagrams affecting. At the SM  $\lambda$  value, the  $\kappa$  is 1.22 and 1.47 at  $\sqrt{s} = 1.4 \text{ TeV}$  and  $3 \text{ TeV}$ , respectively. Since  $\kappa$  is extracted from the relation at generator level, the fully simulated reconstruction



**Figure 7.13:** Cross section for the  $e^-e^+ \rightarrow HH\nu_e \bar{\nu}_e$  process as a function of the ratio  $\lambda/\lambda_{SM}$  at  $\sqrt{s} = 1.4 \text{ TeV}$  and  $3 \text{ TeV}$ , taken from [15].

selection, as presented above, may favours certain Feynman diagrams, hence affecting the sensitivity to  $\lambda$ .

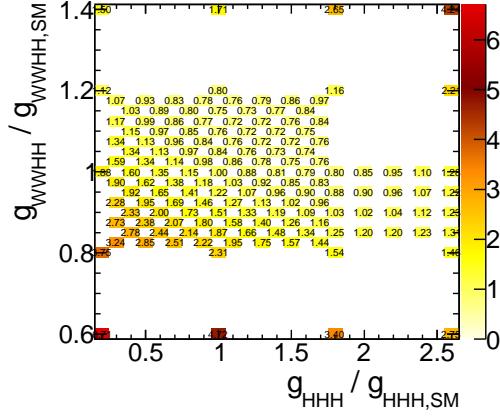
Without electron polarisation, the uncertainty on the Higgs triple self coupling  $\lambda$ , from  $e^-e^+ \rightarrow HH\nu_e \bar{\nu}_e \rightarrow b\bar{b}W^+W^-\nu_e \bar{\nu}_e$  analysis is:

$$\frac{\Delta\lambda}{\lambda} = \begin{cases} 218\%, & \text{at } \sqrt{s} = 1.4 \text{ TeV}, \\ 135\%, & \text{at } \sqrt{s} = 3 \text{ TeV}. \end{cases} \quad (7.12)$$

Since the Feynman diagrams for the double Higgs boson productions include t-channel WW-fusion, the cross section can be enhanced by polarised electron beam. For  $P(e^-) = 80\%$ , the uncertainty on  $\lambda$  becomes:

$$\frac{\Delta\lambda}{\lambda} = \begin{cases} 163\%, & \text{at } \sqrt{s} = 1.4 \text{ TeV}, \\ 97\%, & \text{at } \sqrt{s} = 3 \text{ TeV}. \end{cases} \quad (7.13)$$

When both  $\sqrt{s}$  are combined, the statistical precision on  $\lambda$  increases to 99% for the unpolarised beam, and 87% for the polarised beam with  $P(e^-) = 80\%$ .



**Figure 7.14:** Normalised cross section for the  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  process as a function of the  $g_{HHH}/g_{HHSMSM}$  and  $g_{WWWH}/g_{WWSMSM}$  at  $\sqrt{s} = 3$  TeV.

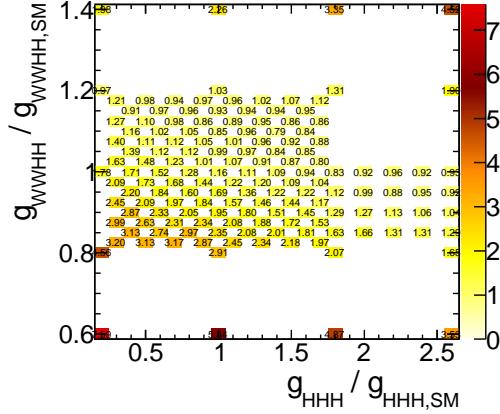
## 7.14 Simultaneous couplings extraction

The double higgs production,  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ , can occur via processes in figure 7.1. As stated in the beginning, Figure 7.1a is sensitive to Higgs triple self coupling  $g_{HHH}$ . Figure 7.1b is sensitive to quartic coupling  $g_{WWWH}$ . Therefore an simultaneous extraction on the coupling uncertainty can be performed by extending the method in the section 7.13. Once a relationship between  $g_{HHH}$ ,  $g_{WWWH}$  and significance is established, a contour of the uncertainty in  $g_{HHH}$  and  $g_{WWWH}$  two dimensional phase space can be obtained.

