

1

2

# THÈSE DE DOCTORAT DE

NANTES UNIVERSITÉ

ÉCOLE DOCTORALE N° 596  
*Matière, Molécules, Matériaux*  
Spécialité : *Physique des particules*

Par

**Léonard Imbert**

**Deep learning methods and Dual Calorimetric analysis for high precision neutrino oscillation measurements at JUNO**

Thèse présentée et soutenue à Nantes, le 2 Decembre 2024

Unité de recherche : Laboratoire SUBATECH, UMR 6457

## Rapporteurs avant soutenance :

Christine Marquet Directrice de recherche au CNRS, LP2I Bordeaux  
David Rousseau Directeur de recherche au CNRS, IJCLab

## Composition du Jury :

|                    |                         |  |
|--------------------|-------------------------|--|
| Président :        | Barbara Erazmus         | Directrice de recherche au CNRS, Subatech        |
| Examinateurs :     | Juan Pedro Ochoa-Ricoux | Full Professor, University of California, Irvine |
|                    | Yasmine Amhis           | Directrice de recherche au CNRS, IJCLab          |
|                    | Christine Marquet       | Directrice de recherche au CNRS, LP2I Bordeaux   |
|                    | David Rousseau          | Directeur de recherche au CNRS, IJCLab           |
| Dir. de thèse :    | Frédéric Yermia         | Professeur des universités, Nantes Université    |
| Co-dir. de thèse : | Benoit Viaud            | Chargé de recherche au CNRS, Subatech            |



# <sup>3</sup> Contents

|               |  |           |
|---------------|--|-----------|
| <sup>4</sup>  | <b>Contents</b>  | <b>1</b>  |
| <sup>5</sup>  | <b>Remerciements</b>   | <b>5</b>  |
| <sup>6</sup>  | <b>Introduction</b>  | <b>7</b>  |
| <sup>7</sup>  | <b>1 Neutrino physics</b>  | <b>9</b>  |
| <sup>8</sup>  | <b>1.1 Standard model</b> . . . . .  | <b>9</b>  |
| <sup>9</sup>  | <b>1.1.1 Limits of the standard model</b> . . . . .                                  | <b>9</b>  |
| <sup>10</sup> | <b>1.2 Historic of the neutrino</b> . . . . .  | <b>9</b>  |
| <sup>11</sup> | <b>1.3 Oscillation</b> . . . . .   | <b>9</b>  |
| <sup>12</sup> | <b>1.3.1 Phenomologies</b> . . . . .   | <b>9</b>  |
| <sup>13</sup> | <b>1.4 Open questions</b> . . . . .  | <b>9</b>  |
| <sup>14</sup> | <b>2 The JUNO experiment</b>   | <b>11</b> |
| <sup>15</sup> | <b>2.1 Neutrinos physics in JUNO</b> . . . . .                                       | <b>12</b> |
| <sup>16</sup> | <b>2.1.1 Reactor neutrino oscillation for NMO and precise measurements</b> . . . . . | <b>12</b> |
| <sup>17</sup> | <b>2.1.2 Other physics</b> . . . . .   | <b>15</b> |
| <sup>18</sup> | <b>2.2 The JUNO detector</b> . . . . .   | <b>17</b> |
| <sup>19</sup> | <b>2.2.1 Detection principle</b> . . . . .   | <b>17</b> |
| <sup>20</sup> | <b>2.2.2 Central Detector (CD)</b> . . . . .   | <b>19</b> |
| <sup>21</sup> | <b>2.2.3 Veto detector</b> . . . . .   | <b>23</b> |
| <sup>22</sup> | <b>2.3 Calibration strategy</b> . . . . .  | <b>23</b> |
| <sup>23</sup> | <b>2.3.1 Energy scale calibration</b> . . . . .                                      | <b>24</b> |
| <sup>24</sup> | <b>2.3.2 Calibration system</b> . . . . .  | <b>25</b> |
| <sup>25</sup> | <b>2.3.3 Instrumental non-linearity calibration</b> . . . . .                        | <b>25</b> |
| <sup>26</sup> | <b>2.4 Satellite detectors</b> . . . . .   | <b>26</b> |
| <sup>27</sup> | <b>2.4.1 TAO</b> . . . . .   | <b>26</b> |
| <sup>28</sup> | <b>2.4.2 OSIRIS</b> . . . . .  | <b>26</b> |
| <sup>29</sup> | <b>2.5 Software</b> . . . . .  | <b>27</b> |
| <sup>30</sup> | <b>2.6 State of the art of the Offline IBD reconstruction in JUNO</b> . . . . .      | <b>28</b> |
| <sup>31</sup> | <b>2.6.1 Interaction vertex reconstruction</b> . . . . .                             | <b>28</b> |
| <sup>32</sup> | <b>2.6.2 Energy reconstruction</b> . . . . .   | <b>33</b> |
| <sup>33</sup> | <b>2.6.3 Machine learning for reconstruction</b> . . . . .                           | <b>36</b> |
| <sup>34</sup> | <b>2.7 JUNO sensitivity to NMO and precise measurements</b> . . . . .                | <b>38</b> |

|    |          |   |           |
|----|----------|---|-----------|
| 35 | 2.7.1    | Theoretical spectrum . . . . .  | 38        |
| 36 | 2.7.2    | Fitting procedure . . . . .   | 39        |
| 37 | 2.7.3    | Physics results . . . . .   | 39        |
| 38 | 2.8      | Summary . . . . .   | 40        |
| 39 | <b>3</b> | <b>Machine learning: Introduction to the methods and algorithms used in this thesis</b> | <b>41</b> |
| 40 | 3.1      | Core concepts in machine learning and neural networks . . . . .                         | 42        |
| 41 | 3.2      | Boosted Decision Tree (BDT) . . . . .   | 42        |
| 42 | 3.2.1    | Artificial Neural Network (NN) . . . . .  | 42        |
| 43 | 3.2.2    | Training procedure . . . . .  | 44        |
| 44 | 3.2.3    | Potential pitfalls . . . . .  | 47        |
| 45 | 3.3      | Neural networks architectures . . . . .   | 50        |
| 46 | 3.3.1    | Fully Connected Deep Neural Network (FCDNN) . . . . .                                   | 50        |
| 47 | 3.3.2    | Convolutional Neural Network (CNN) . . . . .  | 50        |
| 48 | 3.3.3    | Graph Neural Network (GNN) . . . . .  | 52        |
| 49 | 3.3.4    | Adversarial Neural Network (ANN) . . . . .  | 54        |
| 50 | <b>4</b> | <b>Image recognition for IBD reconstruction with the SPMT system</b>                    | <b>55</b> |
| 51 | 4.1      | Motivations . . . . .   | 56        |
| 52 | 4.2      | Method and model . . . . .  | 56        |
| 53 | 4.2.1    | Model . . . . .   | 56        |
| 54 | 4.2.2    | Data representation . . . . .   | 58        |
| 55 | 4.2.3    | Dataset . . . . .   | 59        |
| 56 | 4.2.4    | Data characteristics . . . . .  | 60        |
| 57 | 4.3      | Training . . . . .  | 61        |
| 58 | 4.4      | Results . . . . .   | 61        |
| 59 | 4.4.1    | J21 results . . . . .   | 63        |
| 60 | 4.4.2    | J21 Combination of classic and ML estimator . . . . .                                   | 66        |
| 61 | 4.4.3    | J23 results . . . . .   | 68        |
| 62 | 4.5      | Conclusion and prospect . . . . .   | 69        |
| 63 | <b>5</b> | <b>Graph representation of JUNO for IBD reconstruction</b>                              | <b>73</b> |
| 64 | 5.1      | Motivation . . . . .  | 73        |
| 65 | 5.2      | Data representation . . . . .   | 74        |
| 66 | 5.3      | Message passing algorithm . . . . .   | 76        |
| 67 | 5.4      | Data . . . . .  | 78        |
| 68 | 5.5      | Model . . . . .   | 79        |
| 69 | 5.6      | Training . . . . .  | 80        |
| 70 | 5.7      | Optimization . . . . .  | 80        |
| 71 | 5.8      | Results . . . . .   | 81        |
| 72 | 5.9      | Conclusion . . . . .  | 82        |
| 73 | <b>6</b> | <b>Reliability of machine learning methods</b>  | <b>85</b> |
| 74 | 6.1      | Motivations . . . . .   | 86        |

|     |   |     |
|-----|---|-----|
| 75  | <b>6.2 Method</b>   | 86  |
| 76  | <b>6.3 Architecture</b>   | 86  |
| 77  | 6.3.1 Adversarial Neural Network  | 86  |
| 78  | 6.3.2 Reconstruction Network  | 87  |
| 79  | 6.3.3 Training  | 87  |
| 80  | <b>6.4 Results</b>  | 87  |
| 81  | 6.4.1 Back to identity  | 88  |
| 82  | 6.4.2 Breaking of the reconstruction  | 88  |
| 83  | <b>6.5 Conclusion and prospect</b>  | 88  |
| 84  | <b>7 Joint fit between the SPMT and LPMT spectra</b>                          | 89  |
| 85  | <b>7.1 Motivations</b>  | 90  |
| 86  | 7.1.1 Discrepancies between the SPMT and LPMT results                         | 90  |
| 87  | 7.1.2 Charge Non-Linearity (QNL)  | 91  |
| 88  | <b>7.2 Approach</b>   | 92  |
| 89  | 7.2.1 Data production   | 92  |
| 90  | 7.2.2 Individual fits   | 93  |
| 91  | 7.2.3 Joint fit   | 94  |
| 92  | 7.2.4 Data and theoretical spectrum generation                                | 96  |
| 93  | 7.2.5 Limitations   | 96  |
| 94  | <b>7.3 Fit software</b>   | 97  |
| 95  | 7.3.1 IBD generator   | 97  |
| 96  | 7.3.2 Fit   | 99  |
| 97  | <b>7.4 Technical challenges and development</b>                               | 99  |
| 98  | <b>7.5 Results</b>  | 100 |
| 99  | 7.5.1 Validation  | 100 |
| 100 | 7.5.2 Covariance matrix   | 104 |
| 101 | 7.5.3 Statistical tests   | 108 |
| 102 | <b>7.6 Conclusion and perspectives</b>  | 110 |
| 103 | <b>8 Conclusion</b>   | 115 |
| 104 | <b>A Calculation of optimal <math>\alpha</math> for estimator combination</b> | 117 |
| 105 | A.1 Unbiased estimator  | 117 |
| 106 | A.2 Optimal variance estimator  | 117 |
| 107 | <b>B Charge spherical harmonics analysis</b>                                  | 119 |
| 108 | <b>C Additional spectrum smearing</b>   | 127 |
| 109 | <b>List of Tables</b>   | 129 |
| 110 | <b>List of Figures</b>  | 137 |
| 111 | <b>List of Abbreviations</b>  | 139 |



<sup>113</sup> Remerciements



<sup>114</sup> **Introduction**



<sup>115</sup> **Chapter 1**

<sup>116</sup> **Neutrino physics**

<sup>117</sup> *The neutrino, or  $\nu$  for the close friends, a fascinating and invisible particle. Some will say that dark matter also have those property but at least we are pretty confident that neutrinos exists.*

<sup>118</sup> **Contents**

|  |              |
|--|--------------|
| <sup>119</sup> <b>1.1 Standard model</b> . . . . .                 | <sup>9</sup> |
| <sup>120</sup> <b>1.1.1 Limits of the standard model</b> . . . . . | <sup>9</sup> |
| <sup>121</sup> <b>1.2 Historic of the neutrino</b> . . . . .       | <sup>9</sup> |
| <sup>122</sup> <b>1.3 Oscillation</b> . . . . .                    | <sup>9</sup> |
| <sup>123</sup> <b>1.3.1 Phenomologies</b> . . . . .                | <sup>9</sup> |
| <sup>124</sup> <b>1.4 Open questions</b> . . . . .                 | <sup>9</sup> |
| <sup>125</sup>   |              |
| <sup>126</sup>   |              |
| <sup>127</sup>   |              |

<sup>129</sup> **1.1 Standard model**

Decrire le m  
Regarder th  
Kochebina  
Limite du r  
Interessant,  
les neutrino  
CP ? Pb des

<sup>130</sup> **1.1.1 Limits of the standard model**

<sup>131</sup> **1.2 Historic of the neutrino**

<sup>132</sup> **First theories**

<sup>133</sup> **Discovery**

<sup>134</sup> **Milestones and anomalies**

<sup>135</sup> **1.3 Oscillation**

<sup>136</sup> **1.3.1 Phenomologies**

<sup>137</sup> **1.4 Open questions**



<sup>138</sup> **Chapter 2**

<sup>139</sup> **The JUNO experiment**

<sup>140</sup>

*"Ave Juno, rosae rosam, et spiritus rex". It means nothing but I found it in tone.*

<sup>141</sup>

## Contents

|                |   |               |
|----------------|---|---------------|
| <sup>142</sup> | <b>2.1 Neutrinos physics in JUNO</b>                                  | <sup>12</sup> |
| <sup>143</sup> | 2.1.1 Reactor neutrino oscillation for NMO and precise measurements   | <sup>12</sup> |
| <sup>144</sup> | 2.1.2 Other physics   | <sup>15</sup> |
| <sup>145</sup> | <b>2.2 The JUNO detector</b>  | <sup>17</sup> |
| <sup>146</sup> | 2.2.1 Detection principle   | <sup>17</sup> |
| <sup>147</sup> | 2.2.2 Central Detector (CD)   | <sup>19</sup> |
| <sup>148</sup> | 2.2.3 Veto detector   | <sup>23</sup> |
| <sup>149</sup> | <b>2.3 Calibration strategy</b>                                       | <sup>23</sup> |
| <sup>150</sup> | 2.3.1 Energy scale calibration  | <sup>24</sup> |
| <sup>151</sup> | 2.3.2 Calibration system  | <sup>25</sup> |
| <sup>152</sup> | 2.3.3 Instrumental non-linearity calibration                          | <sup>25</sup> |
| <sup>153</sup> | <b>2.4 Satellite detectors</b>  | <sup>26</sup> |
| <sup>154</sup> | 2.4.1 TAO   | <sup>26</sup> |
| <sup>155</sup> | 2.4.2 OSIRIS  | <sup>26</sup> |
| <sup>156</sup> | <b>2.5 Software</b>   | <sup>27</sup> |
| <sup>157</sup> | <b>2.6 State of the art of the Offline IBD reconstruction in JUNO</b> | <sup>28</sup> |
| <sup>158</sup> | 2.6.1 Interaction vertex reconstruction                               | <sup>28</sup> |
| <sup>159</sup> | 2.6.2 Energy reconstruction   | <sup>33</sup> |
| <sup>160</sup> | 2.6.3 Machine learning for reconstruction                             | <sup>36</sup> |
| <sup>161</sup> | <b>2.7 JUNO sensitivity to NMO and precise measurements</b>           | <sup>38</sup> |
| <sup>162</sup> | 2.7.1 Theoretical spectrum  | <sup>38</sup> |
| <sup>163</sup> | 2.7.2 Fitting procedure   | <sup>39</sup> |
| <sup>164</sup> | 2.7.3 Physics results   | <sup>39</sup> |
| <sup>165</sup> | <b>2.8 Summary</b>  | <sup>40</sup> |

<sup>166</sup>

<sup>167</sup>

The first idea of a medium baseline ( $\sim 52$  km) experiment, was explored in 2008 [1] where it was demonstrated that the Neutrino Mass Ordering (NMO) could be determined by a medium baseline experiment if  $\sin^2(2\theta_{13}) > 0.005$  without the requirements of accurate knowledge of the reactor antineutrino spectra and the value of  $\Delta m_{32}^2$ . From this idea is born the Jiangmen Underground Neutrino Observatory (JUNO) experiment.

<sup>175</sup>

<sup>176</sup>

JUNO is a neutrino detection experiment under construction located in China, in Guangdong province, near the city of Kaiping. Its main objectives are the determination of the mass ordering at the



FIGURE 2.1 – **On the left:** Location of the JUNO experiment and its reactor sources in southern china. **On the right:** Aerial view of the experimental site

177 3-4 $\sigma$  level in 6 years of data taking and the measurement at the sub-percent precision of the oscillation  
 178 parameters  $\Delta m_{21}^2$ ,  $\sin^2 \theta_{12}$ ,  $\Delta m_{32}^2$  and with less precision  $\sin^2 \theta_{13}$ [2].

179 For this JUNO will measure the electronic anti-neutrinos ( $\bar{\nu}_e$ ) flux coming from the nuclear reactors  
 180 of Taishan, Yangjiang, for a total power of 26.6 GW<sub>th</sub>, and the Daya Bay power plant to a lesser  
 181 extent. All of those cores are the second-generation pressurized water reactors CPR1000, which is a  
 182 derivative of Framatome M310. Details about the power plants characteristics and their expected flux  
 183 of  $\bar{\nu}_e$  can be found in the table 2.1. The distance of 53 km has been specifically chosen to maximize  
 184 the disappearance probability of the  $\bar{\nu}_e$ . The data taking is scheduled to start early 2025.

## 185 2.1 Neutrinos physics in JUNO

186 Even if the JUNO design detailed in section 2.2 was optimized for the measurement of the NMO, its  
 187 large detection volume, excellent energy resolution and background level and understanding make it  
 188 also an excellent detector to measure the flux coming from other neutrino sources. Thus the scientific  
 189 program of JUNO extends way over reactor antineutrinos. The following section is an overview of  
 190 the different physics topic JUNO will contribute in the coming years.

### 191 2.1.1 Reactor neutrino oscillation for NMO and precise measurements

Previous works [1, 3] shows that oscillation parameters and the NMO can be observed by looking at the  $\bar{\nu}_e$  disappearance energy spectrum coming from medium baseline nuclear reactor. This disappearance probability can be expressed as [2] :

$$P(\bar{\nu}_e \rightarrow \bar{\nu}_e) = 1 - \sin^2 2\theta_{12} c_{13}^4 \sin^2 \frac{\Delta m_{21}^2 L}{4E} - \sin^2 2\theta_{13} \left[ c_{12}^2 \sin^2 \frac{\Delta m_{31}^2 L}{4E} + s_{12}^2 \sin^2 \frac{\Delta m_{32}^2 L}{4E} \right]$$

Where  $s_{ij} = \sin \theta_{ij}$ ,  $c_{ij} = \cos \theta_{ij}$ ,  $E$  is the  $\bar{\nu}_e$  energy and  $L$  is the baseline. We can see the sensitivity to the NMO in the dependency to  $\Delta m_{32}^2$  and  $\Delta m_{31}^2$  causing a phase shift of the spectrum as we can see in the figure 2.2. By carefully adjusting a theoretical spectrum to the data, one can extract the NMO and the oscillation parameters. The statistic procedure used to adjust the theoretical spectrum is reviewed in more details in the section 2.7. To reach the desired sensitivity, JUNO must meet multiple requirements but most notably:

1. An energy resolution of  $3\%/\sqrt{E(\text{MeV})}$  to be able to distinguish the fine structure of the fast oscillation.
2. An energy precision of 1% in order to not err on the location of the oscillation pattern.

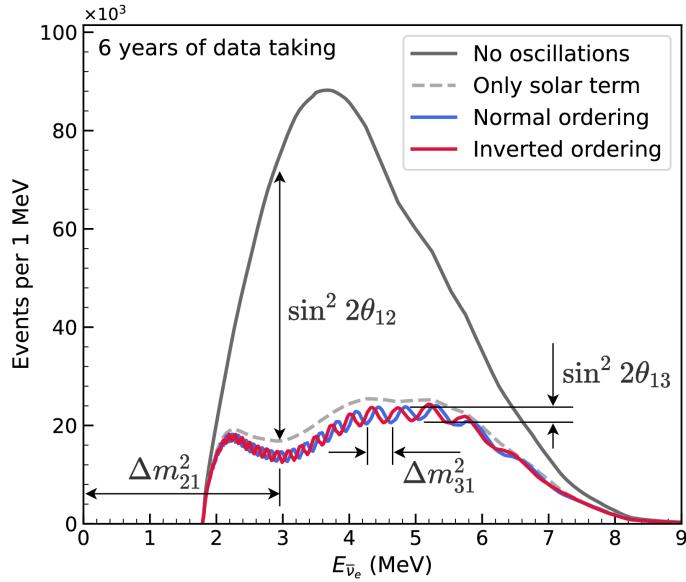


FIGURE 2.2 – Expected number of neutrinos event per MeV in JUNO after 6 years of data taking. The black curve shows the flux if there was no oscillation. The light gray curve shows the oscillation if only the solar terms are taken in account ( $\theta_{12}$ ,  $\Delta m_{21}^2$ ). The blue and red curve shows the spectrum in the case of, respectively, NO and IO. The dependency of the oscillation to the different parameters are schematized by the double sided arrows. We can see the NMO sensitivity by looking at the fine phase shift between the red and the blue curve.

- 201        3. A baseline between 40 and 65 km to maximise the  $\bar{\nu}_e$  oscillation probability. The optimal  
202        baseline would be 58 km and JUNO baseline is 53 km.  
203        4. At least  $\approx 100,000$  events to limit the spectrum distortion due to statistical uncertainties.

204         **$\bar{\nu}_e$  flux coming from nuclear power plants**

205        To get such high measurements precision, it is necessary to have a very good understanding of the  
206        sources characteristics. For its NMO and precise measurement studies, JUNO will observe the energy  
207        spectrum of neutrinos coming from the nuclear power plants Taishan and Yangjiang's cores, located  
208        at 53 km of the detector to maximise the disappearance probability of the  $\bar{\nu}_e$ .

209        The  $\bar{\nu}_e$  coming from reactors are emitted from  $\beta$ -decay of unstable fission fragments. The Taishan  
210        and Yangjiang reactors are Pressurised Water Reactor (PWR), the same type as Daya Bay. In those  
211        type of reactor more than 99.7 % and  $\bar{\nu}_e$  are produced by the fissions of four fuel isotopes  $^{235}\text{U}$ ,  $^{238}\text{U}$ ,  
212         $^{239}\text{Pu}$  and  $^{241}\text{Pu}$ . The neutrino flux per fission of each isotope is determined by the inversion of the  
213        measured  $\beta$  spectra of fission product [4–8] or by calculation using the nuclear databases [9, 10].

214        The neutrino flux coming from a reactor at a time  $t$  can be predicted using

$$\phi(E_\nu, t)_r = \frac{W_{th}(t)}{\sum_i f_i(t)e_i} \sum_i f_i(t) S_i(E_\nu) \quad (2.1)$$

215        where  $W_{th}(t)$  is the thermal power of the reactor,  $f_i(t)$  is the fraction fission of the  $i$ th isotope,  $e_i$  its  
216        thermal energy released in each fission and  $S_i(e_\nu)$  the neutrino flux per fission for this isotope. Using  
217        this method, the flux uncertainty is expected to be of an order of 2-3 % [11].

| Reactor   | Power (GW <sub>th</sub> ) | Baseline (km) |
|-----------|---------------------------|---------------|
| Taishan   | 9.2                       | 52.71         |
| Core 1    | 4.6                       | 52.77         |
| Core 2    | 4.6                       | 52.64         |
| Yangjiang | 17.4                      | 52.46         |
| Core 1    | 2.9                       | 52.74         |
| Core 2    | 2.9                       | 52.82         |
| Core 3    | 2.9                       | 52.41         |
| Core 4    | 2.9                       | 52.49         |
| Core 5    | 2.9                       | 52.11         |
| Core 6    | 2.9                       | 52.19         |
| Daya Bay  | 17.4                      | 215           |
| Huizhou   | 17.4                      | 265           |

TABLE 2.1 – Characteristics of the nuclear power plants observed by JUNO.

218 In addition to those prediction, a satellite experiment named TAO[12] will be setup near the reactor  
 219 core Taishan-1 to measure with an energy resolution of 2% at 1 MeV the neutrino flux coming from  
 220 the core, more details can be found in section 2.4.1. It will help identifying unknown fine structure  
 221 and give more insight on the  $\bar{\nu}_e$  flux coming from this reactor.

222 One the open issue about reactor anti-neutrinos flux is the so-called neutrino anomaly [13], an  
 223 unexpected surplus of neutrino emission in the spectra around 5 MeV. Multiples scientists are trying  
 224 to explain this surplus by advanced recalculation of the nuclei model during beta decay [14, 15] but  
 225 no consensus on this issue has been reached yet.

## 226 Background in the neutrinos reactor spectrum

227 Considering the close reactor neutrinos flux as the main signal, the signals that are considered as  
 228 background are:

- 229 — The geoneutrinos producing background in the 0.511 ~ 2.7 MeV region.
- 230 — The neutrinos coming from the other nuclear reactors around Earth.

231 In addition to all those physics signal, non-neutrinos signal that would mimic an IBD will also be  
 232 present. It is composed of:

- 233 — The signal coming from radioactive decay ( $\alpha$ ,  $\gamma$ ,  $\beta$ ) from natural radioactive isotopes in the  
 material of the detector.
- 235 — Cosmogenic event such as fast neutrons and activated isotopes induced by muons passing  
 through the detector, most notably the spallation on  $^{12}\text{C}$ .

237 All those events represent a non-negligable part of the spectrum as shown in figure 2.3.

## 238 Identification of the mass ordering

239 To identify the mass ordering, we adjust the theoretical neutrino energy spectrum under the two  
 240 hypothesis of NO and IO. Those give us two  $\chi^2$ , respectively  $\chi^2_{NO}$  and  $\chi^2_{IO}$ . By computing the  
 241 difference  $\Delta\chi^2 = \chi^2_{NO} - \chi^2_{IO}$  we can determine the most probable mass ordering and the confidence  
 242 interval: NO if  $\Delta\chi^2 > 0$  and IO if  $\Delta\chi^2 < 0$ . Current studies shows that the expected sensitivity  
 243 the mass ordering would be of  $3.4\sigma$  after 6 years of data taking in nominal setup[2]. More detailed  
 244 explanations about the procedure can be found in the section 2.7.

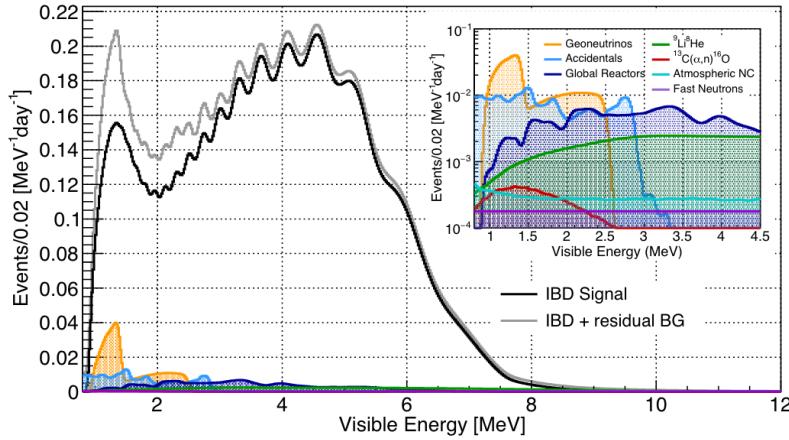


FIGURE 2.3 – Expected visible energy spectrum measured with the LPMT system with (grey) and without (black) backgrounds. The background amount for about 7% of the IBD candidate and are mostly localized below 3 MeV [11]

#### 245 Precise measurement of the oscillations parameters

246 The oscillations parameters  $\theta_{12}$ ,  $\theta_{13}$ ,  $\Delta m_{21}^2$ ,  $\Delta m_{31}^2$  are free parameters in the fit of the oscillation  
 247 spectrum. The precision on those parameters have been estimated and are shown in table 2.2. We  
 248 see that for  $\theta_{12}$ ,  $\Delta m_{21}^2$ ,  $\Delta m_{31}^2$ , precision at 6 years is better than the reference precision by an order of  
 249 magnitude [11]

|   | Central Value | PDG 2020            | 100 days            | 6 years              | 20 years            |
|---|---------------|---------------------|---------------------|----------------------|---------------------|
| $\Delta m_{31}^2 (\times 10^{-3} \text{ eV}^2)$ | 2.5283        | $\pm 0.034$ (1.3%)  | $\pm 0.021$ (0.8%)  | $\pm 0.0047$ (0.2%)  | $\pm 0.0029$ (0.1%) |
| $\Delta m_{21}^2 (\times 10^{-3} \text{ eV}^2)$ | 7.53          | $\pm 0.18$ (2.4%)   | $\pm 0.074$ (1.0%)  | $\pm 0.024$ (0.3%)   | $\pm 0.017$ (0.2%)  |
| $\sin^2 \theta_{12}$                            | 0.307         | $\pm 0.013$ (4.2%)  | $\pm 0.0058$ (1.9%) | $\pm 0.0016$ (0.5%)  | $\pm 0.0010$ (0.3%) |
| $\sin^2 \theta_{13}$                            | 0.0218        | $\pm 0.0007$ (3.2%) | $\pm 0.010$ (47.9%) | $\pm 0.0026$ (12.1%) | $\pm 0.0016$ (7.3%) |

TABLE 2.2 – A summary of precision levels fir the oscillation parameters. The reference value (PDG 2020 [16]) is compared with 100 days, 6 years and 20 years of JUNO data taking.

#### 250 2.1.2 Other physics

251 While the design of JUNO is tailored to measure  $\bar{\nu}_e$  coming from nuclear reactor, JUNO will be able  
 252 to detect neutrinos coming from other sources thus allowing for a wide range of physics studies as  
 253 detailed in the table 2.3 and in the following sub-sections.

#### 254 Geoneutrinos

255 Geoneutrinos designate the antineutrinos coming from the decay of long-lived radioactive elements  
 256 inside the Earth. The 1.8 MeV threshold necessary for the IBD makes it possible to measure geoneu-  
 257 trinos from  $^{238}\text{U}$  and  $^{232}\text{Th}$  decay chains. The studies of geoneutrinos can help refine the Earth  
 258 crust models but is also necessary to characterise their signal, as they are a background to the mass  
 259 ordering and oscillations parameters studies.

| Research             | Expected signal                      | Energy region | Major backgrounds          |
|----------------------|--------------------------------------|---------------|----------------------------|
| Reactor antineutrino | 60 IBDs/day                          | 0–12 MeV      | Radioactivity, cosmic muon |
| Supernova burst      | 5000 IBDs at 10 kpc                  | 0–80 MeV      | Negligible                 |
| DSNB (w/o PSD)       | 2300 elastic scattering              |               |                            |
| Solar neutrino       | 2–4 IBDs/year                        | 10–40 MeV     | Atmospheric $\nu$          |
| Atmospheric neutrino | hundreds per year for ${}^8\text{B}$ | 0–16 MeV      | Radioactivity              |
| Geoneutrino          | hundreds per year                    | 0.1–100 GeV   | Negligible                 |
|                      | $\approx 400$ per year               | 0–3 MeV       | Reactor $\nu$              |

TABLE 2.3 – Detectable neutrino signal in JUNO and the expected signal rates and major background sources

## 260 Atmospheric neutrinos

261 Atmospheric neutrinos are neutrinos originating from the decay of  $\pi$  and  $K$  particles that are pro-  
 262 duced in extensive air showers initiated by the interactions of cosmic rays with the Earth atmosphere.  
 263 Earth is mostly transparent to neutrinos below the PeV energy, thus JUNO will be able to see neu-  
 264 trinos coming from all directions. Their baseline range is large (15km  $\sim$  13000km), they can have  
 265 energy between 0.1 GeV and 10 TeV and will contain all neutrino and antineutrinos flavour. Their  
 266 studies is complementary to the reactor antineutrinos and can help refine the constraints on the NMO  
 267 [2].

## 268 Supernovae burst neutrinos

269 Neutrinos are crucial component during all stages of stellar collapse and explosion. Detection of  
 270 neutrinos coming for core collapse supernovae will provide us important informations on the mech-  
 271 anisms at play in those events. Thanks to its 20 kt sensible volume, JUNO has excellent capabilities  
 272 to detect all flavour of the  $\mathcal{O}(10 \text{ MeV})$  postshock neutrinos, and using neutrinos of the  $\mathcal{O}(1 \text{ MeV})$   
 273 will give informations about the pre-supernovae neutrinos. All those informations will allow to  
 274 disentangle between the multiple hydro-dynamic models that are currently used to describe the  
 275 different stage of core-collapse supernovae.

## 276 Diffuse supernovae neutrinos background

277 Core-collapse supernovae in our galaxy are rare events, but they frequently occur throughout the  
 278 visible Universe sending burst of neutrinos in direction of the Earth. All those events contributes to  
 279 a low background flux of low-energy neutrinos called the Diffuse Supernovae Neutrino Background  
 280 (DSNB). Its flux and spectrum contains informations about the red-shift dependent supernovae rate,  
 281 the average supernovae neutrino energy and the fraction of black-hole formation in core-collapse su-  
 282 pernovae. Depending of the DSNB model, we can expect 2-4 IBD events per year in the energy range  
 283 above the reactor  $\bar{\nu}_e$  signal, which is competitive with the current Super-Kamiokande+Gadolinium  
 284 phase [17].

## 285 Beyond standard model neutrinos interactions

286 JUNO will also be able to probe for beyond standard model neutrinos interactions. After the main  
 287 physics topics have been accomplished, JUNO could be upgraded to probe for neutrinoless beta  
 288 decay ( $0\nu\beta\beta$ ). The detection of such event would give critical informations about the nature of  
 289 neutrinos, is it a majorana or a dirac particle. JUNO will also be able to probe for neutrinos that  
 290 would come for the decay or annihilation of Dark Matter inside the sun and neutrinos from putative

291 primordial black hole. Through the unitary test of the mixing matrix, JUNO will be able to search for  
 292 light sterile neutrinos. Thanks to JUNO sensitivity, multiple other exotic research can be performed  
 293 on neutrino related beyond standard model interactions.

294 **Proton decay**

295 Proton decay is a potential unobserved event where the proton decay by violating the baryon number.  
 296 This violation is necessary to explain the baryon asymmetry in the universe and is predicted  
 297 by multiple Grand Unified Theories which unify the strong, weak and electromagnetic interactions.  
 298 Thanks to its large active volume, JUNO will be able to take measurement of the potential proton  
 299 decay channel  $p \rightarrow \bar{\nu}K^+$  [18] thanks to the timing resolution of the SPMT system. Studies show  
 300 that JUNO should be competitive with the current best limit at  $5.9 \times 10^{33}$  years from Super-K. This  
 301 studies show that JUNO, considering no proton decay events observed, would be able to rule a  
 302 limit of  $9.6 \times 10^{33}$  years at 90 % C.L.

303 **2.2 The JUNO detector**

304 The JUNO detector is a scintillator detector buried 693.35 meters under the ground (1800 meters  
 305 water equivalent). It consists of Central Detector (CD), a water pool and a Top Tracker (TT) as shown  
 306 in figure 2.4a. The CD is an acrylic vessel containing the 20 ktons of Liquid Scintillator (LS). It is  
 307 supported by a stainless steel structure and is immersed in that water pool that is used as shielding  
 308 from external radiation and as a cherenkov detector for the background. The top of the experiment  
 309 is partially covered by the Top Tracker (TT), a plastic scintillator detector which is used to detect the  
 310 atmospheric muons background and is acting as a veto detector.

311 The top of the experiment also host the LS purification system, a water purification system, a ventilation  
 312 system to get rid of the potential radon in the air. The CD is observed by two systems of  
 313 Photo-Multiplier Tubes (PMT). They are attached to the steel structure and their electronic readout  
 314 is submerged near them. A third system of PMT is also installed on the structure but are facing  
 315 outward of the CD, instrumenting the water to be cherenkov detector. The CD and the cherenkov  
 316 detector are optically separated by Tyvek sheet. A chimney for LS filling and purification and for  
 317 calibration operations connects the CD to the experimental hall from the top.

318 The CD has been dimensioned to meet the requirements presented in section 2.1.1:

- 319 — Its 20 ktons monolithic LS provide a volume sizeable enough, in combination with the ex-  
 320 pected  $\bar{\nu}_e$  flux, to reach the desired statistic in 6 years. Its monolithic nature also allows for a  
 321 full containment of most of the events, preventing the energy loss in non-instrumented parts  
 322 that would arise from a segmented detector.
- 323 — Its large overburden shield it from most of the atmospheric background that would pollute  
 324 the signal.
- 325 — The localization of the experiment, chosen to maximize the disappearance with a 53km base-  
 326 line and in a region that allows two nuclear power plants to be used as sources.

327 This section covers in details the different components of the detector and the detection systems.

328 **2.2.1 Detection principle**

The CD will detect the neutrino and measure their energy mainly via an Inverse Beta Decay (IBD) interaction with proton mainly from the  $^{12}\text{C}$  and H nucleus in the LS:

$$\bar{\nu}_e + p \rightarrow n + e^+$$

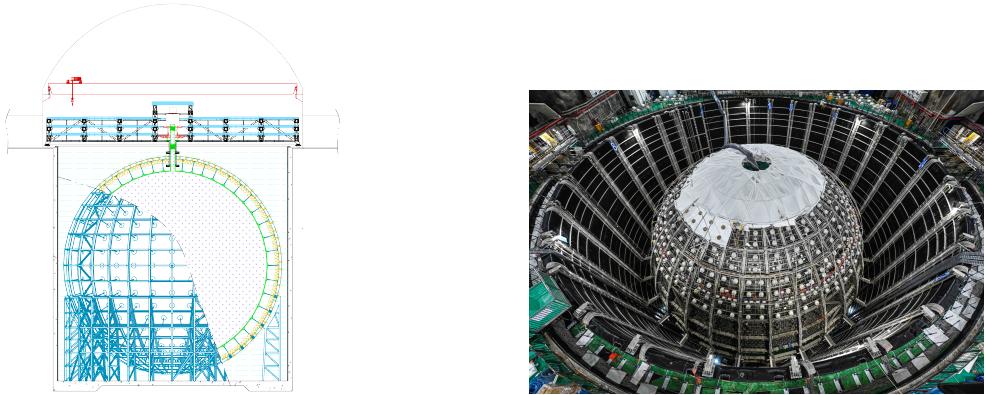


FIGURE 2.4

329 Kinematics calculation shows that this interaction has an energy threshold for the  $\bar{\nu}_e$  of  $(m_n + m_e -$   
 330  $m_p) \approx 1.806$  MeV [19]. This threshold make the experiment blind to very low energy neutrinos.  
 331 The residual energy  $E_\nu - 1.806$  MeV is be distributed as kinetic energy between the positron and the  
 332 neutron. The energy of the emitted positron  $E_e$  is given by [19]

$$E_e = \frac{(E_\nu - \delta)(1 + \epsilon_\nu) + \epsilon_\nu \cos \theta \sqrt{(E_\nu - \delta)^2 + \kappa m_e^2}}{\kappa} \quad (2.2)$$

333 where  $\kappa = (1 + \epsilon_\nu)^2 - \epsilon_\nu^2 \cos^2 \theta \approx 1$ ,  $\epsilon_\nu = \frac{E_\nu}{m_p} \ll 1$  and  $\delta = \frac{m_n^2 - m_p^2 - m_e^2}{2m_p} \ll 1$ . We can see from this  
 334 equation that the positron energy is strongly correlated to the neutrino energy.

335 The positron and the neutron will then propagate in the detection medium, the Liquid Scintillator  
 336 (LS), loosing their kinetic energy by exciting the molecule of the LS (more details in section 2.2.2).  
 337 Once stopped, the positron will annihilate with an electron from the medium producing two 511  
 338 KeV gamma. Those gamma will themselves interact with the LS, exciting it before being absorbed  
 339 by photoelectrical effect. The neutron will be captured by an hydrogen, emitting a 2.2 MeV gamma  
 340 in the process. This gamma will also deposit its energy before being absorbed by the LS.

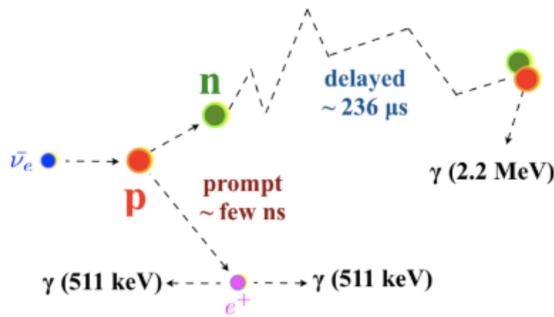


FIGURE 2.5 – Schematics of an IBD interaction in the central detector of JUNO

341 The scintillation photons have frequency in the UV and will propagate in the LS, being re-absorbed  
 342 and re-emitted by compton effect before finally be captured by PMTs instrumenting the acrylic  
 343 sphere. The analog signal of the PMTs digitized by the electronic is the signal of our experiment.

344 The signal produced by the positron is subsequently called the prompt signal, and the signal coming  
 345 from the neutron the delayed signal. This naming convention come from the fact that the positron  
 346 will deposit its energy rather quickly (few ns) where the neutron will take a bit more time ( $\sim 236 \mu\text{s}$ ).

### 347 2.2.2 Central Detector (CD)

348 The central detector, composed of 20 ktons of Liquid Scintillator (LS), is the main part of JUNO. The  
 349 LS is contained in a spherical acrylic vessel supported by a stainless steel structure. The CD and  
 350 its structural support are submerged in a cylindrical water pool of 43.5m diameter and 44m height.  
 351 We're confident that the water pool provide sufficient buffer protection in every direction against the  
 352 rock radioactivity.

#### 353 Acrylic vessel

354 The acrylic vessel is a spherical vessel of inner diameter of 35.4 m and a thickness of 120 mm. It is  
 355 assembled from 265 acrylic panels, thermo bonded together. The acrylic recipes has been carefully  
 356 tuned with extensive R&D to ensure it does not include plasticizer and anti-UV material that would  
 357 stop the scintillation photons. Those panels requires to be pure of radioactive materials to not  
 358 cause background. Current setup where the acrylic panels are molded in cleanrooms of class 10000,  
 359 let us reach a uranium and thorium contamination of <0.5 ppt. The molding and thermoforming  
 360 processes is optimized to increase the assemblage transparency in water to >96%. The acrylic vessel  
 361 is supported by a stainless steel structure via supporting node (fig 2.6). The structure and the nodes  
 362 are designed to be resilient to natural catastrophic events such as earthquake and can support many  
 363 times the effective load of the acrylic vessel.

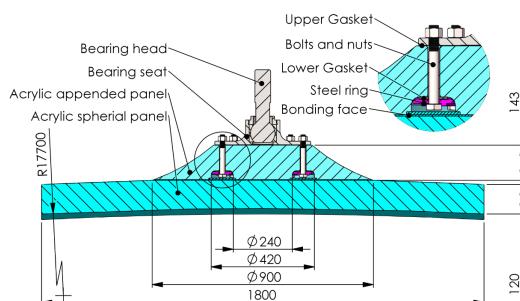


FIGURE 2.6 – Schematics of the supporting node for the acrylic vessel

#### 364 Liquid scintillator

365 The Liquid Scintillator (LS) has a similar recipe as the one used in Daya Bay [20] but without gadolinium  
 366 doping. It is made of three components, necessary to shift the wavelength of emitted photons to  
 367 prevent their reabsorption and to shift their wavelength to the PMT sensitivity region as illustrated  
 368 in figure 2.7:

- 369 1. The detection medium, the *linear alkylbenzene* (LAB). Selected because of its excellent trans-  
 370 parency, high flash point, low chemical reactivity and good light yield. Accounting for  $\sim 98\%$  of the LS, it is the main component with which ionizing particles and gamma interact.  
 371 Charged particles will collide with its electronic cloud transferring energy to the molecules,  
 372 gamma will interact via compton effect with the electronic cloud before finally be absorbed  
 373 via photoelectric effect.

- 375    2. The second component of the LS is the *2,5-diphenyloxazole* (PPO). A fraction of the excitation  
 376    energy of the LAB is transferred to the PPO, mainly via non radiative process [21]. The  
 377    PPO molecules de-excites in the same way, transferring their energy to the bis-MSB. The PPO  
 378    makes for 1.5 % of the LS.
- 379    3. The last component is the *p-bis(o-methylstyryl)-benzene* (bis-MSB). Once excited by the PPO, it  
 380    will emit photon with an average wavelength of  $\sim 430$  nm (full spectrum in figure 2.7) that  
 381    can thus be detected by our photo-multipliers systems. It amount for  $\sim 0.5\%$  of the LS.

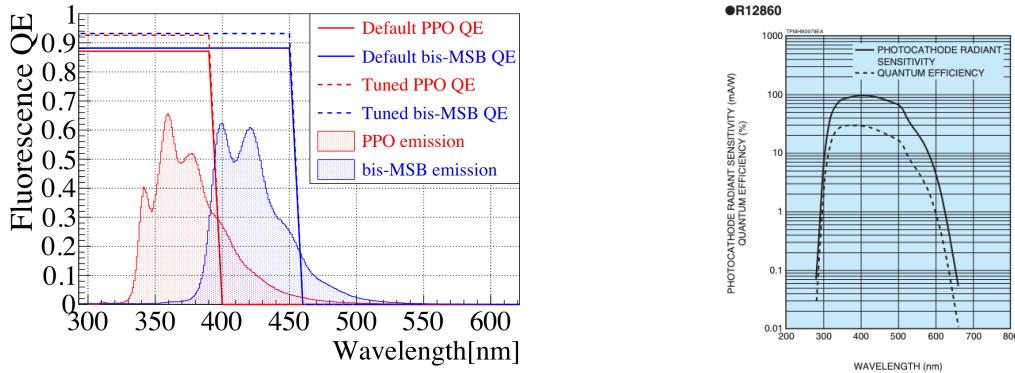


FIGURE 2.7 – On the left: Quantum efficiency (QE) and emission spectrum of the LAB and the bis-MSB [20]. On the right: Sensitivity of the Hamamatsu LPMT depending on the wavelength of the incident photons [22].

This formula has been optimized using dedicated studies with a Daya Bay detector [20, 23] to reach the requirements for the JUNO experiment:

- A light yield / MeV of the amount of  $10^4$  photons to maximize the statistic in the energy measurement.
- An attenuation length comparable to the size of the detector to prevent losing photons during their propagation in the LS. The final attenuation length is 25.8m [24] to compare with the CD diameter of 35.4m.
- Uranium/Thorium radiopurity to prevent background signal. The reactor neutrino program require a contamination fraction  $F < 10^{-15}$  while the solar neutrino program require  $F < 10^{-17}$ .

The LS will frequently be purified and tested in the Online Scintillator Internal Radioactivity Investigation System (OSIRIS) [25] to ensure that the requirements are kept during the lifetime of the experiment, more details to be found in section 2.4.2.

### Large Photo-Multipliers Tubes (LPMTs)

The scintillation light produced by the LS is then collected by Photo-Multipliers Tubes (PMT) that transform the incoming photon into an electric signal. As described in figure 2.8, the incident photons interact with the photocathode via photoelectric effect producing an electron called a Photo-Electron (PE). This PE is then focused on the dynodes where the high voltage will allow it to be multiplied. After multiple amplification the resulting charge - in coulomb [C] - is collected by the anode and the resulting electric signal can be digitalized by the readout electronics from which the charge and timing can be extracted.

The Large Photo-Multipliers Tubes (LPMT), used in the central detector and in the water pool, are 20-inch (50.8 cm) radius PMTs.  $\sim 5000$  dynode-PMTs [22] were produced by the Hamamatsu<sup>®</sup> company and  $\sim 15000$  Micro-Channel Plate (MCP) [26] by the NNVT<sup>®</sup> company. This system is the one responsible for the energy measurement with a energy resolution of  $3\%/\sqrt{E}$ , resolution

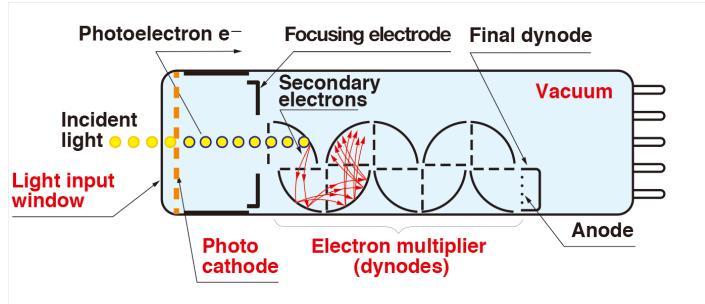


FIGURE 2.8 – Schematic of a PMT

407 necessary for the mass ordering measurement. To reach this precision, the system is composed of  
 408 17612 PMTs quasi uniformly distributed over the detector for a coverage of 75.2% reaching  $\sim 1800$   
 409 PE/MeV or  $\sim 2.3\%$  resolution due to statistic, leaving  $\sim 0.7\%$  for the systematic uncertainties. They  
 410 are located outside the acrylic sphere in the water pool facing the center of the detector. To maintain  
 411 the resolution over the lifetime of the experiment, JUNO require a failure rate  $< 1\%$  over 6 years.