This two dimensional template fitting is performed at  $\sqrt{s} = 3$  TeV, as the precision at  $\sqrt{s} = 1.4$  TeV is too low to support such fitting. The luminosity is assumed to be  $3000\text{fb}^{-1}$  to reflect the updated CLIC running scenario.

The normalised cross section of the  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  as a function of  $g_{HHH}$  and  $g_{WWWH}$  is shown in figure 7.14. The SM cross section is normalised to 1. Around the SM value, the cross section increases with the decrease of  $g_{HHH}$  and the increase of  $g_{WWWH}$ . Therefore, the cross sections along the anti-diagonal do not vary much, which would be difficult to precisely determine the statistical uncertainty on the coupling measurements.

To determine the uncertainty on the coupling measurements, the variables proposed in the theoretical study in section 2.9 are used: the invariant mass of the two Higgs system,  $m_{HH}$ , and the sum of their transverse momenta,  $H_T$ . By choosing the kinematic bins, high-energy behaviour can be disentailed from the physics at threshold, allowing the extraction of the coupling strength  $g_{WWWH}$  and  $g_{HHH}$ .



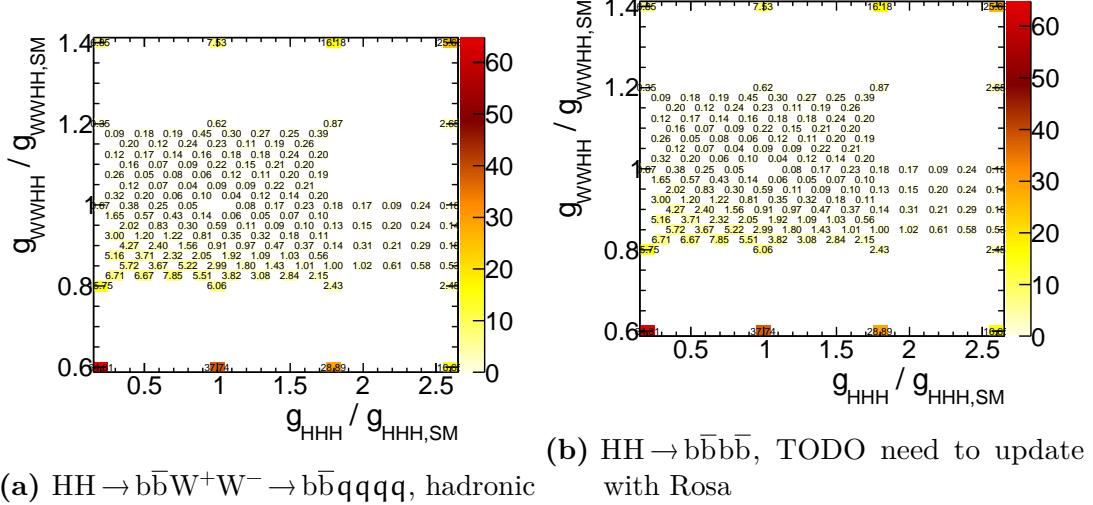
**Figure 7.15:** The significance for the  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  process as a function of the  $g_{HHH}/g_{HHSMSM}$  and  $g_{WWHH}/g_{WWSMSM}$  at  $\sqrt{s} = 3$  TeV, using subchannel hadronic decay  $HH \rightarrow b\bar{b}W^+W^-$ , assuming a luminosity of  $3000\text{fb}^{-1}$ .

The strategy for coupling extraction is described below. Simulated events with non-SM couplings are generated and reconstructed. These events went through the analysis chain discussed in this chapter, with the same cuts and the MVA classifier, trained with the SM coupling sample. The significance of the double higgs production with subchannel hadronic decay  $HH \rightarrow b\bar{b}W^+W^-$  as a function of  $g_{HHH}$  and  $g_{WWHH}$  is shown in figure 7.15.