412 The LPMTs electronic are divided in two parts. One "near", located underwater, in proximity of the  
 413 LPMT to reduce the cable length between the PMT and early electronic. A second one, outside of the  
 414 detector that is responsible for higher level analysis before sending the data to the DAQ.

415 The light yield per MeV induce that a LPMT can collect between 1 and 1000 PE per event, a wide  
 416 dynamic range, causing non linearity in the PMT response that need to be understood and calibrated,  
 417 see section 2.3 for more details.

418 Before performing analysis, the analog readout of the LPMT need to be amplified, digitised and  
 419 packaged by the readout electronics schematized in figure 2.9. This electronic is splitted in two parts:  
 420 *wet* electronic that are located near the LPMTs, protected in an Underwater Box (UWB) and the *dry*  
 421 electronics located in deicated rooms outside of the water pool.

422 The LPMTs are connected to the UWB by groups of three. Each UWB contains:

- 423 — Three high voltage units, each one powering a PMT.
- 424 — A global control unit, responsible for the digitization of the waveform, composed of six analog-digital units that produce digitized waveform and a Field Programmable Gate Array (FPGA) that complete the waveform with metadatas such as the local timestamp trigger, etc... Ths FPGA also act as a data buffer when needed by the DAQ and trigger system.
- 425 — Additional memory in order to temporally store the data in case of sudden burst of the input rate (such as in the case of nearby supernovae).

430 The *dry* electronic synchronize the signals from the UWBS abd centralise the information of the CD  
 431 LPMTs. It act as the Global Trigger by sending the UWB data to DAQ in the case if the LPMT  
 432 multiplicity condition is fulfilled.

### 433 Small Photo-Multipliers Tubes (SPMTs)

434 The Small PMT (SPMTs) system is made of 3-inch (7.62 cm) PMTs. They will be used in the CD  
 435 as a secondary detection system. Those 25600 SPMTs will observe the same events as the LPMTs,  
 436 thus sharing the physics and detector systematics up until the photon conversion. With a detector  
 437 coverage of 2.7%, this system will collect  $\sim 43$  PE/MeV for a final energy resolution of  $\sim 17\%$ .  
 438 This resolution is not enough to measure the NMO,  $\theta_{13}$ ,  $\Delta m^2_{31}$  but will be sufficient to independently  
 439 measure  $\theta_{12}$  and  $\Delta m^2_{21}$ .

440 The benefit of this second system is to be able to perform another, independent measure of the  
 441 same events as the LPMTs, constituting the Dual Calorimetry useful for calibrationa and, as it we

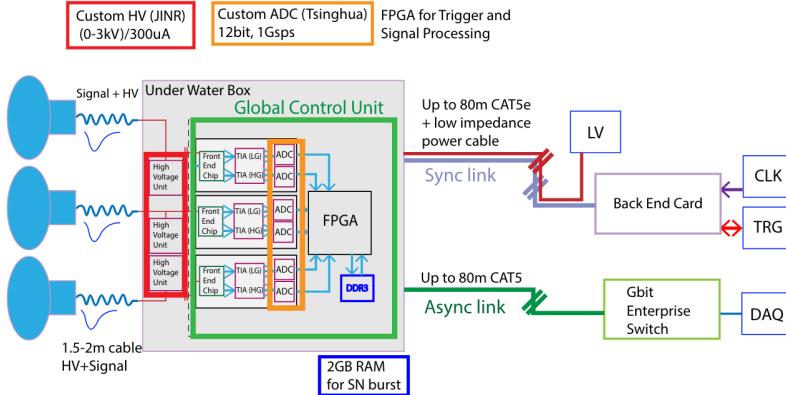


FIGURE 2.9 – The LPMT electronics scheme. It is composed of two part, the *wet* electronics on the left, located underwater and the *dry* electronics on the right. They are connected by Ethernet cable for data transmission and a dedicated low impedance cable for power distribution

will explore in this thesis, for physics analysis. Due to the low PE rate, SPMTs will be running in photo-counting mode in the reactor range and thus will be insensitive to LPMT intrinsic effect (see section 2.3). Using this property, the intrinsic charge non linearity of the LPMTs can be measured by comparing the PE count in the SPMTs and LPMTs [27]. Also, due to their smaller size and electronics, SPMTs have a better timing resolutions than the LPMTs. At higher energy range, like supernovae events, LPMTs will saturate where SPMTs due to their lower PE collection will to produce a reliable measure of the energy spectrum.

The SPMTs will be grouped by pack of 128 to an UWB hosting their electronics as illustrated in figure 2.10. This underwater box host two high voltage splitter boards, each one supplying 64 SPMTs, an ASIC Battery Card (ABC) and a global control unit.

The ABC board will readout and digitize the charge and time of the 128 SPMTs signals and a FPGA will joint the different metadata. The global control unit will handle the powering and control of the board and will be in charge of the transmission of the data to the DAQ.

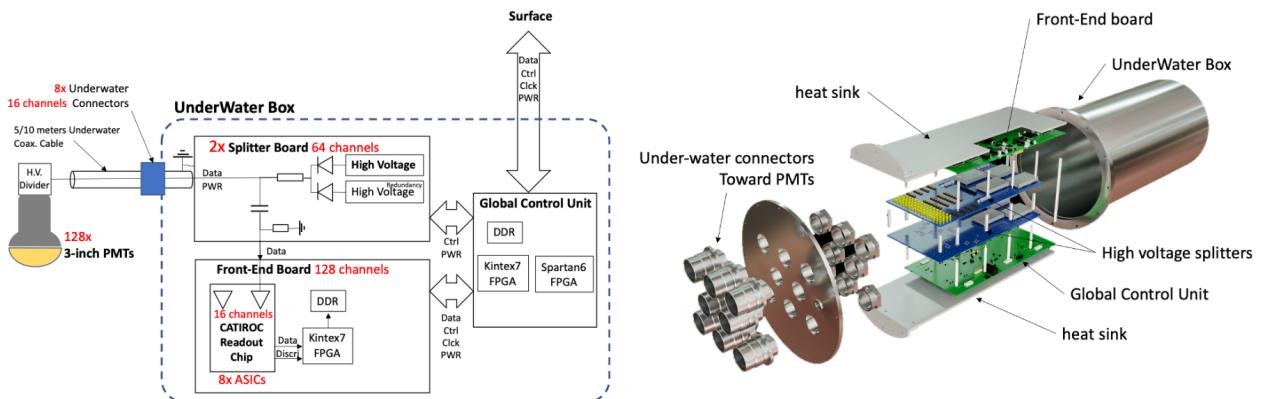


FIGURE 2.10 – Schematic of the JUNO SPMT electronic system (left), and exploded view of the main component of the UWB (right)

### 455 2.2.3 Veto detector

456 The CD will be bathed in constant background noise coming from numerous sources : the radioac-  
 457 tivity from surrounding rock and its own components or from the flux of cosmic muons. This  
 458 background needs to be rejected to ensure the purity of the IBD spectrum. To prevent a big part  
 459 of them, JUNO use two veto detector that will tag events as background before CD analysis.

460 **Cherenkov in water pool**

461 The Water Cherenkov Detector (WCD) is the instrumentation of the water buffer around the CD.  
 462 When high speed charged particles will pass through the water, they will produced cherenkov  
 463 photons. The light will be collected by 2400 MCP LPMTs installed on the outer surface of the CD  
 464 structure. The muons veto strategy is based on a PMT multiplicity condition. WCD PMTs are  
 465 grouped in ten zones: 5 in the top, 5 in the bottom. A veto is raised either when more than 19  
 466 PMTs are triggered in one zone or when two adjacent zones simultaneously trigger more than 13  
 467 PMTs. Using this trigger, we expect to reach a muon detection efficiency of 99.5% while keeping the  
 468 noise at reasonable level.

469 **Top tracker**

470 The JUNO Top Tracker (TT) is a plastic scintillator detector located on the top of the experiment (see  
 471 figure 2.11). Made from plastic scintillator from OPERA [28] layered horizontally in 3 layers on the  
 472 top of the detector, the TT will be able to detect incoming atmospheric muons. With its coverage,  
 473 about 1/3 of the of all atmospheric muons that passing through the CD will also pass through the 3  
 474 layer of the detector. While it does not cover the majority of the CD, the TT is particularly effective  
 475 to detect muons coming through the filling chimney region which might present difficulties from the  
 other subsystems in some classes of events.

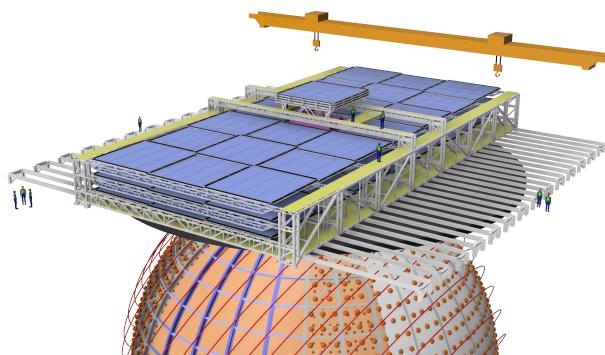


FIGURE 2.11 – The JUNO top tracker

476

## 477 2.3 Calibration strategy

478 The calibration is a crucial part of the JUNO experiment. The detector will continuously bath in  
 479 neutrinos coming from the close nuclear power plant, from other sources such as geo neutrinos,  
 480 the sun and will be exposed to background noise coming from atmospheric muons and natural  
 481 radioactivity. Because of this continuous rate, low frequency signal event, we need high frequency,

recognisable sources in the energy range of interest : [0-12] MeV for the positron signal and 2.2 MeV for the neutron capture. It is expected that the CD response will be different depending on the type of particle, due to the interaction with LS, the position on the event and the optical response of the acrylic sphere (see section 2.6). We also expect a non-linear energy response of the CD due to the LS properties [20] but also due to the saturation of the LPMTs system when collecting a large amount of PE [27].

### 2.3.1 Energy scale calibration

While electrons and positrons sources would be ideal, for a large LS detector thin-walled electrons or positrons sources could lead to leakage of radionuclides causing radioactive contamination. Instead, we consider gamma sources in the range of the prompt energy of IBDs. The sources are reported in table 2.4.

| Sources / Processes             | Type        | Radiation                                |
|---------------------------------|-------------|--|
| $^{137}\text{Cs}$               | $\gamma$    | 0.0662 MeV                               |
| $^{54}\text{Mn}$                | $\gamma$    | 0.835 MeV                                |
| $^{60}\text{Co}$                | $\gamma$    | 1.173 + 1.333 MeV                        |
| $^{40}\text{K}$                 | $\gamma$    | 1.461 MeV                                |
| $^{68}\text{Ge}$                | $e^+$       | annihilation 0.511 + 0.511 MeV           |
| $^{241}\text{Am-Be}$            | $n, \gamma$ | neutron + 4.43 MeV ( $^{12}\text{C}^*$ ) |
| $^{241}\text{Am-}^{13}\text{C}$ | $n, \gamma$ | neutron + 6.13 MeV ( $^{16}\text{O}^*$ ) |
| $(n, \gamma)p$                  | $\gamma$    | 2.22 MeV                                 |
| $(n, \gamma)^{12}\text{C}$      | $\gamma$    | 4.94 MeV or 3.68 + 1.26 MeV              |

TABLE 2.4 – List of sources and their process considered for the energy scale calibration

For the  $^{68}\text{Ge}$  source, it will decay in  $^{68}\text{Ga}$  via electron capture, which will itself  $\beta^+$  decay into  $^{68}\text{Zn}$ . The positrons will be absorbed by the enclosure so only the annihilation gamma will be released. In addition,  $(\alpha, n)$  sources like  $^{241}\text{Am-Be}$  and  $^{241}\text{Am-}^{13}\text{C}$  are used to provide both high energy gamma and neutrons, which will later be captured in the LS producing the 2.2 MeV gamma.

From this calibration we call  $E_{vis}$  the "visible energy" that is reconstructed by our current algorithms and we compare it to the true energy deposited by the calibration source. The results shown in figure 2.12 show the expected response of the detector from calibration sources. The non-linearity is clearly visible from the  $E_{vis}/E_{true}$  shape. See [29] for more details.

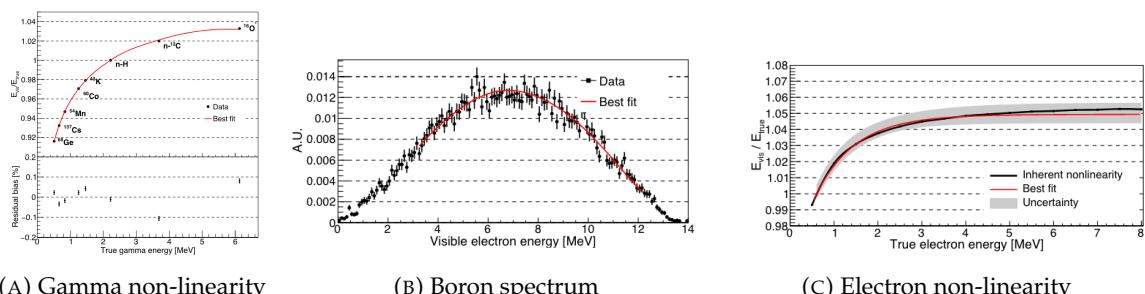


FIGURE 2.12 – Fitted and simulated non linearity of gamma, electron sources and from the  $^{12}\text{B}$  spectrum. Black points are simulated data. Red curves are the best fits. Figures taken from [29].

501 **2.3.2 Calibration system**

502 The non-uniformity due to the event position in the detector (more details in section 2.6) will be  
 503 studied using multiples systems that are schematized in figure 2.13. They allow to position sources  
 504 at different location in the CD.

- 505 — For a one-dimension vertical calibration, the Automatic Calibration Unit (ACU) will be able  
 506 to deploy multiple radioactive sources or a pulse laser diffuser ball along the central axis of  
 507 the CD through the top chimney. The source position precision is less than 1cm.
- 508 — For off-axis calibration, a calibration source attached to a Cable Loop System (CLS) can be  
 509 moved on a vertical half-plane by adjusting the length of two connection cable. Two set of  
 510 CSL will be deployed to provide a 79% effective coverage of a vertical plane.
- 511 — A Guiding Tube (GT) will surround the CD to calibrate the non-uniformity of the response at  
 512 the edge of the detector
- 513 — A Remotely Operated under-LS Vehicle (ROV) can be deployed to desired location inside LS  
 514 for a more precise and comprehensive calibration. The ROV will also be equipped with a  
 515 camera for inspection of the CD.

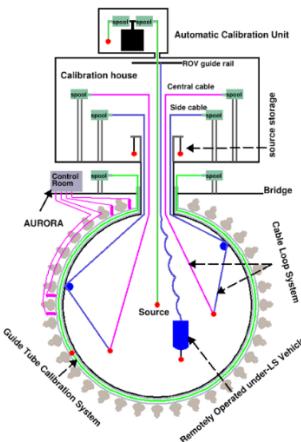


FIGURE 2.13 – Overview of the calibration system

516 The preliminary calibration program is depicted in table 2.5.

517 **2.3.3 Instrumental non-linearity calibration**

518 As mentioned in the introduction of this section, we expect an instrumental non-linearity due to the  
 519 LPMT system saturating. This results in the LPMT underestimating the number of collected photo-  
 520 electrons. This non-linearity is illustrated in figure 2.14. This non-linearity would consequently  
 521 convolve with the LS non-linearity. To correct this effect, the LPMT are first calibrated to the channel  
 522 level using the dual calorimetry calibration technique which consist of comparing the LPMT and  
 523 SPMT calorimetry calibration using a tunable light source covering the range of 0 to 100 PE per  
 524 LPMT channel.

525 Within such range, the SPMT serve as an approximate linear reference since SPMT operate primarily  
 526 operate in photo-counting mode in this range. Using this technique, the residual non-linearity in the  
 527 LPMT response due to the saturation effect is under 0.3 %.

| Program                   | Purpose           | System          | Duration [min] |
|---------------------------|-------------------|-----------------|----------------|
| Weekly calibration        | Neutron (Am-C)    | ACU             | 63             |
|                           | Laser             | ACU             | 78             |
| Monthly calibration       | Neutron (Am-C)    | ACU             | 120            |
|                           | Laser             | ACU             | 147            |
|                           | Neutron (Am-C)    | CLS             | 333            |
|                           | Neutron (Am-C)    | GT              | 73             |
| Comprehensive calibration | Neutron (Am-C)    | ACU, CLS and GT | 1942           |
|                           | Neutron (Am-Be)   | ACU             | 75             |
|                           | Laser             | ACU             | 391            |
|                           | $^{68}\text{Ge}$  | ACU             | 75             |
|                           | $^{137}\text{Cs}$ | ACU             | 75             |
|                           | $^{54}\text{Mn}$  | ACU             | 75             |
|                           | $^{60}\text{Co}$  | ACU             | 75             |
|                           | $^{40}\text{K}$   | ACU             | 158            |

TABLE 2.5 – Calibration program of the JUNO experiment

## 528 2.4 Satellite detectors

529 As introduced in section 2.1.1 and section 2.2.2, the precise knowledge and understanding of the  
 530 detector condition is crucial for the measurements of the NMO and oscillation parameters. Thus two  
 531 satellite detectors will be setup to monitor the experiment condition. TAO to monitor and understand  
 532 the  $\bar{\nu}_e$  flux and spectrum coming from the nuclear reactor and OSIRIS to monitor the LS response.

### 533 2.4.1 TAO

534 The Taishan Antineutrino Observatory (TAO) [12, 30] is a ton-level gadolinium doped liquid scin-  
 535 tillator detector that will be located near the Taishan-1 reactor. It aim to measure the  $\bar{\nu}_e$  spectrum at  
 536 very low distance (44m) from the reactor to measure a quasi-unoscillated spectrum. TAO also aim to  
 537 provide a major contribution to the so-called reactor anomaly [13]. Its requirement are to the level of  
 538 2 % energy resolution at 1 MeV.

#### 539 Detector

540 The TAO detector is close, in concept, to the CD of JUNO. It is composed of an acrylic vessel  
 541 containing 2.8 tons of gadolinium-loaded LS instrumented by an array of silicon photomultipliers  
 542 (SiPM) reaching a 95% coverage. To efficiently reduce the dark count of those sensors, the detector  
 543 is cooled to -50 °C. The  $\bar{\nu}_e$  will interact with the LS via IBD, producing scintillation light, that will  
 544 be detected by the SiPMs. From this signal the  $\bar{\nu}_e$  energy and the full spectrum reconstructed. This  
 545 spectrum will then be used by JUNO to calibrate the unoscillated spectrum, most notably the fission  
 546 product fraction that impact the rate and shape of the spectrum. A schema of the detector is presented  
 547 in figure 2.15a.

### 548 2.4.2 OSIRIS

549 The Online Scintillator Internal Radioactivity Investigation System (OSIRIS) [25] is an ultralow back-  
 550 ground, 20 m<sup>3</sup> LS detector that will be located in JUNO cavern. It aim to monitor the radioactive  
 551 contamination, purity and overall response of the LS before it is injected in JUNO. OSIRIS will

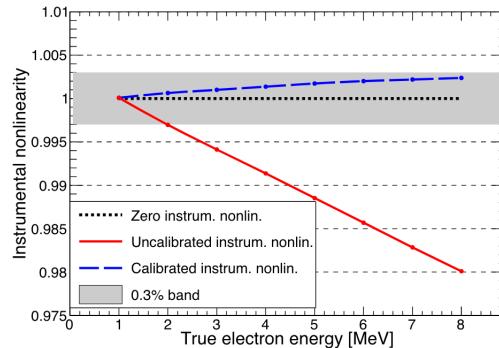


FIGURE 2.14 – Event-level instrumental non-linearity, defined as the ratio of the total measured LPMT charge to the true charge for events uniformly distributed in the detector. The solid red line represents event-level non-linearity without the channel-level correction, with position non-uniformity obtained at 1 MeV applied, in an extreme hypothetical scenario of 50% non-linearity over 100 PEs for the LPMTs. The dashed blue line represents that after the channel-level correction. The gray band shows the residual uncertainty of 0.3%, after the channel-level correction. Figure taken from [29].

552 be located at the end of the purification chain of JUNO, monitoring that the purified LS meet the  
 553 JUNO requirements. The setup is optimized to detect the fast coincidences decay of  $^{214}\text{Bi} - ^{214}\text{Po}$   
 554 and  $^{212}\text{Bi} - ^{212}\text{Po}$ , indicators of the decay chains of U and Th respectively.

555 **Detector**

556 OSIRIS is composed of an acrylic vessel that will contain 17t of LS. The LS is instrumented by  
 557 a PMT array of 64 20 inch PMTs on the top and the side of the vessel. To reach the necessary  
 558 background level required by the LS purity measurements, in addition to being 700m underground  
 559 in the experiment cavern, the acrylic vessel is immersed in a tank of ultra pure water. The water is  
 560 itself instrumented by another array of 20 inch PMTs, acting as muon veto. A schema of the detector  
 561 is presented in figure 2.15b.

562 **2.5 Software**

563 The simulation, reconstruction and analysis algorithms are all packaged in the JUNO software,  
 564 subsequently called the software. It is composed of multiple components integrated in the SNiPER  
 565 [31] framework:

- 566 — Various primary particles simulators for the different kind of events, background and calibra-  
 567 tion sources.
- 568 — A Geant4 [32–34] Monte Carlo (MC) simulation containing the detectors geometries, a custom  
 569 optical model for the LS and the supporting structures of the detectors. The Geant4 simulation  
 570 integrate all relevant physics process for JUNO, validated by the collaboration. This step of the  
 571 simulation is commonly called *Detsim* and compute up to the production of photo-electrons  
 572 in the PMTs. The optics properties of the different materials and detector components have  
 573 been measured beforehand to be used to define the material and surfaces in the simulation.
- 574 — An electronic simulation, simulating the response waveform of the PMTs, tracking it through  
 575 the digitization process, accounting for effects such as non-linearity, dark noise, Time Tran-

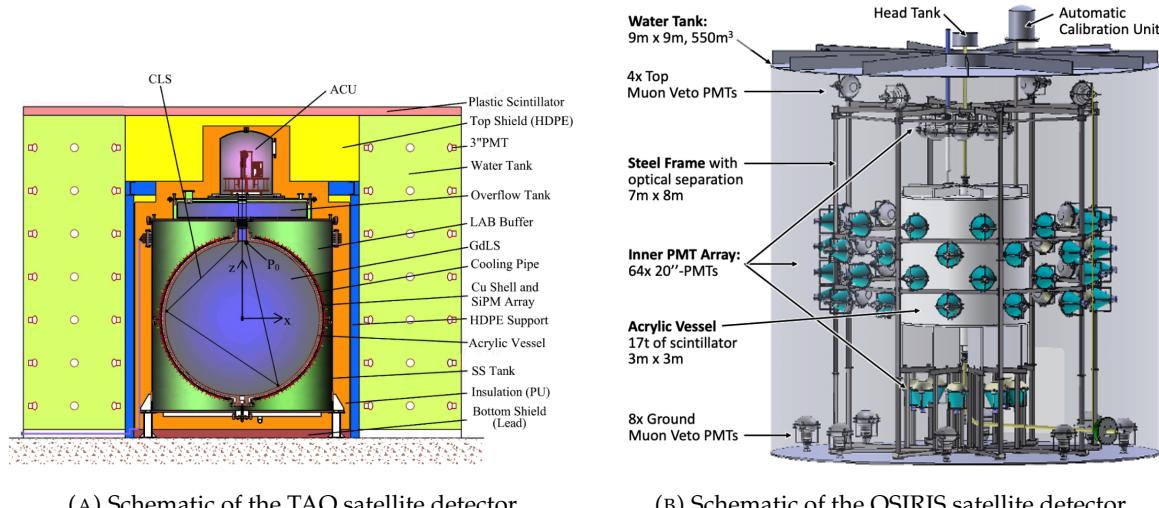


FIGURE 2.15

576 sit Spread (TTS), pre-pulsing, after-pulsing and ringing if the waveform. It's also the step  
 577 handling the event triggers and mixing. This step is commonly referenced as *ElecSim*.

- 578 — A waveform reconstruction where the digitized waveform are filtered to remove high-frequency  
 579 white noise and then deconvoluted to yield time and charge informations of the photons hits  
 580 on the PMTs. This step is commonly referenced as *Calib*.
- 581 — The charge and time informations are used by reconstruction algorithms to reconstruct the  
 582 interaction vertex and the deposited energy. This step is commonly reported as *Reco*. See  
 583 section 2.6 for more details on the reconstruction.
- 584 — Once the singular events are reconstructed, they go through event pairing and classification  
 585 to select IBD events. This step is named Event Classification.
- 586 — The purified signal is then analysed by the analysis framework which depend of the physics  
 587 topic of interest.

588 The steps Reco and Event Classification are divided into two category of algorithm. Fast but less  
 589 accurate algorithms that are running during the data taking designated as the *Online* algorithms.  
 590 Those algorithm are used to take the decision to save the event on tape or to throw it away. More  
 591 accurate algorithms that run on batch of events designated *Offline* algorithms. They are used for the  
 592 physics analysis. The Offline Reco will be one of the main topic of interest for this thesis.

## 593 2.6 State of the art of the Offline IBD reconstruction in JUNO

594 The main reconstruction method currently run in JUNO is a data-driven method based on a like-  
 595 lihood maximization [35, 36] using only the LPMTs. The first step is to reconstruct the interaction  
 596 vertex from which the energy reconstruction is dependent. It is also necessary for event pairing and  
 597 classification.

### 598 2.6.1 Interaction vertex reconstruction

599 To start the likelihood maximization, a rough estimation of the vertex and of the event timing is  
 600 needed. We start by estimating the vertex position using a charge based algorithm.

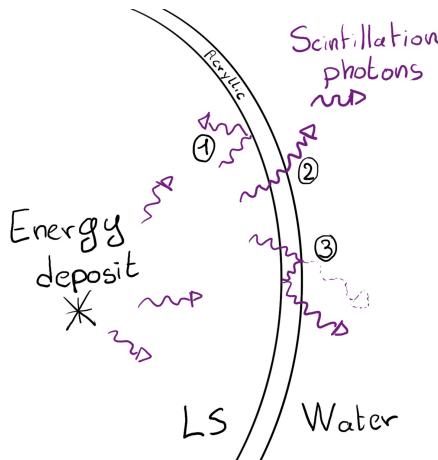
601 **Charge based algorithm**

602 The charge-based algorithm is basically base on the charge-weighted average of the PMT position.

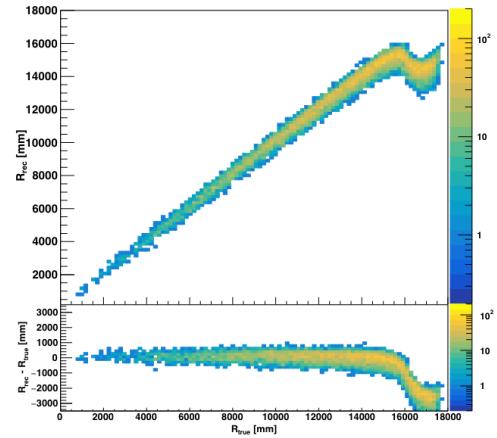
$$\vec{r}_{cb} = a \cdot \frac{\sum_i q_i \cdot \vec{r}_i}{\sum_i q_i} \quad (2.3)$$

603 Where  $q_i$  is the reconstructed charge of the pulse of the  $i$ th PMT and  $\vec{r}_i$  is its position.  $\vec{r}_0$  is the  
 604 reconstructed interaction position.  $a$  is a scale factor introduced because a weighted average over  
 605 a 3D sphere is inherently biased. Using calibration we can estimate  $a \approx 1.3$  [37]. The results in  
 606 figure 2.16b shows that the reconstruction is biased from around 15m and further. This is due to the  
 607 phenomena called “total reflection area” or TR Area.

608 As depicted in the figure 2.16a the optical photons, given that they have a sufficiently large incidence  
 609 angle, can be deviated of their trajectories when passing through the interfaces LS-acrylic and water-  
 610 acrylic due to the optical index difference. This cause photons to be lost or to be detected by PMT  
 611 further than anticipated if we consider their rectilinear trajectories. This cause the charge barycenter  
 612 the be located closer to the center than the event really is.



(A) Illustration of the different optical photons reflection scenarios. 1 is the reflection of the photon at the interface LS-acrylic or acrylic-water. 2 is the transmission of the photons through the interfaces. 3 is the conduction of the photon in the acrylic.



(B) Heatmap of  $R_{rec}$  and  $R_{rec} - R_{true}$  as a function of  $R_{true}$  for 4MeV prompt signals uniformly distributed in the detector calculated by the charge based algorithm

FIGURE 2.16

613 It is to be noted that charge based algorithm, in addition to be biased near the edge of the detector,  
 614 does not provide any information about the timing of the event. Therefore, a time based algorithm  
 615 needs to be introduced to provide initial values.

616 **Time based algorithm**

617 The time based algorithm use the distribution of the time of flight corrections  $\Delta t$  (Eq 2.4) of an event  
 618 to reconstruct its vertex and  $t_0$ . It follow the following iterations:

- 619 1. Use the charge based algorithm to get an initial vertex to start the iteration.

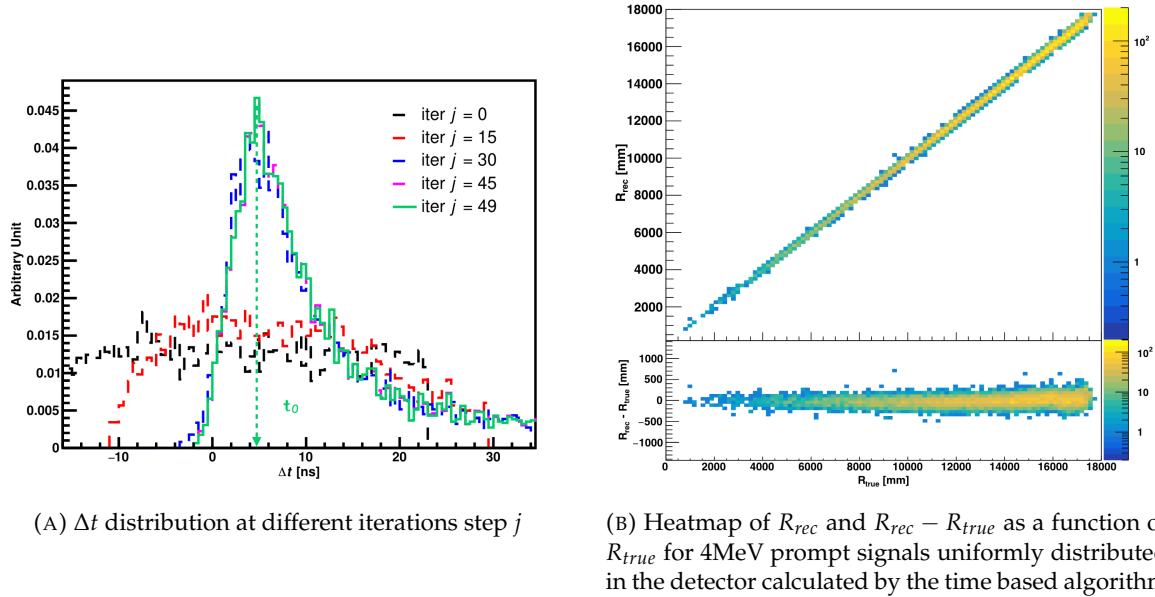


FIGURE 2.17

620 2. Calculate the time of flight correction for the  $i$ th PMT using

$$\Delta t_i(j) = t_i - \text{tof}_i(j) \quad (2.4)$$

621 where  $j$  is the iteration step,  $t_i$  is the timing of the  $i$ th PMT, and  $\text{tof}_i$  is the time-of-flight of the  
622 photon considering an rectilinear trajectory and an effective velocity in the LS and water (see  
623 [37] for detailed description of this effective velocity). Plot the  $\Delta t$  distribution and label the  
624 peak position as  $\Delta t^{\text{peak}}$  (see fig 2.17a).

625 3. Calculate a correction vector  $\vec{\delta}[\vec{r}(j)]$  as

$$\vec{\delta}[\vec{r}(j)] = \frac{\sum_i \left( \frac{\Delta t(j) - \Delta t^{\text{peak}}(j)}{\text{tof}_i(j)} \right) \cdot (\vec{r}_0(j) - \vec{r}_i)}{N^{\text{peak}}(j)} \quad (2.5)$$

626 where  $\vec{r}_0$  is the vertex position at the beginning of this iteration,  $\vec{r}_i$  is the position of the  $i$ th  
627 PMT. To minimize the effect of scattering, dark noise and reflection, only the pulse happening  
628 in a time window (-10 ns, +5 ns) around  $\Delta t^{\text{peak}}$  are considered.  $N^{\text{peak}}$  is the number of PE  
629 collected in this time-window.

630 4. if  $\vec{\delta}[\vec{r}(j)] < 1\text{mm}$  or  $j \geq 100$ , stop the iteration. Otherwise  $\vec{r}_0(j+1) = \vec{r}_0(j) + \vec{\delta}[\vec{r}(j)]$  and go to  
631 step 2.

632 However because the earliest arrival time is used,  $t_i$  is related to the number photoelectrons  $N_i^{\text{pe}}$   
633 detected by the PMT [38–40]. To reduce bias in the vertex reconstruction, the following equation is  
634 used to correct  $t_i$  into  $t'_i$ :

$$t'_i = t_i - p_0 / \sqrt{N_i^{\text{pe}}} - p_1 - p_2 / N_i^{\text{pe}} \quad (2.6)$$

635 The parameters  $(p_0, p_1, p_2)$  were optimized to (9.42, 0.74, -4.60) for Hamamatsu PMTs and (41.31,  
636 -12.04, -20.02) for NNVT PMTs [37]. The results presented in figure 2.17b shows that the time based  
637 algorithm provide a more accurate vertex and is unbiased even in the TR area. This results  $(\vec{r}_0, t_0)$  is  
638 used as initial value for the likelihood algorithm.

639 **Time likelihood algorithm**

640 The time likelihood algorithm use the residual time expressed as follow

$$t_{\text{res}}^i(\vec{r}_0, t_0) = t_i - \text{tof}_i - t_0 \quad (2.7)$$

641 In a first order approximation, the scintillator time response Probability Density Function (PDF) can  
 642 be described as the emission time profile of the scintillation photons, the Time Transit Spread (TTS)  
 643 and the dark noise of the PMTs. The emission time profile  $f(t_{\text{res}})$  is described like

$$f(t_{\text{res}}) = \sum_k \frac{\rho_k}{\tau_k} e^{-\frac{t_{\text{res}}}{\tau_k}}, \sum_k \rho_k = 1 \quad (2.8)$$

644 as the sum of the  $k$  component that emit light in the LS each one characterised by it's decay time  $\tau_k$   
 645 and intensity fraction  $\rho_k$ . The TTS component is expressed as a gaussian convolution

$$g(t_{\text{res}}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_{\text{res}}-\nu)^2}{2\sigma^2}} \cdot f(t_{\text{res}}) \quad (2.9)$$

646 where  $\sigma$  is the TTS of PMTs and  $\nu$  is the average transit time. The dark noise is not correlated with any  
 647 physical events and considered as constant rate over the time window considered  $T$ . By normalizing  
 648 the dark noise probability  $\epsilon(t_{\text{res}})$  as  $\int_T \epsilon(t_{\text{res}}) dt_{\text{res}} = \epsilon_{\text{dn}}$ , it can be integrated in the PDF as

$$p(t_{\text{res}}) = (1 - \epsilon_{\text{dn}}) \cdot g(t_{\text{res}}) + \epsilon(t_{\text{res}}) \quad (2.10)$$

649 The distribution of the residual time  $t_{\text{res}}$  of an event can then be compared to  $p(t_{\text{res}})$  and the best  
 650 fitting vertex  $\vec{r}_0$  and  $t_0$  can be chosen by minimizing

$$\mathcal{L}(\vec{r}_0, t_0) = -\ln \left( \prod_i p(t_{\text{res}}^i) \right) \quad (2.11)$$

651 The parameter of Eq. 2.10 can be measured experimentally. The results shown in figure 2.18 used  
 652 PDF from monte carlo simulation. The results shows that  $R_{\text{rec}} - R_{\text{true}}$  is biased depending on the  
 653 energy. While this could be corrected using calibration, another algorithm based on charge likelihood  
 654 was developed to correct this problem.

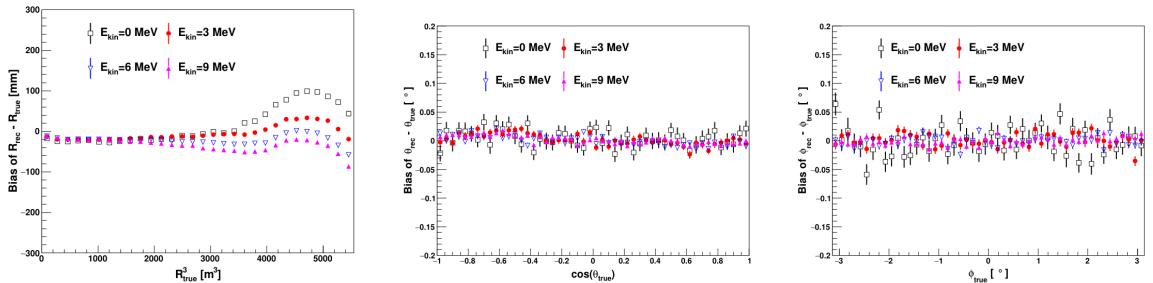


FIGURE 2.18 – Bias of the reconstructed radius  $R$  (left),  $\theta$  (middle) and  $\phi$  (right) for multiple energies by the time likelihood algorithm

655 **Charge likelihood algorithm**

656 Similarly to the time likelihood algorithms that use a timing PDF, the charge likelihood algorithm  
 657 use a PE PDF for each PMT depending on the energy and position of the event. With  $\mu(\vec{r}_0, E)$  the  
 658 mean expected number of PE detected by each PMT, the probability to observe  $N_{pe}$  in a PMT follow  
 659 a Poisson distribution. Thus

660 — The probability to observe no hit ( $N_{pe} = 0$ ) in the  $j$ th PMT is  $P_{nohit}^j(\vec{r}_0, E) = e^{-\mu_j}$

661 — The probability to observe  $N_{pe} \neq 0$  in the  $i$ th PMT is  $P_{hit}^i(\vec{r}_0, E) = \frac{\mu^{N_{pe}} e^{-\mu_i}}{N_{pe}^i!}$

662 Therefore, the probability to observe a specific hit pattern can be expressed as

$$P(\vec{r}_0, E) = \prod_j P_{nohit}^j(\vec{r}_0, E) \cdot \prod_i P_{hit}^i(\vec{r}_0, E) \quad (2.12)$$

663 The best fit values of  $\vec{R}_0$  and  $E$  can then be calculated by minimizing the negative log-likelihood

$$\mathcal{L}(\vec{r}_0, E) = -\ln(P(\vec{r}_0, E)) \quad (2.13)$$

664 In principle,  $\mu_i(\vec{r}_0, E)$  could be expressed

$$\mu_i(\vec{r}_0, E) = Y \cdot \frac{\Omega(\vec{r}_0, r_i)}{4\pi} \cdot \epsilon_i \cdot f(\theta_i) \cdot e^{-\sum_m \frac{d_m}{\zeta_m}} \cdot E + \delta_i \quad (2.14)$$

665 where  $Y$  is the energy scale factor,  $\Omega(\vec{r}_0, r_i)$  is the solid angle of the  $i$ th PMT,  $\epsilon_i$  is its detection  
 666 efficiency,  $f(\theta_i)$  its angular response,  $\zeta_m$  is the attenuation length in the materials and  $\delta_i$  the expected  
 667 number of dark noise.

668 However Eq. 2.14 assume that the scintillation light yield is linear with energy and describe poorly  
 669 the contribution of indirect light, shadow effect due to the supporting structure and the total reflec-  
 670 tion effects. The solution is to use data driven methods to produce the pdf by using the calibra-  
 671 tions sources and position described in section 2.3. In the results presented in figures 2.19, the PDF was  
 672 produced using MC simulation and 29 specific calibrations position [37] along the Z-axis of the  
 detector. We see that the charge likelihood algorithm show little bias in the TR area and a better

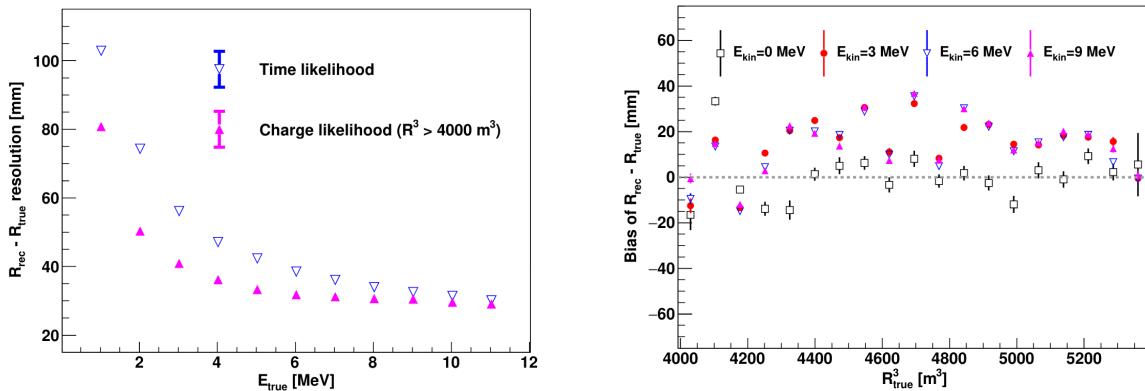


FIGURE 2.19 – On the left: Resolution of the reconstructed  $R$  as a function of the energy in the TR area ( $R^3 > 4000 \text{ m}^3 \equiv R > 16 \text{ m}$ ) by the charge and time likelihood algorithms. On the right: Bias of the reconstructed  $R$  in the TR area for different energies by the charge likelihood algorithm

673 resolution than the time likelihood. The figure 2.20 shows the radial resolution of the different  
 674

675 algorithm presented for this section, we can see the refinement at each step and that the charge  
 676 likelihood yield the best results.

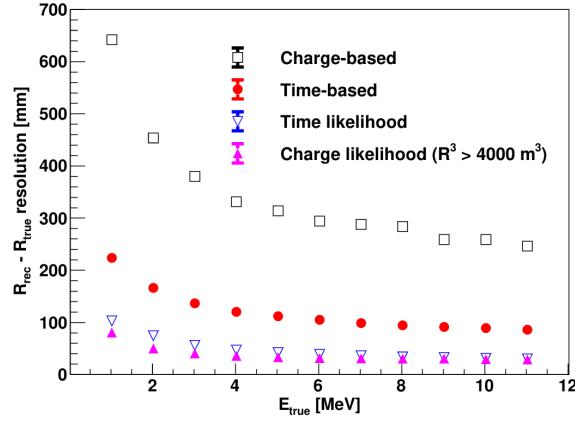
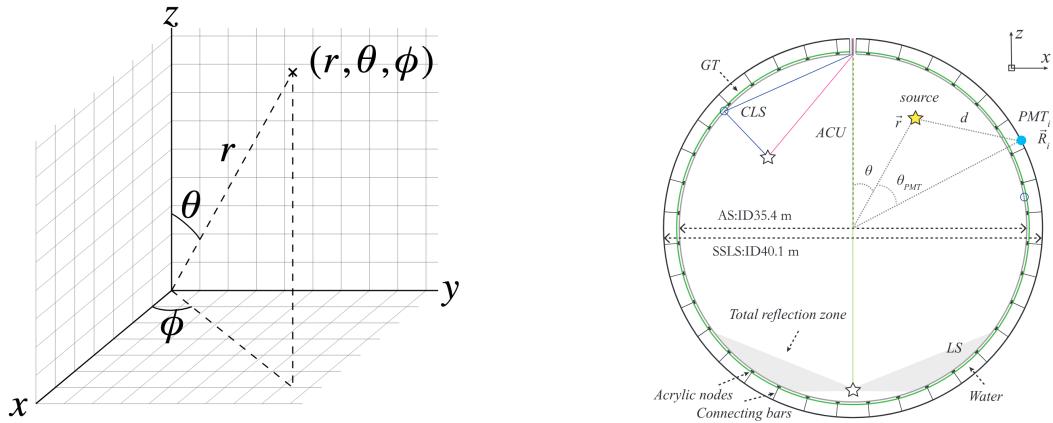


FIGURE 2.20 – Radial resolution of the different vertex reconstruction algorithms as a function of the energy

677 The charge based likelihood algorithms already give some information on the energy as Eq. 2.13  
 678 is minimized but the energy can be further refined as shown in the next section.

## 679 2.6.2 Energy reconstruction

680 As explained in section 2.1.1, energy resolution is crucial for the NMO and oscillation parameters  
 681 measurements. Thus the energy reconstruction algorithm should take into consideration as much  
 682 detector effect as possible. The following method is a data driven method based on calibration  
 683 samples inspired by the charge likelihood algorithm described above [41].



(A) Spherical coordinate system used in JUNO for reconstruction

(B) Definition of the variables used in the energy reconstruction

FIGURE 2.21

684 **Charge estimation**

685 The most important element in the energy reconstruction is  $\mu_i(\vec{r}_0, E)$  described in Eq. 2.14. For  
 686 realistic cases, we also need to take into account the electronics effect that were omitted in the  
 687 previous section. Those effect will cause a charge smearing due to the uncertainties in the  $N_{pe}$   
 688 reconstruction. Thus we define  $\hat{\mu}^L(\vec{r}_0, E)$  which is the expected  $N_{pe}/E$  in the whole detector for an  
 689 event with visible energy  $E_{vis}$  and position  $\vec{r}_0$ . The position of the event and PMTs are now defined  
 690 using  $(r, \theta, \theta_{pmt})$  as defined in figure 2.21b.

$$\hat{\mu}(r, \theta, \theta_{pmt}, E_{vis}) = \frac{1}{E_{vis}} \frac{1}{M} \sum_i^M \frac{\bar{q}_i - \mu_i^D}{\text{DE}_i}, \quad \mu_i^D = \text{DNR}_i \cdot L \quad (2.15)$$

691 where  $i$  runs over the PMTs with the same  $\theta_{pmt}$ ,  $\text{DE}_i$  is the detection efficiency of the  $i$ th PMT.  $\mu_i^D$   
 692 is the expected number of dark noise photoelectrons in the time window  $L$ . The time window have  
 693 been optimized to  $L = 280$  ns [41].  $\bar{q}_i$  is the average recorded photoelectrons in the time window  
 694 and  $\hat{Q}_i$  is the expected average charge for 1 photoelectron. The  $N_{pe}$  map is constructed following the  
 695 procedure described in [36].

696 **Time estimation**

697 The second important observable is the hit time of photons that was previously defined in Eq. 2.7. It  
 698 is here refined as

$$t_r = t_h - \text{tof} - t_0 = t_{LS} + t_{TT} \quad (2.16)$$

699 where  $t_h$  is the time of hit,  $t_{LS}$  is the scintillation time and  $t_{TT}$  the transit time of PMTs that is described  
 700 by a gaussian

$$t_{TT} = \mathcal{N}(\overline{\mu_{TT} + t_d}, \sigma_{TT}) \quad (2.17)$$

701 where  $\mu_{TT}$  is the mean transit time in PMTs,  $\sigma_{TT}$  is the Transit Time Spread (TTS) of the PMTs and  $t_d$   
 702 is the delay time in the electronics. The effective refraction index of the LS is also corrected to take  
 703 into account the propagation distance in the detector.

704 The timing PDF  $P_T(t_r | r, d, \mu_l, \mu_d, k)$  can now be generated using calibration sources [41]. This PDF  
 705 describe the probability that the residual time of the first photon hit is in  $[t_r, t_r + \delta]$  with  $r$  the radius  
 706 of the event vertex,  $d = |\vec{r} - \vec{r}_{PMT}|$  the propagation distance,  $\mu_l$  and  $\mu_d$  the expected number of PE  
 707 and dark noise in the electronic reading window and  $k$  is the detected number of PE.