The selected events are classified into 8 kinematic bins. 2 bins in  $H_T$  are cut at 200 GeV. 4 bins in  $m_{HH}$  are cut at 500, 700, 1000 GeV. A  $\chi^2$  function is constructed to access the difference of the  $m_{HH}$  and  $H_T$  distributions for non-SM coupling comparing to SM coupling sample.  $\chi^2$  is:

$$\chi^2 = \sum_i^{\text{bins}} \frac{(N_i - N_{i,\text{observed}})^2}{N_i}, \quad (7.14)$$

where  $N_i$  is the number of event expected in a kinematic bin  $i$  for a non-SM coupling sample.  $N_{i,\text{observed}}$  is the number of event observed in a kinematic bin  $i$ . Here the observed set is the SM coupling sample. The expression is summing over all kinematic bins. By construction, the SM coupling point has a  $\chi^2$  of 0. Figure 7.16 shows the  $\chi^2$  as a function of  $g_{HHH}$  and  $g_{WWHH}$ , for two subchannels, hadronic decay  $HH \rightarrow b\bar{b}W^+W^-$  and  $HH \rightarrow b\bar{b}b\bar{b}$ . The  $\chi^2$  values for the  $HH \rightarrow b\bar{b}b\bar{b}$  subchannel is larger, which is due to larger significance obtained with that subchannel.

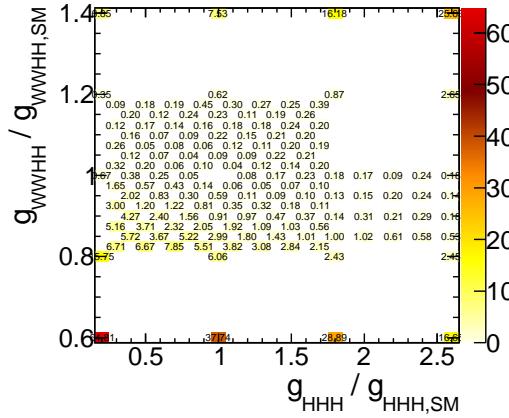


**Figure 7.16:** The  $\chi^2$  for the  $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$  process as a function of the  $g_{HHH}/g_{HHH,\text{SM}}$  and  $g_{WWHH}/g_{WWHH,\text{SM}}$  at  $\sqrt{s} = 3 \text{ TeV}$ , using subchannel hadronic decay  $HH \rightarrow b\bar{b}W^+W^-$  and subchannel, assuming a luminosity of  $3000 \text{ fb}^{-1}$ .

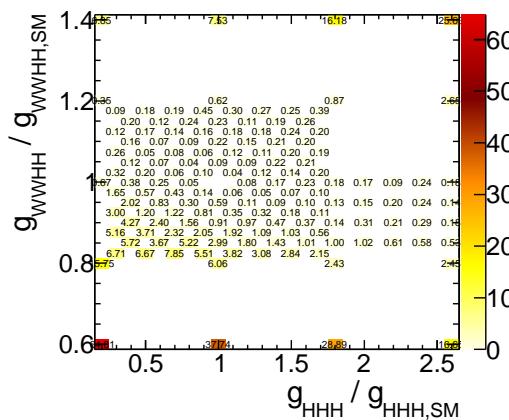
Two subchannels, hadronic decay  $HH \rightarrow b\bar{b}W^+W^-$  and  $HH \rightarrow b\bar{b}b\bar{b}$ , are combined to increase the statistical precision on the coupling measurements. Two  $\chi^2$  surfaces are summed. To avoid the statistical fluctuation in the sample, a toy MC experiment is performed. The SM coupling samples is treated as data template set. 100000 data sets are generated by fluctuate the event number in each kinematic bin according to poisson distribution. The  $\chi^2$  is performed and summed using these data sets as the observed data. The summed  $\chi$  is averaged by the number of data sets (100000), and normalised such that the  $\chi^2$  at the SM coupling is 0. Since only the difference between the non-SM and SM  $\chi^2$  is used for the coupling measurements, the normalisation does not affect the measurements and helps to ease the visualisation. Figure 7.17 shows the normalised  $\chi^2$ , after averaging the toy MC experiments, as a function of  $g_{HHH}/g_{HHH,\text{SM}}$  and  $g_{WWHH}/g_{WWHH,\text{SM}}$ . The  $\chi^2$  is flat along the anti-diagonal, which is similar to the cross section plot.