708 Now let denote  $f(t, r, d)$  the probability density function of "photoelectron hit a time t" for an event  
 709 happening at  $r$  where the photons traveled the distance  $d$  in the LS

$$F(t, r, d) = \int_t^L f(t', r, d) dt' \quad (2.18)$$

710 Based on the PDF for one photon  $k = 1$ , one can define

$$P_T^l(t | k = n) = I_n^l [f_l(t) F_l^{n-1}(t)] \quad (2.19)$$

711 where the indicator  $l$  means that the photons comes from the LS and  $I_n^l$  a normalisation factor. To this  
 712 pdf we add the probability to have photons coming from the dark noise indicated by the indicator  $d$   
 713 using

$$f_d(t) = 1/L, \quad F_d(t) = 1 - \frac{t}{L} \quad (2.20)$$

<sup>714</sup> and so for the case where only one photon is detected by the PMT ( $k = 1$ )

$$P_T(t|\mu_l, \mu_d, k=1) = I_1[P(1, \mu_l)P(0, \mu_d)f_l(t) + P(0, \mu_l)P(1, \mu_d)f_d(t)] \quad (2.21)$$

<sup>715</sup> where  $P(k_\alpha, \mu_\alpha)$  is the Poisson probability to detect  $k_\alpha$  PE from  $\alpha \in \{l, d\}$  with the condition  $k_l + k_d = k$ .  
<sup>716</sup>

<sup>717</sup> Now that we have the individual timing and charge probability we can construct the charge likelihood referred as QMLE:  
<sup>718</sup>

$$\mathcal{L}(q_1, q_2, \dots, q_N | \vec{r}, E_{vis}) = \prod_{j \in \text{unfired}} e^{-\mu_j} \prod_{i \in \text{fired}} \left( \sum_{k=1}^{\infty} P_Q(q_i|k) \cdot P(k, \mu_i) \right) \quad (2.22)$$

<sup>719</sup> where  $\mu_i = E_{vis}\hat{\mu}_i^L + \mu_i^D$  and  $P(k, \mu_i)$  is the Poisson probability of observing  $k$  PE.  $P_Q(q_i|k)$  is the  
<sup>720</sup> charge pdf for  $k$  PE. And we can also construct the time likelihood referred as TMLE:

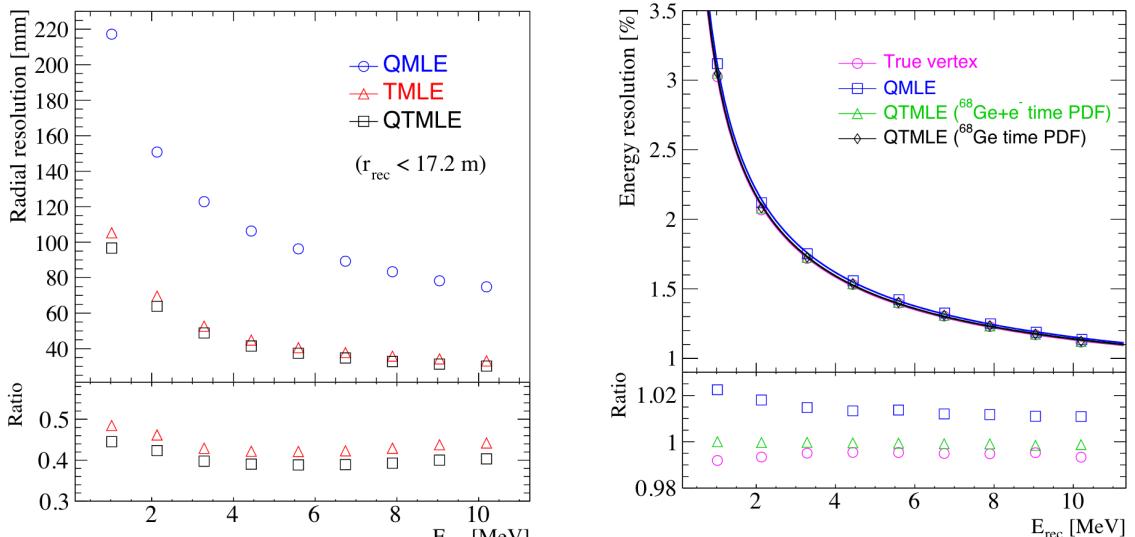
$$\mathcal{L}(t_{1,r}, t_{2,r}, \dots, t_{N,r} | \vec{r}, t_0) = \prod_{i \in \text{hit}} \frac{\sum_{k=1}^K P_T(t_{i,r}|r, d, \mu_i^L, \mu_i^d, k) \cdot P(k, \mu_i^L + \mu_i^d)}{\sum_{k=1}^K P(k, \mu_i^L + \mu_i^d)} \quad (2.23)$$

<sup>721</sup> where  $K$  is cut to 20 PE and hit is the set of hits satisfying  $-100 < t_{i,r} < 500$  ns.  
<sup>722</sup>

Merging those two likelihood give the charge-time likelihood QTMLLE

$$\mathcal{L}(q_1, q_2, \dots, q_N; t_{1,r}, t_{2,r}, \dots, t_{N,r} | \vec{r}, t_0, E_{vis}) = \mathcal{L}(q_1, q_2, \dots, q_N | \vec{r}, E_{vis}) \cdot \mathcal{L}(t_{1,r}, t_{2,r}, \dots, t_{N,r} | \vec{r}, t_0) \quad (2.24)$$

<sup>723</sup> The radial and energy resolutions of the different likelihood are presented in figure 2.22 (from [41]).  
<sup>724</sup> We can see the improvement of adding the time information to the vertex reconstruction and that  
<sup>725</sup> an increase in vertex precision can bring improvement in the energy resolution, especially at low  
<sup>726</sup> energies.



(A) Radial resolutions of the likelihood-based algorithm TMLE, QMLE and QTMLLE

(B) Energy resolution of QMLE and QTMLLE using different vertex resolutions

FIGURE 2.22

<sup>727</sup> Data driven methods prove to be performant in the energy and vertex reconstruction given that we

728 have enough calibrations sources to produce the PDF. In the next section, we'll see another type of  
 729 data-driven method based on machine learning.

### 730 2.6.3 Machine learning for reconstruction

731 Machine learning (ML) is family of data-driven algorithms that are inferring behavior and results  
 732 from a training dataset. A overview of methods and detailed explanation of the Neural Network  
 733 (NN) subfamily can be found in Chapter 3.

734 The power of ML is the ability to model complex response to a specific problem. In JUNO the  
 735 reconstruction problematic can be expressed as follow: knowing that each PMT, large or small,  
 736 detected a given number of PE  $Q$  at a given time  $t$  and their position is  $x, y, z$  where did the energy  
 737 was deposited and how much energy was it, modeling a function that naively goes:

$$\mathbb{R}^{5 \times N_{pmt}} \mapsto \mathbb{R}^4 \quad (2.25)$$

738 It is worth pointing that while this is already a lot in informations, this is not the rawest representa-  
 739 tion of the experiment. We could indeed replace the charge and time by the waveform in the time  
 740 window of the event but that would lead to an input representation size that would exceed our  
 741 computational limits. Also, due to those computational limits, most of the ML algorithm reduce this  
 742 input phase space either by structurally encoding the information (pictures, graph), by aggregating  
 743 it (mean, variance, ...) or by exploiting invariance and equivariance of the experiment (rotational  
 744 invariance due to the sphericity, ...).

745 For machine learning to converge to performant algorithm, a large dataset exploring all the phase  
 746 space of interest is needed. For the following studies, data from the monte carlo simulation presented  
 747 in section 2.5 are used for training. When the detector will be finished calibrations sources will be  
 748 complementarily be used.

#### 749 Boosted Decision Tree (BDT)

750 On of the most classic ML method used in physics in last years is the Boosted Decision Tree (see  
 751 chapter 3.2). They have been explored for vertex reconstruction [42] et for energy reconstruction [42,  
 752 43].

753 For vertex and energy reconstruction a BDT was developed using the aggregated informations pre-  
 754 sented in 2.6.

| Parameter                        | description                          |
|----------------------------------|--------------------------------------|
| $n_{Hits}$                       | Total number of hits                 |
| $x_{cc}, y_{cc}, z_{cc}, R_{cc}$ | Coordinates of the center of charge  |
| $ht_{mean}, ht_{std}$            | Hit time mean and standard deviation |

TABLE 2.6 – Features used by the BDT for vertex reconstruction

755 Its reconstruction performances are presented in figure 2.24.

756 A second and more advanced BDT, subsequently named BDTE, that only reconstruct energy use a  
 757 different set of features [43]. They are presented in the table 2.7

|                  |                  |
|------------------|------------------|
| AccumCharge      | $ht_{5\%-2\%}$   |
| $R_{cht}$        | $pe_{mean}$      |
| $z_{cc}$         | $J_{cht}$        |
| $pe_{std}$       | $\phi_{cc}$      |
| nPMTs            | $ht_{35\%-30\%}$ |
| $ht_{kurtosis}$  | $ht_{20\%-15\%}$ |
| $ht_{25\%-20\%}$ | $pe_{35\%}$      |
| $R_{cc}$         | $ht_{30\%-25\%}$ |

TABLE 2.7 – Features used by the BDTE algorithm.  $pe$  and  $ht$  reference the charge and hit-time distribution respectively and the percentages are the quantiles of those distributions.  $cht$  and  $cc$  reference the barycenters of hit time and charge respectively

### 758 Neural Network (NN)

759 The physics have shown a rising for Neural Network (NN) in the past years for event reconstruction,  
760 notably in the neutrino community [44–47]. Three type of neural networks have explored for event  
761 reconstruction in JUNO Deep Neural Network (DNN), Convolutional Neural Network (CNN) and  
762 Graph Network (GNN). More explanation about those neural network can be found in chapter 3.

763 The CNN are using 2D projection of the detector representing it as an image with two channel, one  
764 for the charge  $Q$  and one for the time  $t$ . The position of the PMTs is structurally encoded in the pixel  
765 containing the information of this PMT. In [42], the pixel is chosen based on a transformation of  $\theta$   
766 and  $\phi$  coordinates to the 2D plane and rounded to the nearest pixel. A sufficiently large image has  
767 been chosen to prevent two PMT to be located in the same pixel. An example of this projection can  
768 be found in figure 2.23. The performances of the CNN can be found in figure 2.24.

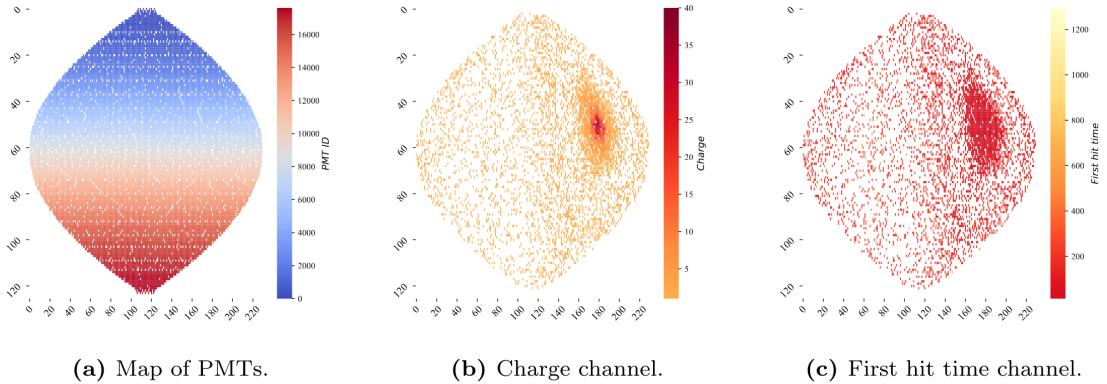


FIGURE 2.23 – Projection of the LPMTs in JUNO on a 2D plane. (a) Show the distribution of all PMTs and (b) and (c) are example of what the charge and time channel looks like respectively

769 Using 2D have the upside of encoding a large part of the informations structurally but loose the rotat-  
770ional invariance of the detector. It also give undefined information to the neural network (what is a  
771pixel without PMT ? What should be its charge and time ?), cause deformation in the representation  
772of the detector (sides of projection) and loose topological informations.

773 One of the way to present structurally the sphericity of JUNO to a NN is to use a graph: A collection  
774 of objects  $V$  called nodes and relations  $E$  called edges, each relation associated to a couple  $v_1, v_2$   
775 forming the graph  $G(E, V)$ . Nodes and edges can hold informations or features. In [42] the nodes,  
776 are geometrical region of the detector as defined by the HealPix [48]. The features of the nodes are

777 aggregated informations from the PMTs it contains. The edges contains geographic informations of  
 778 the nodes relative positions.

779 This data representation has the advantages to keep the topology of the detector intact. It also permit  
 780 the use of rotational invariant algorithms for the NN, thus taking advantage of the symmetries of the  
 781 detector.

782 The neural network then process the graph using Chebyshev Convolutions [49]. The performances  
 783 of the GNN are presented in figure 2.24.

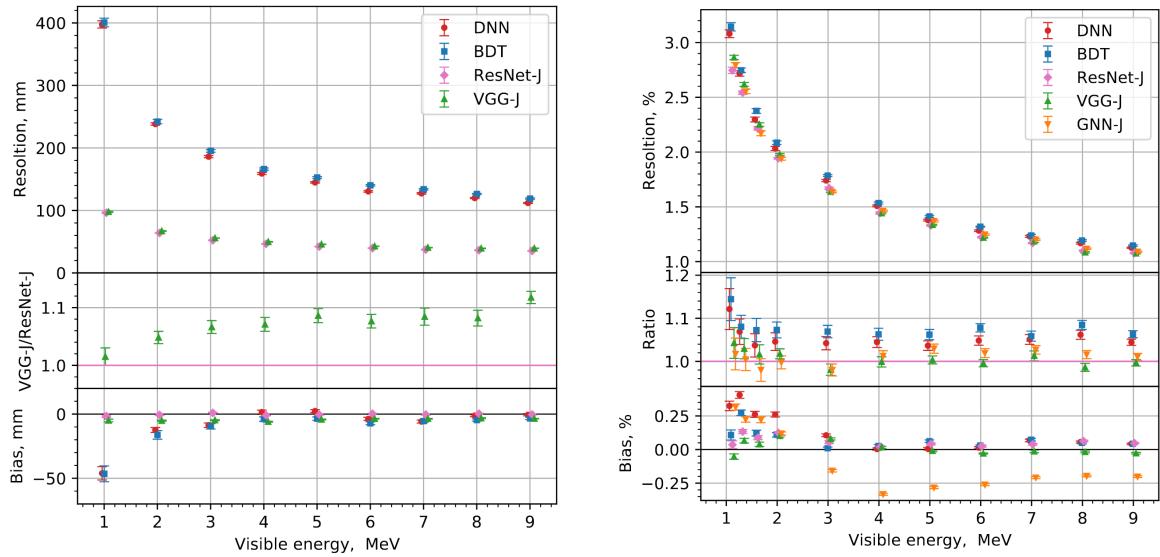


FIGURE 2.24 – Radial (left) and energy (right) resolutions of different ML algorithms.  
 The results presented here are from [42]. DNN is a deep neural network, BDT is a BDT,  
 ResNet-J and VGG-J are CNN and GNN-J is a GNN.

784 Overall ML algorithms show similar performances as classical algorithms in term of energy recon-  
 785 structions with the more complex structure CNN and GNN showing better performances than BDT  
 786 and DNN. For vertex reconstruction, the BDT and DNN show poor performance while CNN are on  
 787 the level of the classical algorithms.

## 788 2.7 JUNO sensitivity to NMO and precise measurements

789 Now that the event have been reconstructed, selected and that the non-IBD background have been  
 790 rejected, we have access to the measured energy flux from JUNO. We consider two spectra, the  
 791 one measured by the LPMT system and the one measured by the SPMT system. This give rise to  
 792 three possible analysis: A LPMT only analysis, a SPMT only analysis and a joint analysis. This joint  
 793 analysis is the subject of the chapter 7 of this thesis.

794 The following details about JUNO measurement is common to the three analysis. The details and  
 795 specific of the joint analysis are detailed in chapter 7.

### 796 2.7.1 Theoretical spectrum

797 To extract the oscillation parameters and the NMO from the measured spectrum, it is compared to a  
 798 theoretical spectrum. This theoretical spectrum is produced based on the theory of the three flavour

799 oscillation (see section 1.3), the measurements produced by the calibration, the input from TAO and  
 800 adjusted Monte Carlo simulations:

- 801 — The absolute flux and the fission product fraction yield calibrated by TAO.
- 802 — The estimation of the neutrinos flux from other sources, such as the geoneutrinos, by theoret-  
 803 ical model.
- 804 — The computed cross-section of  $\bar{\nu}_e$  and the LS.
- 805 — The estimation of mislabelled event, such as fast neutron events from cosmic muons, using  
 806 Monte Carlo simulation.
- 807 — The measured bias and resolution of the LPMT and SPMT system by the calibration.
- 808 — The time dependent reactor parameters (age of fuel, instantaneous power of the reactors, etc...)

809 These systematics parameters come with their uncertainties that need to be taken into account by the  
 810 fitting framework. This theoretical spectrum will, in the end, depend of the oscillation parameters of  
 811 interest  $\theta_{13}, \theta_{12}, \Delta m_{21}^2, \Delta m_{31}^2$ . Noise parameters can be included in the parameters spectrum such as  
 812 the earth density  $\rho$  between the power plants and JUNO.

### 813 2.7.2 Fitting procedure

814 The theoretical and measured spectra are represented as two histograms depending on the energy.  
 815 The theoretical spectrum is adjusted with the data using a  $\chi^2$  minimization where  $\chi^2$  is naively  
 816 defined as

$$\chi^2 = \sum_i \frac{(N_{th}^i - N_{data}^i)^2}{\sigma_i^2} \quad (2.26)$$

817 where  $N_{th}^i$  is the number event in the  $i$ th bin of the theoretical spectrum,  $N_{data}^i$  is the number of event  
 818 in the  $i$ th bin of the measured spectrum and  $\sigma_i$  is the uncertainty of this bin. Two classic statistic test  
 819 exist Pearson and Neyman where the difference is the estimation of  $\sigma_i$  parameters.

820 This  $\sigma_i$  is composed of the systematics uncertainties discussed above but also from the statistic  
 821 uncertainty of the spectrum. Considering a Poisson process, the statistic uncertainty is estimated  
 822 as  $\sigma_{stat}^i = \sqrt{N^i}$ . In a Pearson test,  $N^i \equiv N_{th}^i$  whereas in a Neyman test  $N^i \equiv N_{data}^i$ . Under the  
 823 assumption that the content of each bin follow a Gaussian distribution (a Poisson with high enough  
 824 statistic), the two test are equivalent. But studies on Monte Carlo spectrum showed that the Pearson  
 825 and Neyman statistic are biased in opposite direction. It is easily visible where, for the same data,  
 826 Pearson will prefer a higher  $N_{th}^i$  to reduce the ration  $\frac{1}{N_{th}^i}$  whereas Neyman will prefer a lower  $N_{th}^i$  to  
 827 reduce the  $(N_{th}^i - N_{data}^i)$  term.

828 This problematic can be circumvented by summing the two test, yielding the CNP statistic test  
 829 and/or by adding a term

$$\chi^2 = \sum_i \frac{(N_{th}^i - N_{data}^i)^2}{\sigma_i^2} - \ln |V| \quad (2.27)$$

830 where  $V$  is the covariance matrix of the theoretical spectrum yielding the PearsonV and CNPV  
 831 statistic test.

832 The  $\chi^2$  is minimized by exploring the parameter phase space via gradient descent.

### 833 2.7.3 Physics results

834 The oscillation parameters are directly extracted from the minimization procedure and the error can  
 835 be estimated directly from the procedure. For the NMO, the data are fitted under the two assumption  
 836 of NO and IO. The difference in  $\chi^2$  give us the preferred ordering and the significance of our test.  
 837 Latest studies show that the precision on oscillation parameters after six year of data taking will be

of 0.2%, 0.3%, 0.5% and 12.1% for  $\Delta m_{31}^2$ ,  $\Delta m_{21}^2$ ,  $\sin^2 \theta_{12}$  and  $\sin^2 \theta_{13}$  respectively [11]. The expected sensitivity to mass ordering is  $3\sigma$  after 6.5 years [50].

## 2.8 Summary

JUNO is one the biggest new generation neutrino experiment. Its goal, the measurements of oscillation parameters with unprecedented precision and an NMO preference at the 3 sigma confidence level, needs an in depth knowledge and understanding of the detector and the physics at hand. The characterisation and calibration of the detector are of the utmost importance and the understanding of the detector response in its resolution and bias is capital to be able to correctly carry the high precision physics analysis of the neutrino oscillation.

In this thesis, I explore the usage of data-driven reconstruction methods to validate and optimize the reconstruction of IBD events in JUNO in the chapters 4, 5 and 6 and the usage of the dual calorimetry in the detection of possible mis-modelisation in the theoretical spectrum 7.

850 **Chapter 3**

851 **Machine learning: Introduction to the  
852 methods and algorithms used in this  
853 thesis**

854 “I have the shape of a human being and organs equivalent to those of a  
855 human being. My organs, in fact, are identical to some of those in a  
856 prosthетized human being. I have contributed artistically, literally, and  
857 scientifically to human culture as much as any human being now  
858 alive. What more can one ask?”

Isaac Asimov, *The Complete Robot*

855 **Contents**

---

|                  |  |    |
|------------------|--|----|
| 856       3.1    | <b>Core concepts in machine learning and neural networks</b> | 42 |
| 857       3.2    | <b>Boosted Decision Tree (BDT)</b>                           | 42 |
| 858        3.2.1 | Artificial Neural Network (NN)                               | 42 |
| 859        3.2.2 | Training procedure   | 44 |
| 860        3.2.3 | Potential pitfalls   | 47 |
| 861       3.3    | <b>Neural networks architectures</b>                         | 50 |
| 862        3.3.1 | Fully Connected Deep Neural Network (FCDNN)                  | 50 |
| 863        3.3.2 | Convolutional Neural Network (CNN)                           | 50 |
| 864        3.3.3 | Graph Neural Network (GNN)                                   | 52 |
| 865        3.3.4 | Adversarial Neural Network (ANN)                             | 54 |

---

866 Machine Learning (ML) and more specifically Neural Network (NN) are families of data-driven  
867 algorithms. They are used in a wide variety of domains including natural language processing,  
868 computer vision, speech recognition and, the subject of this thesis, scientific studies.

869 They are used to model complex distributions from a finite dataset to extract a generalist behavior.  
870 For example, in our case, it could be an algorithm that would differentiate the nature of a particle  
871 interacting in the liquid scintillator, between a positron and an electron, based on the readout charge  
872 and time ( $Q, t$ ) of the 17612 LPMT of the JUNO experiment. During a first training phase, it would  
873 learn the discriminative features between the two in the 35224-dimensional charge and time distri-  
874 bution, built from samples of  $e^+$  and  $e^-$  events.

875 It would learn to derive from a complex, highly dimensional set of data the essential few informations  
876 characterizing the interactions: a three body energy deposition (the positron and two annihilation  
877 gammas) and the single deposit from an electron.

878 Ideally, the algorithm would learn to recognize those informations on its own, regardless of the input  
879 size and complexity. In practice, however, these algorithms are guided by human design through

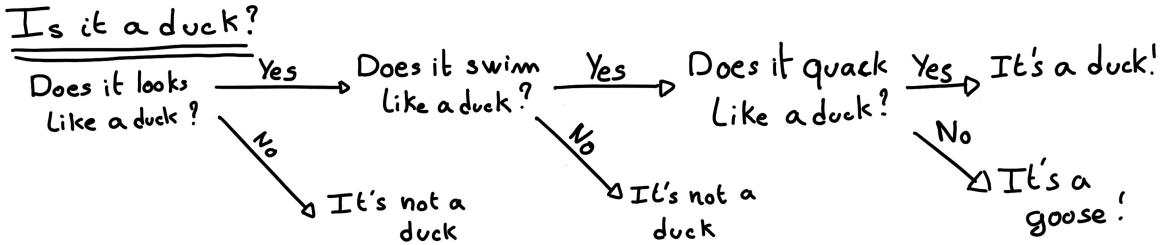


FIGURE 3.1 – Example of a BDT that determine if the given object is a duck

their architectures and training conditions. We can still hope that they can use more thoroughly the detector informations while traditional methods are often subject to assumptions or simplifications to make the task easier (see for instance the algorithm in section 2.6).

The role of machine learning algorithms has expanded rapidly in the past decade, either as the main or secondary algorithm for a wide variety of tasks: event reconstruction, event classification, waveform reconstruction and so on. In particular in domains where the underlying physic and detector processes are complex and highly dimensional, and when large amount of data must be processed quickly.

This chapter present an overview of the different kind of machine learning methods and neural networks that will be discussed in this thesis.