Since there are two couplings in this  $\chi^2$  surface, the degree of freedom for this fit is 2. A contour of 68% confidence ( $\chi^2 = 2.3$ ) can be drawn by interpolating between points on the surface. Figure 7.18 shows the contour. The counter can be sliced one dimensionally to extract the uncertainty of one coupling for a given value of the other coupling. For example:

$$\frac{\Delta g_{WWHH}}{g_{WWHH}} \simeq 5\% \text{ for } g_{HHH} = g_{HHH,\text{SM}} \quad (7.15)$$



**Figure 7.17:** TODO update Normalised  $\chi^2$ , after averaging the toy MC experiments, as a function of  $g_{HHH}/g_{HHH,SM}$  and  $g_{WWHH}/g_{WWHH,SM}$ , combining subchannel hadronic decay  $HH \rightarrow bbW^+W^-$  and subchannel, assuming a luminosity of  $3000\text{fb}^{-1}$



**Figure 7.18:** TODO update Contour plot of  $\chi^2$ , after averaging the toy MC experiments, as a function of  $g_{HHH}/g_{HHH,SM}$  and  $g_{WWHH}/g_{WWHH,SM}$ , combining subchannel hadronic decay  $HH \rightarrow bbW^+W^-$  and subchannel, assuming a luminosity of  $3000\text{fb}^{-1}$

$$\frac{\Delta g_{\text{HHH}}}{g_{\text{HHH}}} \simeq 33\% \text{ for } g_{\text{WWHH}} = g_{\text{WWHHSM}} \quad (7.16)$$

The statistical precisions on  $g_{\text{WWHH}}$  and  $g_{\text{HHH}}$  are much better than that with LHC, or high luminosity upgraded LHC [12].

# Colophon

This thesis was made in L<sup>A</sup>T<sub>E</sub>X 2 <sub>$\epsilon$</sub>  using the “heptesis” class [65].



# Bibliography

- [1] ATLAS Collaboration, G. Aad *et al.*, Phys.Lett. **B716**, 1 (2012), 1207.7214.
- [2] Particle Data Group, K. A. Olive *et al.*, Chin. Phys. **C38**, 090001 (2014).
- [3] M. Thomson, *Modern particle physics* (Cambridge University Press, New York, 2013).
- [4] D. Tong, Lectures on quantum field theory, 2006.
- [5] B. Gripaios, Lectures on gauge field theory, 2017.
- [6] SLD Electroweak Group, DELPHI, ALEPH, SLD, SLD Heavy Flavour Group, OPAL, LEP Electroweak Working Group, L3, S. Schael *et al.*, Phys. Rept. **427**, 257 (2006), hep-ex/0509008.
- [7] S. Weinberg, Phys. Rev. Lett. **19**, 1264 (1967).
- [8] D. Rainwater, Searching for the Higgs boson, in *Proceedings of Theoretical Advanced Study Institute in Elementary Particle Physics : Exploring New Frontiers Using Colliders and Neutrinos (TASI 2006): Boulder, Colorado, June 4-30, 2006*, pp. 435–536, 2007, hep-ph/0702124.
- [9] D. B. Kaplan and H. Georgi, Phys. Lett. **B136**, 183 (1984).
- [10] W. D. Goldberger, B. Grinstein, and W. Skiba, Phys. Rev. Lett. **100**, 111802 (2008), 0708.1463.
- [11] G. F. Giudice, C. Grojean, A. Pomarol, and R. Rattazzi, JHEP **06**, 045 (2007), hep-ph/0703164.
- [12] R. Contino, C. Grojean, M. Moretti, F. Piccinini, and R. Rattazzi, JHEP **05**, 089 (2010), 1002.1011.
- [13] R. Contino, C. Grojean, D. Pappadopulo, R. Rattazzi, and A. Thamm, JHEP **02**,