### 3.1 Core concepts in machine learning and neural networks

In this section, we discuss the core concepts in machine learning that will be used thorough this thesis. We place particular emphasis on Neural Networks, as it's the family of the algorithms described in chapters 4, 5 and 6.

### 3.2 Boosted Decision Tree (BDT)

One of the most classic machine learning algorithm used in particle physics is Boosted Decision Tree (BDT) [51] (or more recently Gradient Boosting Machine [52]). The principle of a BDT is fairly simple : based on a set of observables, a serie of decisions, represented as node in a tree, are taken by the algorithm. Each decision point, or node, takes its decision based on a set of trainable parameters leading to a subtree of decisions. The process is repeated until it reach the final node, yielding the prediction. A simplistic example is given in figure 3.1.

The training procedure follow a simple score reward procedure. During the training phase the prediction of the BDT is compared to a known truth about the data. The score is then used to backpropagate corrections to the parameters of the tree. Modern BDT use gradient boosting where the gradient of the loss is calculated for each of the BDT parameters. Following the gradient descent, we can reach the, hopefully, global minima of the loss for our set of parameters.

#### 3.2.1 Artificial Neural Network (NN)

One of the modern ML family is the Neural Network, historical name as their design was inspired by the behaviour of biological neurons in the brain. As schematized in figure 3.2, the input, output and steps inside the NN is described as neuron *layers*. The neurons of the layers take as input a

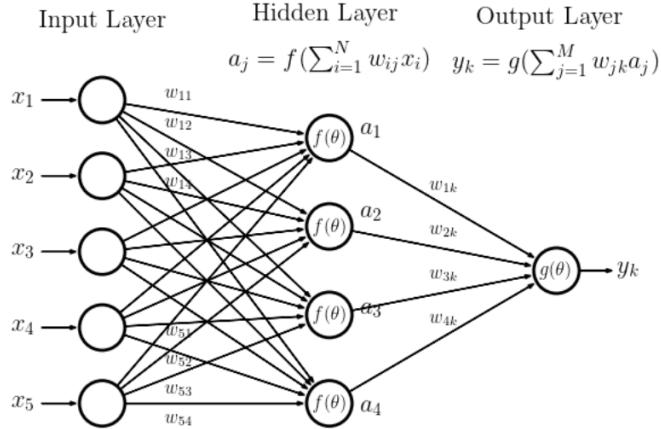


FIGURE 3.2 – Schema of a simple neural network

set of values from the preceding layer, here the  $a_i$  takes every informations of the  $x_i$  input layer, and aggregate those values following learnable *parameters*  $w_{ij}$ . The aggregation procedure is core of defining the architecture of the NN. The different architectures used in this thesis will be discussed in section 3.3. The process is repeated until reaching the output layer.

For example, let's take the network in figure 3.2 and say that  $a_1$ ,  $a_2$  and  $a_3$  are the neurons of the output layer. We try to produce a vertex reconstruction algorithm that will approach the charge barycentre. Let's limit the input  $x_i$  to the charge of the  $i$ th PMT, one of the solution is to aggregate on  $a_1$  the  $x$  coordinate of the barycenter. The network would thus adapt the  $w_{i1}$  parameters so they correspond to the  $x$  coordinates of the  $i$ th PMT. Same for the  $y$  and  $z$  coordinate on  $a_2$  and  $a_3$  respectively.

The layers used in the example above are designated as *Fully connected* layers, where every neurons of the layer is connected to the every neurons of the preceding layer. The layer can be expressed using the Einstein summation and in bold the learnable parameters

$$O_j = I_i + \mathbf{W}_j^i \quad (3.1)$$

where  $O_j$  is the output neurons vector (the  $a_i$ ),  $I_i$  is the preceding layer neurons vector (the  $x_i$ ) and  $\mathbf{W}$  is the parameters, or weights, matrix (composed of the  $w_{ij}$ ). In practice, this fully connected layer is often adjoined a bias  $B$  and an *activation function*  $F$ .

$$I_j = F(I_i \mathbf{W}_j^i + B_j) \quad (3.2)$$

This is the fundamental component of the Fully Connected Deep NN (FCDNN) family presented in section 3.3.1.

This description of neural networks as layers introduce the principles of *depth* and *width*, the number of layers in the NN and the number of neurons in each layer respectively. Those quantities that not directly used for the computation of the results but describes the NN or its training are designated as *hyperparameters*.

Now we just need to adapt the parameters so that this network learn that  $w_{ij}$  are the PMT coordinate. We describe the space produced by the parameters of the network as the *parameter phase space* or *latent space*. The optimization of the network and exploration of this phase space is done through training as described in next section.

940 **3.2.2 Training procedure**

941 To adapt the parameters we need an object that describe how well the network perform. This is  
 942 the *loss* of our neural networks  $\mathcal{L}$ . In our barycenter example, it could be the distance between the  
 943 reconstructed and real barycenter. Using this metric we can adjust the parameters of our network.

944 Depending if we try to minimize or maximize it, it need to posses a minima or a maxima. For example  
 945 when doing *regression*, i.e. produce a scalar result like the coordinates of a barycenter, a common loss  
 946 is the Mean Square Error (MSE). Let  $i$  be our dataset, the  $N$  events considered for training,  $y_i$  be the  
 947 target scalar, the barycenter positions of each events,  $x_i$  the input data, the charge vector, and  $f(x_i, \theta)$   
 948 the result of the network. The network here is modelled by  $f$ , and its parameter  $\theta$

$$\mathcal{L} \equiv MSE = \frac{1}{N} \sum_i^N (y_i - f(x_i, \theta))^2 \quad (3.3)$$

949 Another common loss function is the Mean Absolute Error (MAE)

$$\mathcal{L} \equiv MAE = \frac{1}{N} \sum_i^N |y_i - f(x_i, \theta)| \quad (3.4)$$

950 We see that those loss function possess a minima when  $f(x_i, \theta) = y_i$ .

951 Most of the modern neural networks use gradient descent to optimize their parameters, i.e. the  
 952 gradient of the parameter  $w$ , designated in literature as  $\theta$ , with respect of the loss function  $\mathcal{L}$  is  
 953 subtracted each optimisation step  $t$

$$\theta_{t+1} = \theta_t - \frac{\partial \mathcal{L}}{\partial \theta} \quad (3.5)$$

954 This induce  $\mathcal{L}$  needs to be differentiable with respect to  $\theta$ , thus the layers and their activation func-  
 955 tions also need to be differentiable. This simple gradient descent, designated as Stochastic Gradient  
 956 Descent (SGD), can be extended with first and second order momentums like in the Adam optimizer  
 957 [53]. More details about the optimizers can be found in section 3.2.2.

958 **Training lifecycle**

959 The training of NN does not follow strict rules, you could imagine totally different lifecycle but I will  
 960 describe here the one used in this thesis, the most common one.

961 As illustrated in figure 3.3, the training is split into *epochs*. Each epochs is split into *step* where the  
 962 NN will optimize its parameters over a *batch*, a sub-sample of the training datasets. The ideal batch  
 963 size, number of event in a batch, would be the entire dataset, as the NN optimization would not be  
 964 biased by the specificity of a sub-sample, but due to memory limitations the batch size is driven by  
 965 technical limitations.

966 At the end of each epochs, the neural network is evaluated over a validation dataset, a dataset from  
 967 which no optimisation is done. It is used as reference for the network performance as and monitor  
 968 overtraining (see section 3.2.3).

969 Hyperparameters that can be optimized during the training can be optimized at each epoch, for  
 970 example the learning rate, or each step, the optimizer momentum for example.

971 There is not really a typical number of epochs or steps for the training. The number steps can be  
 972 defined such as in one epoch, the NN see the entirety of the dataset but the number of steps and  
 973 epochs are hyperparameters that are optimized over the each subsequent training. We adjust them  
 974 by looking at the loss evolution profile over time.

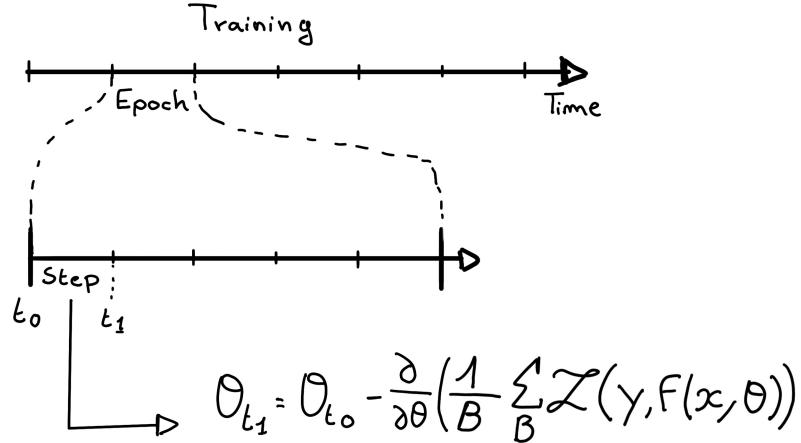


FIGURE 3.3 – Illustration of the training lifecycle

975 Most training are started with a fixed number of epochs, i.e. from what we've seen from precedent  
 976 training, the network stop learning, the loss is constant, after  $N$  epoch so we run the training for  
 977  $N + \delta$  epochs to see if the modification brings improvements to the loss profile. We can setup what's  
 978 called *early stopping policies* that'll stop the training early in specific cases like loss explosion or loss  
 979 stability but this require fine tuning and don't bring much in our case as we are not really limited in  
 980 training time.

### 981 The optimizer

982 As briefly introduced at the beginning of this section, the parameters of the neural network are  
 983 optimized using the gradient descent method. We compute the gradient of the mean loss over the  
 984 batch with respect of each parameters and we update the parameters in accord to minimize the loss.  
 985 The gradient is computed backward from the loss up to the first layer parameters using the chain  
 986 rule:

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \theta_2}{\partial \theta_1} \frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \theta_2}{\partial \theta_1} \frac{\partial \theta_3}{\partial \theta_2} \frac{\partial \mathcal{L}}{\partial \theta_3} = \frac{\partial \theta_2}{\partial \theta_1} \prod_{i=2}^{N-1} \frac{\partial \theta_{i+1}}{\partial \theta_i} \frac{\partial \mathcal{L}}{\partial \theta_N} \quad (3.6)$$

987 where  $\theta$  is a parameter,  $i$  is the layer index. We see here that the gradient of the first layer is  
 988 dependent of the gradient of all the following layers. Because the only value known at the start  
 989 of the optimization procedure is  $\mathcal{L}$  we compute  $\frac{\partial \mathcal{L}}{\partial \theta_N}$  then,  $\frac{\partial \theta_N}{\partial \theta_{N-1}}$ , etc... This is called the *backward  
 990 propagation*.

991 This update of the parameters is done following an optimizer policy. Those optimizers depends on  
 992 hyperparameters. The ones used in this thesis are:

- 993 1. SGD (Stochastic Gradient Descent). This is the simplest optimizer, it depend on only one  
 994 hyperparameter, the learning rate  $\lambda$  (LR) and update the parameters  $\theta$  following

$$\theta_{t+1} = \theta_t - \lambda \frac{\partial \mathcal{L}}{\partial \theta} \Big|_{\theta_t} \quad (3.7)$$

995 where  $t$  is the step index. It is a powerful optimizer but is very sensible to local minima of the  
 996 loss in the parameters phase space as illustrated in figure 3.4a.

- 997 2. Adam [53]. The concept is, in short, to have and SGD but with momentum. Adam possess  
 998 two momentum  $m(\beta_1)$  and  $v(\beta_2)$  which are respectively proportional to  $\frac{\partial \mathcal{L}}{\partial \theta}$  and  $(\frac{\partial \mathcal{L}}{\partial \theta})^2$ .  $\beta_1$

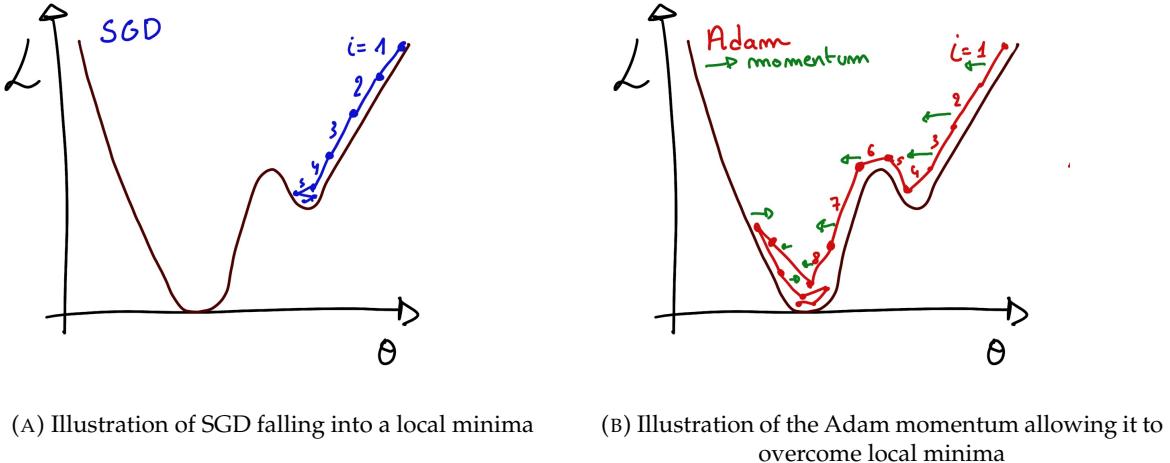


FIGURE 3.4

and  $\beta_2$  are hyperparameters that dictate the moment update at each optimization step. The parameters are then upgraded following

$$m_{t+1} = \beta_1 m_t + (1 - \beta_1) \frac{\partial \mathcal{L}}{\partial \theta} \quad (3.8)$$

$$v_{t+1} = \beta_2 v_t + (1 - \beta_2) \left( \frac{\partial \mathcal{L}}{\partial \theta} \right)^2 \quad (3.9)$$

$$\theta_{t+1} = \theta_t - \lambda \frac{m_{t+1}}{\sqrt{v_{t+1}} + \epsilon} \quad (3.10)$$

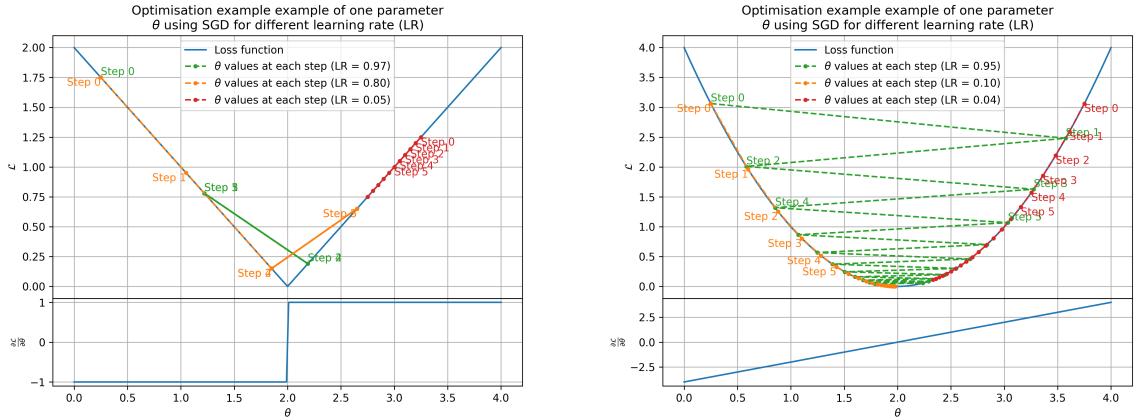
where  $\epsilon$  is a small number to prevent divergence when  $v$  is close to 0. These momentums allow to overcome small local minima in the parameters phase. Imagine ball going down a slope as illustrated in 3.4a, if you ignore the stored momentum you get SGD and get stuck as on the left plot. Now if you consider the momentum you get over the hill and end up in the global minima.

The LR is a crucial parameter in the training of NN. You see that in case of MAE in figure 3.5a that if the LR is too high, you can end up missing the minima. Is the LR is too low, even with MSE as in figure 3.5b, you never reach the minima in the allocated number of epochs. To prevent possible issues, we setup scheduler policies.

#### Scheduler policies

Sometimes we want to update our hyperparameters or take a set of action during the training procedure. We use for this scheduler policies, for example a common policy is a decrease of the learning rate after each epochs. We want to get the closest possible in early epochs before refining the training with a smaller learning rate, finer step. By reducing the learning rate, we allow it to make more fine steps in the parameters phase space, hopefully converging to the true minima.

Another policy that is often used is the save of the best model. In some situations, the loss value after each epoch will strongly oscillate or can even worsen. This policy allows us to keep the best version of the model attained during the training phase.



(A) Illustration of the SGD optimizer on one parameter  $\theta$  on the MAE Loss. We see here that it has trouble reaching the minima due to the gradient being constant.

(B) Illustration of the SGD optimizer on one parameter  $\theta$  on the MSE Loss. We see two different behavior: A smooth one (orange and red) when the LR is small enough and a more chaotic one when the LR is too high.

FIGURE 3.5 – Illustration of the SGD optimizer. In blue is the value of the loss function, orange, green and red are the path taken by the optimized parameter during the training for different LR.

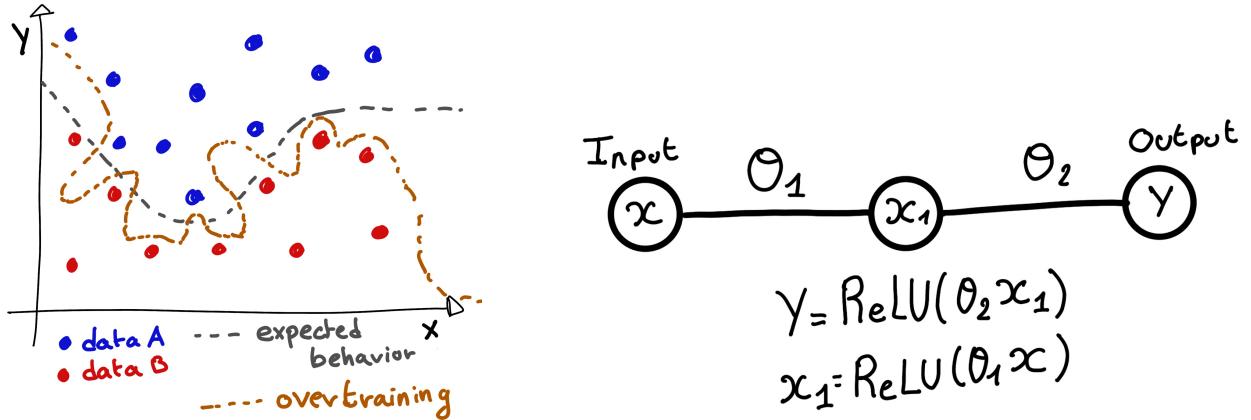
### 3.2.3 Potential pitfalls

Apart from being stuck in local minima, there is also other behaviors and effects we want to prevent during training.

#### Overtraining

This happens when the network learns the specificities of the training dataset instead of a more general representation of the underlying data distribution. This can happen if there is not enough data in comparison to the number of learning parameters, if the training data possess specific features that are not representative of the application dataset or if the NN trains for too long on the same dataset. This behavior is illustrated in figure 3.6a. Overtraining can be fought in multiple ways, for example:

- **More data.** By having more data in the training dataset, the network will not be able to learn the specificities of every data.
- **Less parameters.** By reducing the number of parameters, we reduce the computing and learning capacities of the network. This will force it to fallback to generalist behaviours.
- **Dropout.** This technique implies to randomly set some neurons to 0, i.e. cutting the relation between two neurons in a layer. By doing this, we force the network to allocate more of its parameters to the features learning, preventing those parameters to be used for overtraining.
- **Early stopping.** During the training we monitor the network performance over a validation dataset. The network does not train on this dataset and thus cannot learn its specificities. If the loss on the training dataset diverges too much from the loss on the validation dataset, we can stop the training earlier to prevent it from overtraining.



(A) Illustration of overtraining. The task at hand is to determine depending on two input variable  $x$  and  $y$  if the data belong to the dataset  $A$  or the dataset  $B$ . The expected boundary between the two dataset is represented in grey. A possible boundary learnt by overtraining is represented in brown.

(B) Illustration of a very simple NN

FIGURE 3.6

### 1035 Gradient vanishing

1036 Gradient vanishing is the effect of the gradient being so small for the early layers that the parameters  
 1037 are barely updated after each step. This cause the network to be unable to converge to the minima.

1038 This comes from the way the gradient descent is calculated. Imagine a simple network composed of  
 1039 three fully connected layers: the input layer, a intermediate layer and the output layer. Let  $L$  be the  
 1040 loss,  $\theta_1$  the parameter between the input and the intermediate layer and  $\theta_2$  the parameter between  
 1041 the intermediate and output layer. This network is schematized in figure 3.6b.

1042 The gradient for  $\theta_1$  will be computed using the chain rule presented in equation 3.6. Because  $\theta_1$   
 1043 depends on  $\theta_2$ , if the gradient of  $\theta_2$  is small, so will be the gradient of  $\theta_1$ . Now if we would have  
 1044 much more layer, we can see how the subsequent multiplication of small gradients would lead to  
 1045 very small update of the parameters thus "vanishing gradient".

1046 Multiple actions can be taken to prevent this effect such as:

- 1047 — **Batch normalization:** In this case we apply a normalization layer that will normalize the data.  
 1048 It means that we transform the input variable  $X$  into a variable  $D$  which distribution follow  
 1049  $\langle D \rangle = 0$  and  $\sigma_D = 1$ . This helps the parameters of the network to maintain an appropriate  
 1050 scale.
- 1051 — **Residual Network (ResNet)** [54]: Residual network is a technique for neural network in  
 1052 which, instead of just sequentially feeding the results of each layer to the next one, you  
 1053 compute a residual over the input data. This technique is illustrated in figure 3.7. The  
 1054 reference [54] show empirical evidence of its relevance.

### 1055 Gradient explosion

Gradient explosion happens when the consecutive multiplication of gradient cause exponential grow in the parameter value or if the training lead the network in part of the parameter space where the gradient is significantly higher than usual. For illustration, consider that the loss dependency in  $\theta$

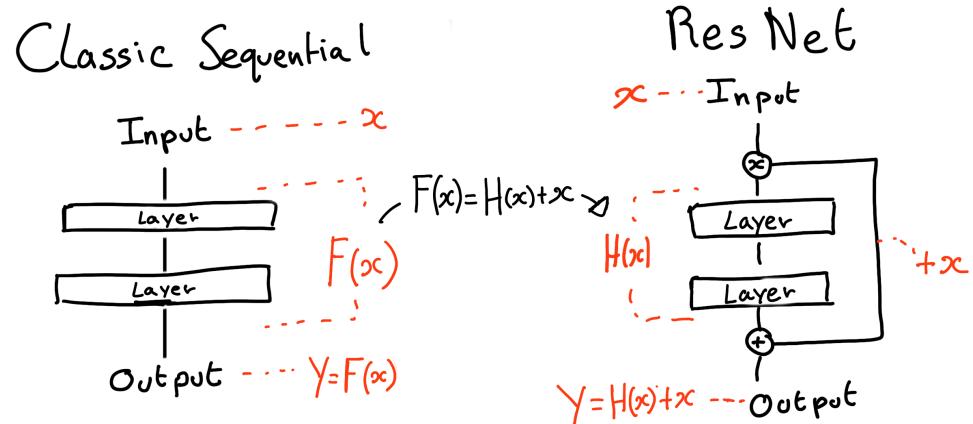


FIGURE 3.7 – Illustration of the ResNet framework

follow

$$\begin{aligned}\mathcal{L}(\theta) &= \frac{\theta^2}{2} + e^{4\theta} \\ \frac{\partial \mathcal{L}}{\partial \theta} &= \theta + 4e^{4\theta}\end{aligned}$$

The explosion is illustrated in figure 3.8 where we can see that the loss degrades with each step of optimization. In this illustration it is clear that reducing the learning rate suffice but this behaviour can happens in the middle of the training where the learning rate schedule does not permit reactivity.

There exist solutions to prevent this explosions:

- **Gradient clipping:** Is this case we work on the gradient so that the norm of gradient vector does not exceed a certain threshold. In our illustration in figure 3.8 the gradient for  $\theta > 0$  could be clipped at 3 for example.
- **Batch normalization:** For the same reasons as for gradient vanishing, normalizing the input data help reduce erratic behaviour.