- 006 (2014), 1309.7038.
- [14] V. Barger, T. Han, P. Langacker, B. McElrath, and P. Zerwas, Phys. Rev. **D67**, 115001 (2003), hep-ph/0301097.
- [15] H. Abramowicz *et al.*, (2016), 1608.07538.
- [16] CMS Collaboration, S. Chatrchyan *et al.*, Phys.Lett. **B716**, 30 (2012), 1207.7235.
- [17] J. Brau *et al.*, (2007).
- [18] L. Linssen, A. Miyamoto, M. Stanitzki, and H. Weerts, (2012), 1202.5940.
- [19] H. Baer *et al.*, (2013), 1306.6352.
- [20] M. Aicheler *et al.*, (2012).
- [21] M. Thomson, Nucl.Instrum.Meth. **A611**, 25 (2009), 0907.3577.
- [22] J. S. Marshall, A. Mílznich, and M. A. Thomson, Nucl. Instrum. Meth. **A700**, 153 (2013), 1209.4039.
- [23] I. G. Knowles and G. D. Lafferty, J. Phys. **G23**, 731 (1997), hep-ph/9705217.
- [24] M. Green, *Electron-Positron Physics at the ZStudies in high energy physics, cosmology, and gravitation* (Taylor & Francis, 1998).
- [25] CALICE, C. Adloff, (2011), 1105.0511.
- [26] B. Parker *et al.*, (2009).
- [27] CALICE, JINST **7**, P04015 (2012), 1201.1653.
- [28] A. Sailer, *Radiation and Background Levels in a CLIC Detector Due to Beam-beam Effects: Optimisation of Detector Geometries and Technologies* (Humboldt Universität zu Berlin, Mathematisch-Naturwissenschaftliche Fakultät I, 2012).
- [29] W. Kilian, T. Ohl, and J. Reuter, European Physical Journal C **71** (2011).
- [30] M. Moretti, T. Ohl, and J. Reuter, p. 1981 (2001), hep-ph/0102195.
- [31] T. Sjostrand, (1995), hep-ph/9508391.
- [32] OPAL, G. Alexander *et al.*, Z. Phys. **C69**, 543 (1996).
- [33] S. Jadach, Z. Was, R. Decker, and J. H. Kuhn, Comput. Phys. Commun. **76**, 361

- (1993).
- [34] D. Schulte, (1999).
  - [35] GEANT4, S. Agostinelli *et al.*, Nucl.Instrum.Meth. **A506**, 250 (2003).
  - [36] P. Mora de Freitas and H. Videau, p. 623 (2002).
  - [37] F. Gaede, Nucl. Instrum. Meth. **A559**, 177 (2006).
  - [38] J. S. Marshall and M. A. Thomson, Eur. Phys. J. **C75**, 439 (2015), 1506.05348.
  - [39] J. S. Marshall, Presentation on pandorapfa with lc reconstruction, [https://github.com/PandoraPFA/Documentation/blob/master/Pandora\\_LC\\_Reconstruction.pdf](https://github.com/PandoraPFA/Documentation/blob/master/Pandora_LC_Reconstruction.pdf), 2017.
  - [40] F. Gaede and J. Engels, EUDET Report (2007).
  - [41] B. Xu, Improvement of photon reconstruction in PandoraPFA, in *Proceedings, International Workshop on Future Linear Colliders (LCWS15): Whistler, B.C., Canada, November 02-06, 2015*, 2016, 1603.00013.
  - [42] E. Farhi, Phys. Rev. Lett. **39**, 1587 (1977).
  - [43] TMVA Core Developer Team, J. Therhaag, AIP Conf.Proc. **1504**, 1013 (2009).
  - [44] S. Dittmaier *et al.*, (2012), 1201.3084.
  - [45] A. Míznich, CERN Report No. LCD-Note-2010-009, 2010 (unpublished).
  - [46] CLICdp, A. Sailer and A. Sapronov, (2017), 1702.06945.
  - [47] S. Lukić, Forward electron tagging in the  $h \rightarrow \mu \mu$  analysis at 1.4 tev, <http://indico.cern.ch/event/262809/contributions/1595499/attachments/464689/643931/electronTagging.pdf>, 2013.
  - [48] G. Milutinović-Dumbelović *et al.*, (2014), 1412.5791.
  - [49] C. Grefe, T. Lastovicka, and J. Strube, Light Higgs Studies for the CLIC CDR, in *Helmholtz Alliance Linear Collider Forum: Proceedings of the Workshops Hamburg, Munich, Hamburg 2010-2012, Germany*, pp. 258–264, Hamburg, 2013, DESY, DESY, 1205.3908.
  - [50] G. F. Sterman and S. Weinberg, Phys. Rev. Lett. **39**, 1436 (1977).