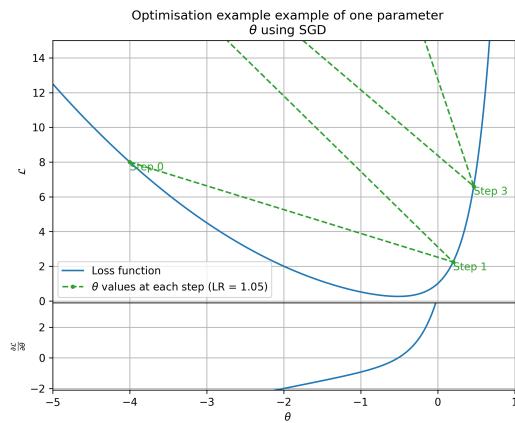


FIGURE 3.8 – Illustration of the gradient explosion. Here it can be solved with a lower learning rate but its not always the case.

1065 **3.3 Neural networks architectures**

1066 **3.3.1 Fully Connected Deep Neural Network (FCDNN)**

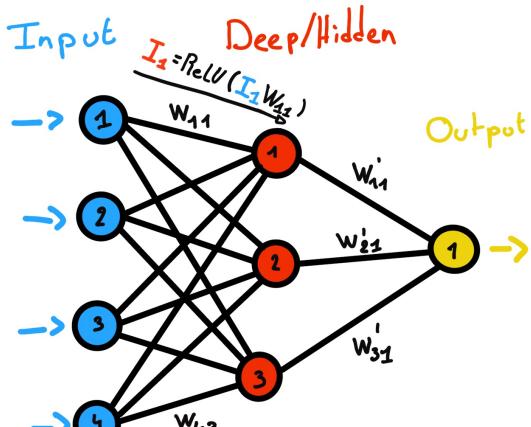
1067 The Fully Connected Deep Neural Network (FCDNN) architecture is the stack of multiple fully  
 1068 connected layers as presented in the figure 3.9a. Most of the time, the classic ReLU function

$$\text{ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (3.11)$$

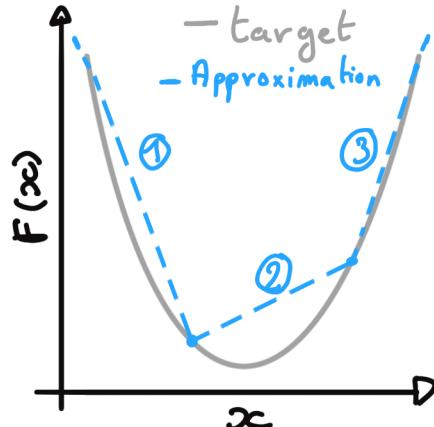
1069 is used as activation function. PReLU and Sigmoid are also popular choices:

$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (3.12) \quad \text{PReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{otherwise} \end{cases} \quad (3.13)$$

1071 The reasoning behind ReLU and PReLU is that with enough of them, you can mimic any continuous  
 1072 function as illustrated in figure 3.9b. Sigmoid is more used in case of classification, its behavior going  
 1073 hand in hand with the Cross Entropy loss function used in classification problems.



(A) Schema of a FCDNN



(B) Illustration of a composition of ReLU "approximating" a function. (1) No ReLU is taking effect (2) One ReLU is activating (3) Another ReLU is activating

FIGURE 3.9

1074 Due to its simplicity, FCDNN are also used as basic pieces for more complex architectures such as  
 1075 the CNN and GNN that will be presented in the next sections.

1076 **3.3.2 Convolutional Neural Network (CNN)**

1077 It's not trivial to describe in text the principles of Convolutional Neural Network (CNN) and how  
 1078 they works. We try a general description below followed by a step by step description of a concrete  
 1079 example.

1080 Convolutional Neural Networks are a family of neural networks that use discrete convolution filters,  
 1081 as illustrated in an example in figure 3.10, to process the input data, often images. They are com-  
 1082 monly used in image recognition [55] for classification or regression problematics. Concretely, you  
 1083 multiply element-wise a portion of the input data, in the case of an image, a small part of the image,

1084 with a kernel of same dimension. In figure 3.10, we multiply the  $3 \times 3$  pixels sub-image with the  
 1085  $3 \times 3$  kernel.

1086 Their filters scan the input data, highlighting patterns of interest, this scanning procedure making  
 1087 them translation-invariant. In the concrete case of figure 3.10, for each pixel of the input image, we  
 1088 group it with the 8 neighbours pixel and produce a new pixel that correspond to the output image.  
 1089 For the pixel on the edges that do not have neighbours, we either create “imaginary” pixel with the  
 1090 value 0 or we just ignore them. If we ignore them, the output image will posses fewer pixels than the  
 1091 input image. We see that the operation do not care where is the pattern of interest in the images, the  
 1092 filter output will be *invariant* whatever *translation* is applied to the image.

1093 This invariance mean that they are capable of detecting oriented features independently of their  
 1094 location on the image. Again taking 3.10 as an example, with only the 9 parameters composing the  
 1095 kernel, we can highlight the contour of the duck by looking at the “yellowness” of the pixels.

1096 The learning parameters of CNNs are the kernels components, the network thus learn the optimal  
 1097 filters to extract the desired features.

1098 The convolution layers are commonly chained [56], reducing the input dimension while increasing  
 1099 the number of filters. The idea behind is that the first layers will process local informations and  
 1100 the latest layers will process more global informations, as the latest convolution filters will process  
 1101 the results of the preceding that themself have processed local information. To try to preserve the  
 1102 amount of information, we tend to grow the numbers of filters for each division of the input data.  
 1103 The results of the convolution filters is commonly then flattened and feed to a smaller FCDNN which  
 1104 will process the filters results to yield the desired output.

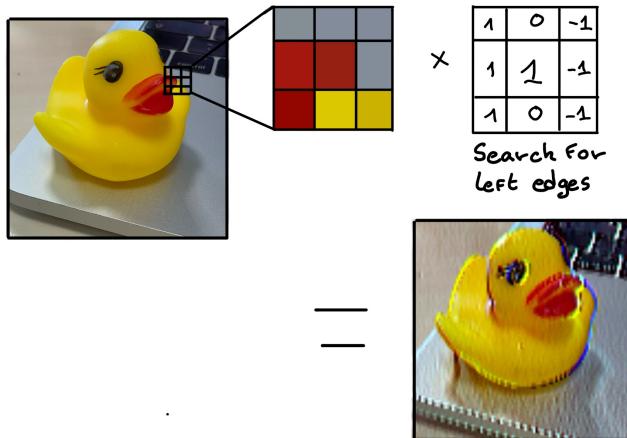


FIGURE 3.10 – Illustration of the effect of a convolution filter. Here we apply a filter with the aim do detect left edges. We see in the resulting image that the left edges of the duck are bright yellow where the right edges are dark blue indicating the contour of the object. The convolution was calculated using [57].

1105 As an example, let’s take the Pytorch [58] example for the MNIST [59], a dataset of black and white  
 1106 images of handwritten digits. Those images are  $28 \times 28$  pixels with only one channel corresponding  
 1107 to the grey level of the pixel. Example of images from this dataset are presented in figure 3.11a

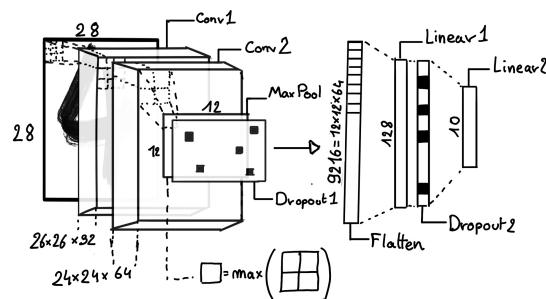
1108 A schema of the CNN used in the Pytorch example is presented in figure 3.11b. Using this schema  
 1109 as a reference, the trained network is made of:

- 1110 1. A convolutional layer of  $(3 \times 3)$  filters yielding 32 channels. A bias parameter is applied  
 1111 to each channel for a total of  $(32 \cdot (3 \times 3) + 32) = 320$  parameters. The resulting image is  
 1112  $(26 \times 26 \times 32)$  (26 per 26 pixels with 32 channels). The ReLU activation function is applied to  
 1113 each pixel.

- 1114 2. A second convolutional layer of  $(3 \times 3)$  filters yielding 64 channels. This channel also posses  
 1115 a bias parameter for a total of  $(64 \cdot (3 \times 3) + 64) = 640$  parameters. Resulting image is  $(24 \times$   
 1116  $24 \times 64)$ . This channel also apply a ReLU activation function.
- 1117 3. Then comes a  $(2 \times 2)$  max pool layer with a stride of 1 meaning that for each channel the max  
 1118 value of pixels in a  $(2 \times 2)$  block is condensed in a single resulting pixel. The resulting image  
 1119 is  $(12 \times 12 \times 64)$ .
- 1120 4. This image goes through a dropout layer which will set the pixel to 0 with a probability of  
 1121 0.25. This help prevent overtraining the neural network (see section 3.2.3 for more details).
- 1122 5. The data is the flattened i.e. condensed into a vector of  $(12 \times 12 \times 64) = 9216$  values.
- 1123 6. Then comes a fully connected linear layer (Eq. 3.2) with a ReLU activation that output 128  
 1124 feature. It needs  $(9216 \cdot 128) + 128 = 1'179'776$  parameters.
- 1125 7. This 128 item vector goes through another dropout layer with a probability of 0.5
- 1126 8. The vector is then transformed through a linear layer with ReLU activation. It output 10  
 1127 values, one for each digit class  $(0, 1, 2, \dots, 9)$ . It need  $(128 \cdot 10) + 128 = 1408$  parameters.
- 1128 9. Finally the 10 values are normalized using a log softmax function  $\text{LogSoftmax}(x_i) = \log \left( \frac{\exp(x_i)}{\sum_j \exp(x_j)} \right)$ .
- 1129 Each of those values are the probability of the input image to be a certain digit.



(A) Example of images in the MNIST dataset



(B) Schema of the CNN used in Pytorch example to process the MNIST dataset

FIGURE 3.11

1130 The final network needs 1'182'144 parameters or, if we consider each parameters to be a double  
 1131 precision floating point, 9.45 MB of data. To gives a order of magnitude, such neural network is  
 1132 considered "simple", train in a matter of minutes on T4 GPU [60] (14 epochs) and reach an accuracy  
 1133 in its prediction of 99%.

### 1134 3.3.3 Graph Neural Network (GNN)

1135 As seen in the previous section, the CNNs are powerful for image processing, and more generally  
 1136 any data that can be expressed as a regular, discrete space and from which the information reside  
 1137 in the dispersion in this space. For an image, the edges of an object and how they assemble. A red  
 1138 square, straight edges with a sharp angle between them, is much less representative of a duck than  
 1139 an yellow sphere, round edges without sharp angles.

1140 This "image" projection is not fitted for every problematics. The signals produced by a detector does  
 1141 not always have the properties of images. In the case of JUNO for example, we can create an image  
 1142 of two channels, one for the charge  $Q$  and one for the timing  $t$  but this image should be spheric.  
 1143 Furthermore JUNO is by nature inhomogeneous, using two different systems : The LPMT and the

SPMT. Those two systems have different regime, and thus should be processed differently. We could imagine images with four channels, two for the LPMT and two for the SPMT, or even a branched CNN with one convolution branch for the LPMT and another one for the SPMT. Anyway, the CNN will need to combine the two systems.

To get around the restrictions of data representation imposed by CNNs, we can use the more flexible *graph* representation. A graph  $G(\mathcal{N}, \mathcal{E})$  is composed of vertex or node  $n \in \mathcal{N}$  and edges  $e \in \mathcal{E}$ . The

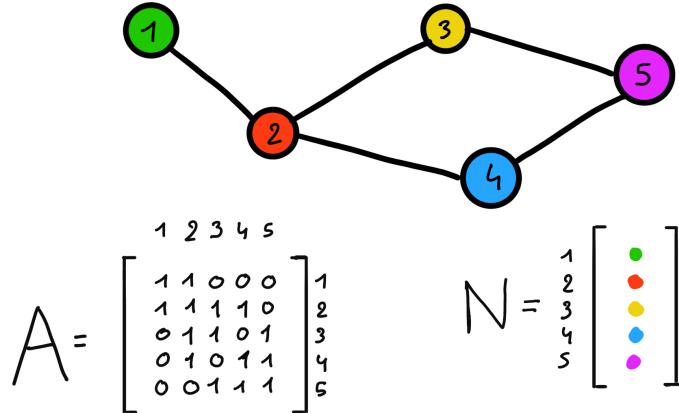


FIGURE 3.12 – Illustration of a graph and its tensor representation.

edges are associated to two nodes  $(u, v) \in \mathcal{N}^2$ , “connecting” them. The node and the edges can hold features, commonly represented as vector  $n \in \mathbb{R}^{k_n}$ ,  $e \in \mathbb{R}^{k_e}$  with  $k_n$  and  $k_e$  the number of features on the nodes and edges respectively. We can thus define a graph using two tensors  $A_{ij}^{ij}$  the adjacency tensor that hold the features  $e \in [0, k_e]$  of the edge connecting the node  $i$  and  $j$  and the tensor  $N_v^i$  that hold the features  $v \in [0, k_n]$  of a node  $i$ .

More figuratively, using the example in figure 3.12, we have a graph of 5 nodes with a color as feature. The edges have no features, we thus encode their existences as 0 or 1. In a realistic examples as JUNO we could represent each PMTs as nodes and the edges between them as their relation such as distance, timing difference, etc... There no strict rules about what is a node or how they should be linked together. This abstraction allow us to represent virtually any type of detector of any geometry.

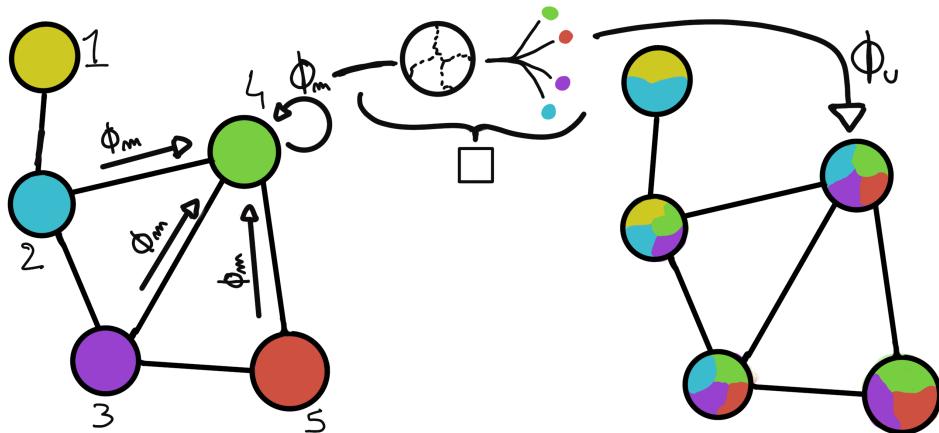


FIGURE 3.13 – Illustration of the message passing algorithm. The detailed explanation can be found in section 3.3.3

1160 To process such object we need specific machine learning algorithms we call Graph neural network.  
 1161 To efficiently manipulate graph we need to structurally encode their property in the neural network  
 1162 computing architecture: each node is equivalent (as opposite to ordered data in a vector), each node  
 1163 has a set of neighbours, ... One of this method is the message passing algorithm presented historically  
 1164 in "Neural Message Passing for Quantum Chemistry" [61]. In this algorithm, with each layer of  
 1165 message passing a new set of features is computed for each node following

$$n_i^{k+1} = \phi_u(n_i^k, \square_j \phi_m(n_i^k, n_j^k, e_{ij}^k)); n_j \in \mathcal{N}'_i \quad (3.14)$$

1166 where  $\phi_u$  is a differentiable *update* function,  $\square_j$  is a differentiable *aggregation* function and  $\phi_m$  is a  
 1167 differentiable *message* function.  $\mathcal{N}'_i = \{n_j \in \mathcal{N} | (n_i, n_j) \in \mathcal{E}\}$  is the set of neighbours of  $n_i$ , i.e. the  
 1168 nodes  $n_j$  from which it exist an edge  $e_{i,j} \rightarrow (n_i, n_j)$ .  $k$  is the layer on which the message passing  
 1169 algorithm is applied. The update function need also a few other property if we want to keep the  
 1170 graph property, most notably the permutational invariance of its parameters (example: mean, std,  
 1171 sum, ...).

1172 The edges features can also be updated, either by directly taking the results of  $\phi_m$  or by using another  
 1173 message function  $\phi_e$ .

1174 To explain this process, let's take the situation presented in figure 3.13. We start with an input graph  
 1175 on left, in this case the message passing algorithm is mixing the color on each nodes and produce  
 1176 nodes of mixed color. For simplicity, the  $\phi_m$  and  $\phi_u$  function are the identity, they take a color and  
 1177 output the same color.

1178 Let's look at what's happening in the node 4. It has 3 neighbours and is a neighbour of itself. The four  
 1179 resulting  $\phi_m$  extract the color of each nodes and then feed them to the  $\square$  function. The  $\square$  function  
 1180 just equally distribute the color in the node. Finally the  $\phi_u$  function just update the node with the  
 1181 output of  $\square$ .

1182 Interestingly we see that the new node 4 does not have any yellow, the color of node 1. But if we were  
 1183 to run the message passing algorithm again, it would get some as node 2 is now partially yellow. If  
 1184 color here represent information, we see that multiple step are needed so that each node is "aware"  
 1185 of the informations the other nodes possess.

1186 Message passing is a very generic way of describing the process of GNN and it can be specialized  
 1187 for convolutional filtering [49], diffusion [62] and many other specific operation. GNN are used in a  
 1188 wide variety of application such as regression problematics, node classification, edge classification,  
 1189 node and edge prediction, ...

1190 It is a very versatile but complex tool.

### 1191 3.3.4 Adversarial Neural Network (ANN)

1192 The adversarial machine learning, Adversarial Neural Networks (ANN) in the case of neural net-  
 1193 work, is a family of unsupervised machine learning algorithms where the learning algorithm (gen-  
 1194 erator) is competing against another algorithm (discriminator). Taking the example of Generative  
 1195 Adversarial Networks, concept initially developed by Goodfellow et al. [63], the discriminator goal  
 1196 is to discriminate between data coming from a reference dataset and data produced by the generator.  
 1197 The generator goal, on the other hand, is to produce data that the discriminator would not be able to  
 1198 differentiate from data from the reference dataset. The expression of duality between the two models  
 1199 is represented in the loss where, at least a part of it, is driven by the results of the discriminator.

1200 **Chapter 4**

1201 **Image recognition for IBD  
reconstruction with the SPMT system**

1203

*Dave - Give me the position and momentum, HAL.  
HAL - I'm afraid I can't do that Dave.  
Dave - What's the problem ?  
HAL - I think you know what the problem is just as well as I do.  
Dave - What are you talking about, HAL?  
HAL -  $\sigma_x \sigma_p \geq \frac{\hbar}{2}$*

1204

## Contents

---

|                     |   |                   |
|---------------------|---|-------------------|
| <small>1205</small> | <b>4.1 Motivations</b>                            | <small>56</small> |
| <small>1206</small> | <b>4.2 Method and model</b>                       | <small>56</small> |
| <small>1207</small> | 4.2.1 Model                                       | <small>56</small> |
| <small>1208</small> | 4.2.2 Data representation                         | <small>58</small> |
| <small>1209</small> | 4.2.3 Dataset                                     | <small>59</small> |
| <small>1210</small> | 4.2.4 Data characteristics                        | <small>60</small> |
| <small>1211</small> |   |                   |
| <small>1212</small> | <b>4.3 Training</b>                               | <small>61</small> |
| <small>1213</small> | <b>4.4 Results</b>                                | <small>61</small> |
| <small>1214</small> | 4.4.1 J21 results                                 | <small>63</small> |
| <small>1215</small> | 4.4.2 J21 Combination of classic and ML estimator | <small>66</small> |
| <small>1216</small> | 4.4.3 J23 results                                 | <small>68</small> |
| <small>1217</small> | <b>4.5 Conclusion and prospect</b>                | <small>69</small> |

---

1218 1220 As explained in chapter 2, JUNO is an experiment composed of two systems, the Large Photomultiplier (LPMT) system and the Small Photomultiplier (SPMT) system. Both of them observe the same physics events inside of the same medium but they differ in their photo-coverage, respectively 75.2% and 2.7%, their dynamic range (see section 2.2.2), a thousands versus a few dozen, and their front-end electronics (see section 2.2.2).

1221 They are complementary in their strengths and weaknesses and support each other, this is what we call *Dual Calorimetry*. One important point is their differences in expected resolution, the LPMT system outperform largely the SPMT system but is subject to effects such as charge non linearity [29] that could bias the reconstruction. Effects that the SPMT system is impervious to. This topic will be studied in more detail in chapter 7. Also, due to the dynamic range of the LPMT, in case of high energy and high density event such as core-collapse supernova, the LPMT system could saturate and the lower photo-coverage become a benefit.

1233 Thus, although event reconstruction algorithm and physics analysis combines both LPMT and SPMT systems, individual approach are key studies to understand the detector and ensure their reliability.

1234

1235 This topic will also be studied in more details in chapter 7. The subject of this chapter is to propose  
 1236 a machine learning algorithm for the SPMT reconstruction based on Convolutional Neural Network  
 1237 (CNN).

## 1238 4.1 Motivations

1239 As explained in chapter 3, Machine Learning (ML) algorithms shine when modeling highly dimen-  
 1240 sional data from a given dataset. In our case, we have access to complete monte-carlo simulation of  
 1241 our detector to produce arbitrary large datasets that could represent multiple years of data taking.  
 1242 Ideally ML algorithms would be able to consider the entirety of the information in the detector and  
 1243 converge on the best parameters to yield optimal results, while classical methods could be biased by  
 1244 the prior knowledge of the detector and physics processes. To study this potential phenomena, we  
 1245 will compare our machine algorithm to a classical reconstruction method developed for energy and  
 1246 vertex reconstruction [64].

1247 We have access to a very detailed simulation of the detector (section 2.5) that will allow us to simulate  
 1248 arbitrary large dataset while giving access to the all the physics parameters of the event. Those  
 1249 parameters include the target of our reconstruction algorithms: the vertex and energy of our event.  
 1250 As introduced above, we hope that the ML algorithm will be able to used all the informations in the  
 1251 event, but that could lead that potential mismodelings in our simulation could be exploited by the  
 1252 algorithm. This specific subject will be studied in chapter 6.

## 1253 4.2 Method and model

1254 One of simplest way to look at JUNO data is to consider the detector as an array of geometrically  
 1255 distributed sensors on a sphere. Their repartition is almost homogeneous, on this sphere surface  
 1256 providing an almost equal amount of information per unit surface on this sphere. It is then tempting  
 1257 to represent the detector as a spherical image with the PMTs in place of pixels. Two events with two  
 1258 different energy or position would produce two different images.

1259 The most common approach in machine learning for image processing and image recognition is the  
 1260 Convolutional Neural Network (CNN). It is widely used in research and industry [56, 65–67] due to  
 1261 its strengths (see section 3.3.2) and has proven its relevance in image processing.

1262 Some CNN are developed to process spherical images [68] but for the sake of simplicity and as a  
 1263 first approach we decided to go with a planar projection of the detector, approach that has proven its  
 1264 efficiency using the LPMT system (see section 2.6.3). The details about this planar projection will be  
 1265 discussed in section 4.2.2.

### 1266 4.2.1 Model

1267 The architecture we use is derived from the VGG-16 architecture [56] illustrated in figure 4.1. We  
 1268 define a set of hyperparameters that will define the size, complexity and computational power of the  
 1269 NN. The chose hyperparameters are detailed below and their values are presented in table 4.1.

- 1270 —  $N_{\text{blocks}}$ : the number of convolution blocks, a block being composed of two convolutional  
 1271 layers with  $3 \times 3$  filters using ReLU activation function, a  $3 \times 3$  max-pooling layer (except for  
 1272 the last block).
- 1273 —  $N_{\text{channels}}$ : The number of channels in the first block. The number of channels in the subsequent  
 1274 blocks is computed using  $N_{\text{channels}}^i = i * N_{\text{channels}}, i \in [1..N_{\text{blocks}}]$ .