- 
- [51] S. Moretti, L. Lonnblad, and T. Sjostrand, JHEP **08**, 001 (1998), hep-ph/9804296.
  - [52] G. P. Salam, Eur. Phys. J. **C67**, 637 (2010), 0906.1833.
  - [53] A. Ali and G. Kramer, Eur. Phys. J. **H36**, 245 (2011), 1012.2288.
  - [54] M. Cacciari, G. P. Salam, and G. Soyez, Eur. Phys. J. **C72**, 1896 (2012), 1111.6097.
  - [55] M. Cacciari and G. P. Salam, Phys. Lett. **B641**, 57 (2006), hep-ph/0512210.
  - [56] S. Catani, Y. L. Dokshitzer, M. H. Seymour, and B. R. Webber, Nucl. Phys. **B406**, 187 (1993).
  - [57] S. D. Ellis and D. E. Soper, Phys. Rev. **D48**, 3160 (1993), hep-ph/9305266.
  - [58] M. Battaglia and F. P., CERN Report No. LCD-Note-2010-006, 2010 (unpublished).
  - [59] M. Boronat, J. Fuster, I. Garcia, E. Ros, and M. Vos, Phys. Lett. **B750**, 95 (2015), 1404.4294.
  - [60] T. Suehara and T. Tanabe, Nucl. Instrum. Meth. **A808**, 109 (2016), 1506.08371.
  - [61] Linear Collider ILD Concept Group -, T. Abe *et al.*, (2010), 1006.3396.
  - [62] H. Aihara *et al.*, (2009), 0911.0006.
  - [63] A. Hocker *et al.*, PoS **ACAT**, 040 (2007), physics/0703039.
  - [64] G. Hanson *et al.*, Phys. Rev. Lett. **35**, 1609 (1975).
  - [65] A. Buckley, The heptesis L<sup>A</sup>T<sub>E</sub>X class.

# List of figures

2.1	SM Higgs boson decay width and branching ratios . . . . .	10
4.1	Illustration of the cone clustering algorithm, taken from [39] . . . . .	34
4.2	Example of MVA overtraining . . . . .	43
4.3	Example of a decision tree . . . . .	45
5.1	Two 500 GeV photons (yellow and blue), just resolved in the transverse plane perpendicular to the direction of the flight, of their energy deposition in electromagnetic calorimeter. U and V axis are two arbitrary axis perpendicular to each other in the plane. Z axis is the sum of the calorimeter hit energy in each particular bin in 2D plane in GeV. . . . .	54
5.3	An event display of a typical 500 GeV photon, reconstructed into a main photon in the ECAL (yellow) and a neutral hadron fragment in the HCAL (blue). . . . .	62
7.1	Feynman diagrams of leading-order $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ processes at CLIC . . . . .	87
7.2	BeamCAL and LumiCAL electron tagging efficiency. . . . .	96
7.3	Example MC mass fit for jet optimisation in double Higgs analysis . . . . .	100
7.4	Fitted mass, and resolution of $H_{bb}$ , $H_{WW^*}$ and $W$ for $\sqrt{s} = 1.4$ TeV . . . . .	101
7.5	Performance of b-jet tagging with training samples at $\sqrt{s} = 1.4$ TeV. . . . .	103
7.6	Discriminative pre-selection variables for $\sqrt{s} = 1.4$ TeV. . . . .	106
7.7	Sum of b tag against $y_{34}$ at $\sqrt{s} = 1.4$ TeV . . . . .	108

---

7.8	Stacked plots for discriminative variables for the MVA at $\sqrt{s} = 1.4$ TeV after all pre-selection cuts. . . . .	112
7.9	Fitted mass and relative mass resolution of $H_{bb}$ , $H_{WW^*}$ and $W$ at $\sqrt{s} = 3$ TeV. . . . .	118
7.10	Performance of b-jet tagging with training samples at $\sqrt{s} = 3$ TeV. . . . .	119
7.11	Discriminative pre-selection variables at $\sqrt{s} = 3$ TeV. . . . .	120
7.12	Sum of b tag against $y_{34}$ at $\sqrt{s} = 3$ TeV . . . . .	122
7.13	Cross section for the $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ process as a function of the ratio $\lambda/\lambda_{SM}$ . . . . .	128
7.14	Normalised cross section for the $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ process as a function of the $g_{HHH}/g_{HHHSM}$ and $g_{WWHH}/g_{WWHHS}$ at $\sqrt{s} = 3$ TeV. . . . .	129
7.15	The significance for the $e^-e^+ \rightarrow HH\nu_e\bar{\nu}_e$ process as a function of the $g_{HHH}/g_{HHHSM}$ and $g_{WWHH}/g_{WWHHS}$ at $\sqrt{s} = 3$ TeV, using subchannel hadronic decay $HH \rightarrow b\bar{b}W^+W^-$ , assuming a luminosity of $3000\text{fb}^{-1}$ . . . . .	130
7.16	$\chi^2$ as a function of $g_{HHH}/g_{HHHSM}$ and $g_{WWHH}/g_{WWHHS}$ at $\sqrt{s} = 3$ TeV . . . . .	131
7.17	TODO update Normalised $\chi^2$ , after averaging the toy MC experiments, as a function of $g_{HHH}/g_{HHHSM}$ and $g_{WWHH}/g_{WWHHS}$ , combining subchannel hadronic decay $HH \rightarrow b\bar{b}W^+W^-$ and subchannel, assuming a luminosity of $3000\text{fb}^{-1}$ . . . . .	131
7.18	TODO update Contour plot of $\chi^2$ , after averaging the toy MC experiments, as a function of $g_{HHH}/g_{HHHSM}$ and $g_{WWHH}/g_{WWHHS}$ , combining subchannel hadronic decay $HH \rightarrow b\bar{b}W^+W^-$ and subchannel, assuming a luminosity of $3000\text{fb}^{-1}$ . . . . .	132

# List of tables

4.1	The attribute of samples for the decision tree example.	46
4.2	A toy example to demonstrate definitions of efficiency and purity.	50
6.1	Branching ratios of the seven major $\tau^-$ decays, taken from [2]. $\tau^+$ decays similarly to $\tau^-$ .	73
7.1	Signal and background samples with the corresponding cross sections at $\sqrt{s} = 1.4 \text{ TeV}$ .	89
7.2	Optimised selection criterion for IsolatedLeptonFinderProcessor	91
7.3	Optimised selection criterion for IsolatedLeptonFinderProcessor	92
7.4	Optimised selection criterion for TauFinderProcessor	93
7.5	Optimised selection criterion for BonoTauFinderProcessor	94
7.6	Isolated lepton finder processors performance on the signal and selected background samples at $\sqrt{s} = 1.4 \text{ TeV}$ .	97
7.7	Very forward electron and photon finder performance on the signal and selected background samples at $\sqrt{s} = 1.4 \text{ TeV}$ .	97
7.8	The extracted fitted parameters of optimal jet reconstructions, at $\sqrt{s} = 1.4 \text{ TeV}$	102
7.9	Pre-selection cuts at $\sqrt{s} = 1.4 \text{ TeV}$ .	105
7.10	Pre-selection efficiency at $\sqrt{s} = 1.4 \text{ TeV}$ .	107
7.11	Mutually exclusive cuts at $\sqrt{s} = 1.4 \text{ TeV}$ .	109

7.12 List of signal and background samples of loose cuts and mutually exclusive cuts at $\sqrt{s} = 1.4$ TeV . . . . .	110
7.13 Variables used in MVA at $\sqrt{s} = 1.4$ TeV . . . . .	111
7.14 Signal and background selection efficiency and event numbers at $\sqrt{s} = 1.4$ TeV . . . . .	115
7.15 Cross sections of samples at $\sqrt{s} = 3$ TeV. . . . .	116
7.16 Isolated lepton finder processors performance on the signal and selected background samples at $\sqrt{s} = 3$ TeV. . . . .	117
7.17 The extracted fitted parameters of optimal jet reconstructions at $\sqrt{s} = 3$ TeV. . . . .	119
7.18 Pre-selection cuts at $\sqrt{s} = 3$ TeV. . . . .	120
7.19 Signal and background events with selection efficiency and event numbers after the MVA at $\sqrt{s} = 3$ TeV . . . . .	121
7.20 List of signal and background samples of loose cuts and mutually exclusive cuts at $\sqrt{s} = 3$ TeV. . . . .	123
7.21 List of signal and background selection efficiencies and event numbers after MVA at $\sqrt{s} = 3$ TeV. . . . .	124
7.22 Pre-selection cuts at $\sqrt{s} = 3$ TeV. . . . .	125
7.23 List of signal and background selection efficiencies and event numbers after MVA for semi-leptonic analysis at $\sqrt{s} = 3$ TeV. . . . .	126
7.24 Number of signal and background events, and significance after MVA. . .	127