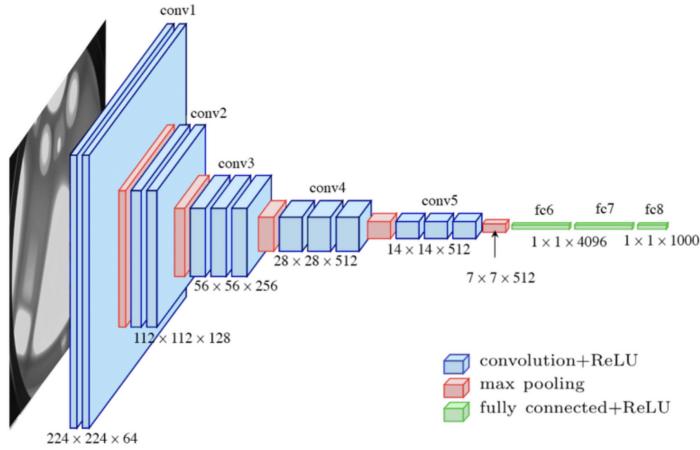


FIGURE 4.1 – Graphic representation of the VGG-16 architecture, presenting the different kind of layer composing the architecture.

- **FCDNN configuration:** The result of the last convolution layer is flattened then fed to a FCDNN. Its configuration is expressed as a sequence of fully connected linear layer using the PReLU activation function. For example  $2 * 1024 + 2 * 512$  is the sequence of 2 layers with a width of 1024 followed by 2 other layers with a width of 512. Finally the last layer is a 4 neurons wide linear layers without activation function. Each neurons of the last layer represent a component of the interaction vertex: Energy, X, Y, Z.
- **Loss:** The loss function. In this work we study two different loss function  $(E + V)$  and  $(E_r + V_r)$  detailed below.

$$(E + V)(E, x, y, z) = \left\langle (E - E_{true})^2 + 0.85 \sum_{\lambda \in [x, y, z]} (\lambda - \lambda_{true})^2 \right\rangle \quad (4.1)$$

$$(E_r + V_r)(E, x, y, z) = \left\langle \frac{(E - E_{true})^2}{E_{true}} + \frac{10}{R} \sum_{\lambda \in [x, y, z]} (\lambda - \lambda_{true})^2 \right\rangle \quad (4.2)$$

where  $R$  is the radius of the CD. With the energy in MeV and the distance in meters, we use the factor 0.85 and 10 to equilibrate the two term of the loss function so they have the same magnitude.

- The loss function  $(E + V)$  is close to a simple Mean Squared Error (MSE). MSE is one of the most basic loss function, the derivative is simple and continuous in every point. It is a strong starting point to explore the possibility of CNNs.
- $(E_r + V_r)$  can be seen as a relative MSE.

The idea is that: due to the inherent statistic uncertainty over the number of collected Number of Photo Electrons (NPE), the absolute resolution  $\sigma(E - E_{true})$  will be larger at higher energy than at low energy. But we expect the *relative* energy resolution  $\frac{\sigma(E - E_{true})}{E_{true}}$  to be smaller at high energy than lower energy as illustrated in figure 2.22. Because of this, by using simple MSE the most important part in the loss come from the high energy part of the dataset whereas with a relative MSE, the most important part become the low energy events in the dataset. We hope that by using a relative MSE, the neural network will focus on low energy events where the reconstruction is considered the hardest.

Each combination of those hyperparameters (for example  $(N_{blocks} = 2, N_{channels} = 32, \text{FCDNN} = (2 * 1024), \text{Loss} = (E + V))$ ), subsequently designated as configurations, is then tested and compared to each other over an analysis sample.

On top those generated models, we define 4 hand tailored models:

- 1301 — “gen\_0”:  $N_{blocks} = 4$ ,  $N_{channels} = 64$ , FCDNN configuration:  $1024 * 2 + 512 * 2$ , Loss :=  $E + V$   
 1302 — “gen\_1”:  $N_{blocks} = 4$ ,  $N_{channels} = 64$ , FCDNN configuration:  $1024 * 2 + 512 * 2$ , Loss :=  $E_r + V_r$   
 1303 — “gen\_2”:  $N_{blocks} = 5$ ,  $N_{channels} = 64$ , FCDNN configuration:  $4096 * 2 + 1024 * 2$ , Loss :=  $E + V$   
 1304 — “gen\_3”:  $N_{blocks} = 5$ ,  $N_{channels} = 64$ , FCDNN configuration:  $4096 * 2 + 1024 * 2$ , Loss :=  $E_r + V_r$

|                      |   |
|----------------------|---|
| $N_{blocks}$         | {2, 3, 4}   |
| $N_{channels}$       | {32, 64, 128}   |
| FCDNN configurations | $2 * 1024$<br>$2 * 2048 + 2 * 1024$<br>$3 * 2048 + 3 * 512$<br>$2 * 4096$ |
| Loss                 | { $E + V$ , $E_r + V_r$ }   |

TABLE 4.1 – Sets of hyperparameters values considered in this study

1305 We cannot use the mean loss because we consider multiple loss functions, there is no guarantee that  
 1306 comparison of their numerical value will be meaningful. We use multiple observables to rank the  
 1307 performances of each configuration:

- 1308 — The mean absolute energy error  $\langle E \rangle = \langle |E - E_{true}| \rangle$ . It is an indicator of the energy bias of our  
 1309 reconstruction.  
 1310 — The standard deviation of the energy error  $\sigma E = \sigma(E - E_{true})$ . This the indicator on our  
 1311 precision in energy reconstruction.  
 1312 — The mean distance between the reconstructed vertex and the true vertex  $\langle V \rangle = \langle |\vec{V} - \vec{V}_{true}| \rangle$ .  
 1313 This an indicator of the bias and precision of our vertex reconstruction.  
 1314 — The standard deviation of the distance between the true and reconstructed vertex  $\sigma V = \sigma |\vec{V} -$   
 1315  $\vec{V}_{true}|$ . This is an indicator if the precision in our vertex reconstruction.

1316 The models were developped in Python using the pytorch framework [58] using NVIDIA A100 [69]  
 1317 and NVIDIA V100 [70] gpus. The A100 was split in two, thus the accessible gpu memory was 20 Gb  
 1318 making it impossible to train some of the architectures due to memory consumption.

1319 The training was monitored in realtime by a custom tooling that was developed during this thesis,  
 1320 DataMo [71].

1321 The training of one model takes between 4h and 15h depending of its size, overall training the full  
 1322 72 model takes around 500 GPU hours. Even with parallel training, this random search hyper-  
 1323 optimisation was time consuming.

## 1324 4.2.2 Data representation

1325 This data is represented as  $240 \times 240$  images with a charge  $Q$  channel and a time  $t$  channel. The  
 1326 SPMTs are then projected on the plane as illustrated in figure 4.2. The  $x$  position is proportional to  $\theta$   
 1327 and the  $y$  position is defined by  $\phi \sin \theta$  in spherical coordinates.  $\theta = 0$  is defined as being the top of  
 1328 the detector and  $\phi = 0$  is defined as an arbitrary direction in the detector. In practice,  $\phi = 0$  is given  
 1329 by the MC simulation.

$$x = \left\lfloor \frac{\theta \cdot H}{\pi} \right\rfloor, \theta \in [0, \pi] \quad (4.3)$$

$$y = \left\lfloor \frac{(\phi + \pi) \sin \theta \cdot W}{2\pi} \right\rfloor, \phi \in [-\pi, \pi], \theta \in [0, \pi] \quad (4.4)$$

1330 where  $H$  is the height of the image,  $W$  the width of the image and  $(0, 0)$  the top left corner of the  
 1331 image.

When two SPMTs are in the same pixel, the charges are summed and the lowest of the hit-time is chosen. The SPMTs being located close to each other, we expect the time difference between two successive physics signals, two photons being collected, to be small. The first hit time is chosen because it can be considered as the relative propagation time of the photons that went the "straightest", i.e. that went under the less perturbation of the two. The only potential problem in using this first time come from the Dark Noise (DN). Its time distribution is uniform over the signal and could come before a physics signal on the other SPMT in the pixel. In that case, the time information in the pixel become irrelevant and we lose the timing information for this part of the detector. As illustrated in figure 4.2 the image dimension have been optimized so that at most two SPMTs are in the same pixel while keeping the number of empty pixels relatively low to prevent this kind of issue.

While it could be possible to use larger images (more pixel) to prevent overlapping, keeping image small images gives multiple advantages:

- As presented in section 4.2.1, the convolution filter we use are  $3 \times 3$  convolution filter, meaning that if SPMTs would be separated by more than one pixel, the first filter would only see one SPMT per filter. This behavior would be kind of counterproductive as the first convolution block would basically be a transmission layer and would just induce noise in the data.
- It keep the network relatively small, while this do not impact the convolution layers, the flatten operation just before the FCDNN make the number parameters in the first layer of it dependent on the size of the image.
- It reduce the number of empty pixel in the image.

The question of empty pixel is an important question in this data representation. There is two kind of empty pixels in the data.

The first kind is pixel that contain a SPMT but the SPMT did not get hit nor registered any dark noise during the event. In this case, the charge channel is zero, which have a physical meaning but then come the question of the time layer. One could argue that the correct time would be infinity (or the largest number our memory allows us) because the hit "never" happened, so extremely far from the time of the event. This cause numerical problem as large number, in the linear operation that are happening in the convolution layers, are more significant than smaller value. We could try to encode this feature in another way but no number have any significance due to our time being relative to the trigger of the experiment so -1 for example is out of question. Float and Double gives us access to special value such as NaN (Not a Number) [72] but the behavior is to propagate the NaN which leaves us with NaN for energy and position. We choose to keep the value 0 because it's the absorbing element of multiplication, absorbing the "information" of the parameter it would be multiplied by. It also can be though as no activation in the ReLU activation function.

The second kind of pixel is pixel that do not represent parts of the detector such as the corners of the image. The question is basically the same, what to put in the charge and the time channel. The decision is to set the charge and time to 0 following the above reasoning. It's important to keep in mind the fact that a part of the detector that has not been hit is also an information: There is no signal in this part of the detector. This problematic will be explored in more details in chapter 5.

Another problematic that happens with this representation, and this is not dependent of the chosen projection, is the deformation in the edges of the image and the loss of the neighbouring information in the for the SPMTs at the edge of the image  $\phi \sim 180^\circ$ . This deformation and neighbouring loss could be partially circumvented as explained in section 4.5

### 4.2.3 Dataset

In this study we will discuss two datasets of one millions events:

- **J21**: The first one comes from the JUNO official mc simulation J21v1r0-Pre2 (released the 18th August 2021). This historical version is the one on which the classical algorithm presented in [64] was developed. This dataset is used as a reference for comparison to classical algorithm.

The data in this dataset is *detsim* level (see section 2.5), where only the physic is simulated. The charge and time biases and uncertainties are implemented using toy MC adjusted using [26, 73]. The time window is not based on a selection algorithm but  $t_0 := t = 0$  is defined as the first PMT hit. The window goes up to  $t_0 + 1000$  ns.

- J23: The second comes from the JUNO official monte-carlo simulations J23.0.1-rc8.dc1 (released the 7th January 2024). The data is *calib* level (see section 2.5). Here the charge comes from the waveform integration, the time window resolution and trigger decision are all simulated inside the software. This dataset is more realistic and is used to confirm the performance of our algorithm.

To put in perspective this amount of data, the expected IBD rate in JUNO is 47 / days. Taking into account the calibration time, and the source reactor shutdown, it amount to  $\sim 94'000$  IBD events in 6 years. With this million of event, we are training the equivalent of  $\sim 10$  years of data. With this amount we reach a density of  $4783 \frac{\text{event}}{\text{m}^3 \cdot \text{MeV}}$ , meaning our dataset is representative of the multiple event scenarios that could be happening in the detector.

While we expect and hope the monte-carlo simulation to give use a realistic representation of the detector, there could be effect, even after the fine-tuning on calibration data, that the simulation cannot handle. Thus, once the calibration will be available, we will need to evaluate, and if needed retrain, the network on calibration data to establish definitive performances.

The simulated data is composed of positron events, uniformly distributed in the CD volume and in kinetic energy over  $E_k \in [0; 9]$  MeV producing a deposited energy  $E_{dep} \in [1.022; 10.022]$  MeV. This is done to mimic the signal produced by the IBD prompt signal. Uniform distributions are used so that the CNN does not learn a potential energy distribution, favoring some part of the energy spectrum instead of other.

Those events can be considered as “optimistic” as there is no pile-up with potential background or other IBD.

#### 4.2.4 Data characteristics

To delve a bit into the kind of data we will use, you can find in figure 4.2 the repartition of the SPMTs in the image. The color represent the number of SPMTs per pixel.

In figures 4.3, 4.4, 4.5 and 4.6 are presented events from J23 for different positions and energies. We see some characteristics and we can instinctively understand how the CNN could discriminate different situations.

To give an idea of the strength of the signal in comparison to the dark noise background, figure 4.7a present the distribution of the ratio of NPE per deposited energy. Assuming a linear response of the LS we can model:

$$NPE_{tot} = E_{dep} \cdot P_{mev} + D_N \quad (4.5)$$

$$\frac{NPE_{tot}}{E_{dep}} = P_{mev} + \frac{D_N}{E_{dep}} \quad (4.6)$$

where  $NPE_{tot}$  is the total number of PE detected by the event,  $P_{mev}$  is the mean number of PE detected per MeV and  $D_N$  is the dark noise contribution that is considered energy independent. In the case where the readout time window is dependent of the energy the dark noise contribution become energy dependant, also the LS response is realistically energy dependant but figure 4.7a shows that we have heavily dominated by statistical uncertainties which is why we are using this simple model.

The fit shows a light yield of 40.78 PE/MeV and a dark noise contribution of 4.29 NPE. As shown in figure 4.7b, the physics makes for 90% of the signal at low energy.

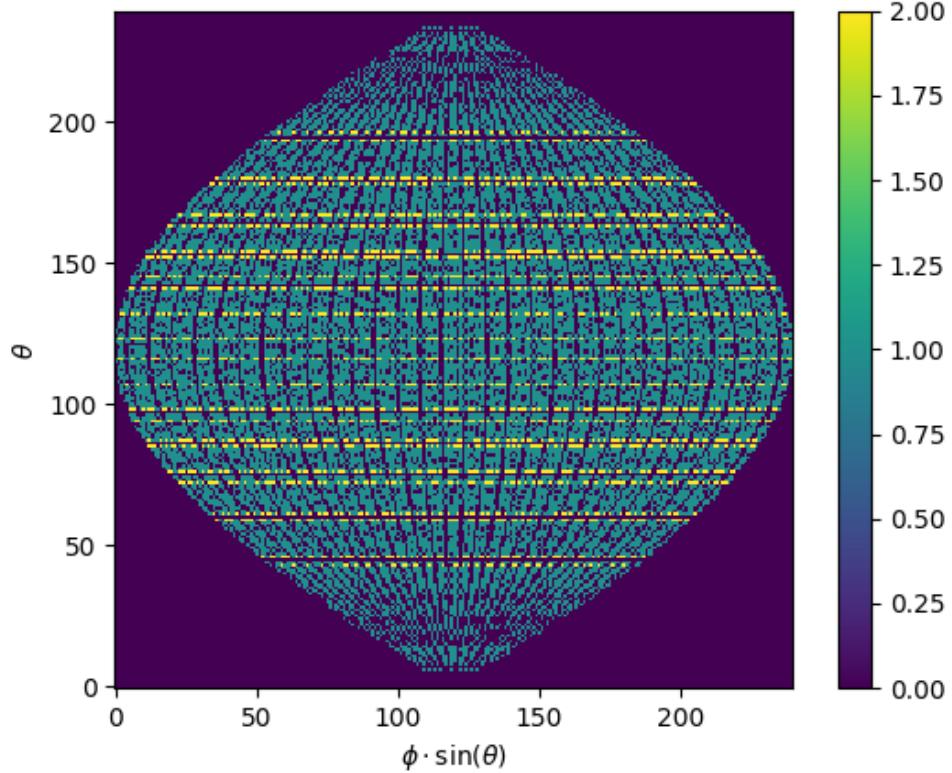


FIGURE 4.2 – Repartition of SPMTs in the image projection. The color scale is the number of SPMTs per pixel

### 4.3 Training

The optimizer used for the training is the Adam [53] optimizer, with a learning rate  $\lambda$  of  $1e-3$ . The other hyperparameters were left to their default value ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 1e^{-8}$ ). The learning rate was reduced exponentially during the training at a rate of  $\gamma = 0.95$ , thus  $\lambda_{i+1} = 0.95\lambda_i$  where  $i$  is the epoch.

The training was composed of 30 epochs, each epoch constituted of 10k steps using a batch size of 64 events. The validation was computed over a 100 steps on the validation dataset.

### 4.4 Results

Before presenting the results, let's discuss the different observables.

The event are considered point like in this study. The target truth position, or vertex, is the mean position of the energy deposits of the positron and the two annihilation gammas. Due to the symmetries of the detector, we mainly consider and discuss the bias and precision evolution depending of the radius  $R$  but we will still monitor the performances depending of the spheric angle  $\theta$  and  $\phi$ . From the detector construction and effect we expect dependency in radius due to the TR area effect presented in section 2.6 and the possibility for the positron or the gammas to escape from the CD for near the

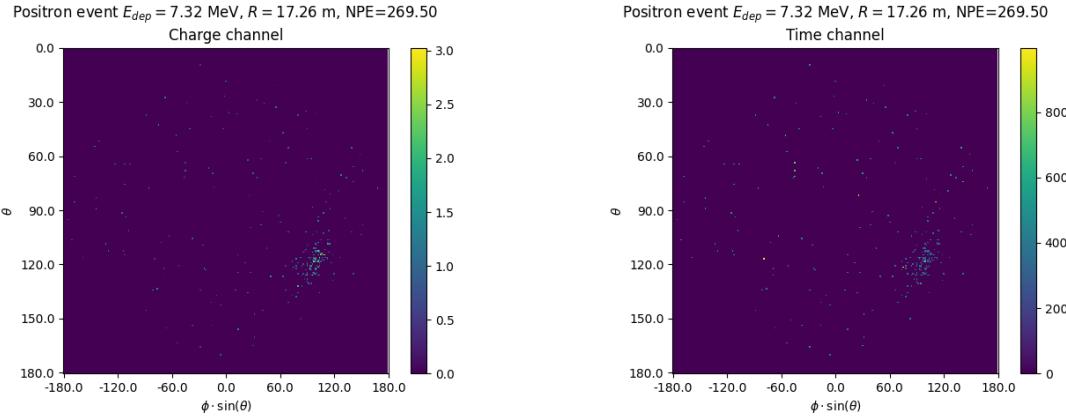


FIGURE 4.3 – Example of a high energy, radial event. We see a concentration of the charge on the bottom right of the image, clear indication of a high radius event. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

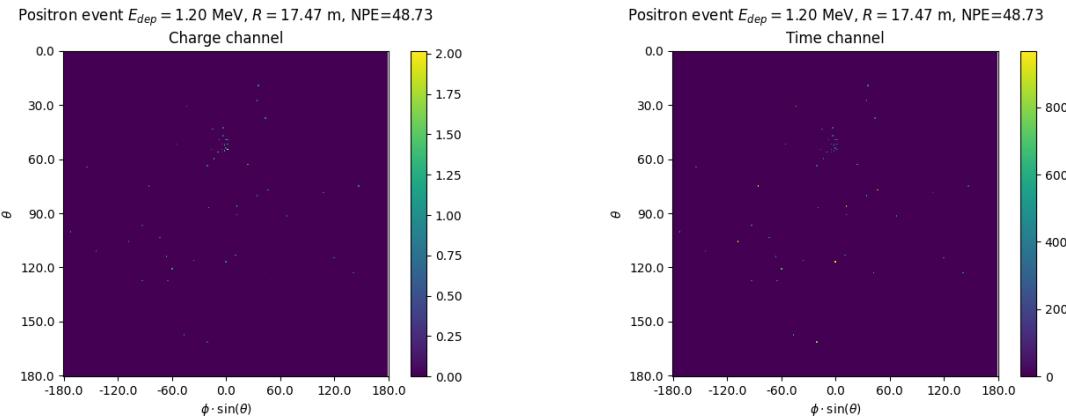


FIGURE 4.4 – Example of a low energy, radial event. The signal here is way less explicit, we can kind of guess that the event is located in the top middle of the image. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

1433 edge events. We also expect dependency in  $\theta$ , the top of the experiment being non-instrumented due  
 1434 to the filling chimney. It is also to be noted that the events in the dataset are uniformly distributed in  
 1435 the CD, and so are uniformly distributed in  $R^3$  and  $\phi$ . The  $\theta$  distribution is not uniform and we will  
 1436 have more event for  $\theta \sim 90^\circ$  than  $\theta \sim 0^\circ$  or  $\theta \sim 180^\circ$ .

1437 We define multiple energy in JUNO:

- 1438 —  $E_\nu$ : The energy of the neutrino.
- 1439 —  $E_k$ : The kinetic energy of the resulting positron from the IBD.
- 1440 —  $E_{dep}$ : The deposited energy of the positron and the two annihilation gammas.
- 1441 —  $E_{vis}$ : The equivalent visible energy, so  $E_{dep}$  after the detector effect such as the absorption of  
 1442 scintillation photons by the LS and the LS response non-linearity.
- 1443 —  $E_{rec}$ : The reconstructed energy by the reconstruction algorithm. The expected value depend  
 1444 on the algorithm we discuss about. For example the algorithm presented in section 2.6 is  
 1445 reconstructing  $E_{vis}$  while the ones presented in section 2.6.3 reconstruct  $E_{dep}$ .

1446 In this study, we will set  $E_{dep}$  as our target for energy reconstruction. This choice is motivated by the  
 1447 ease with which we can retrieve this information in the monte-carlo data while  $E_{vis}$  is less trivial to

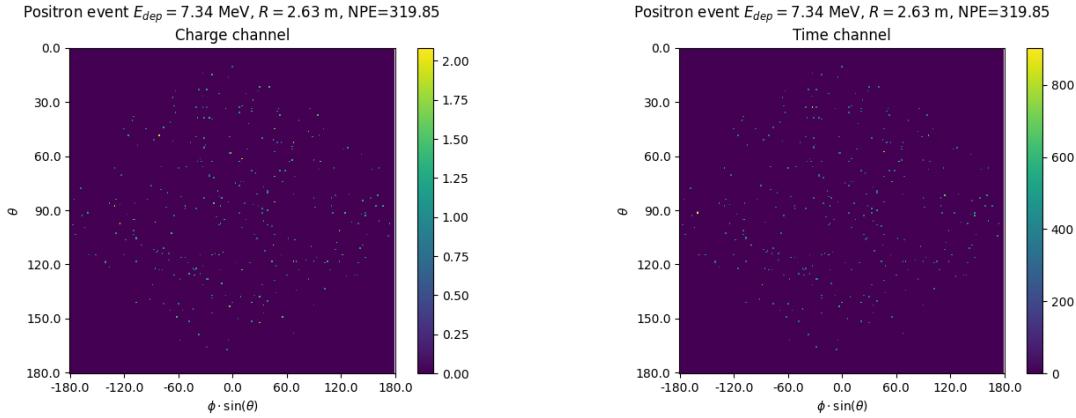


FIGURE 4.5 – Example of a high energy, central event. In this image we can see a lot of signal but uniformly spread, this is indicative of a central event. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

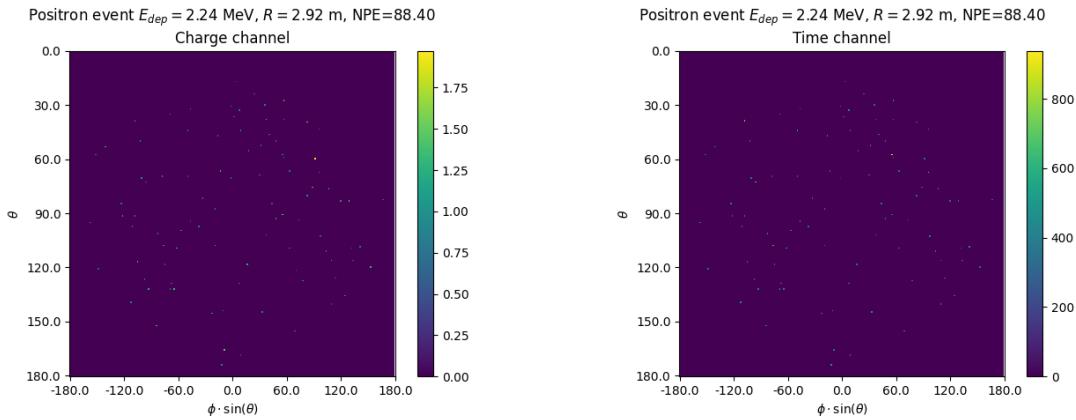
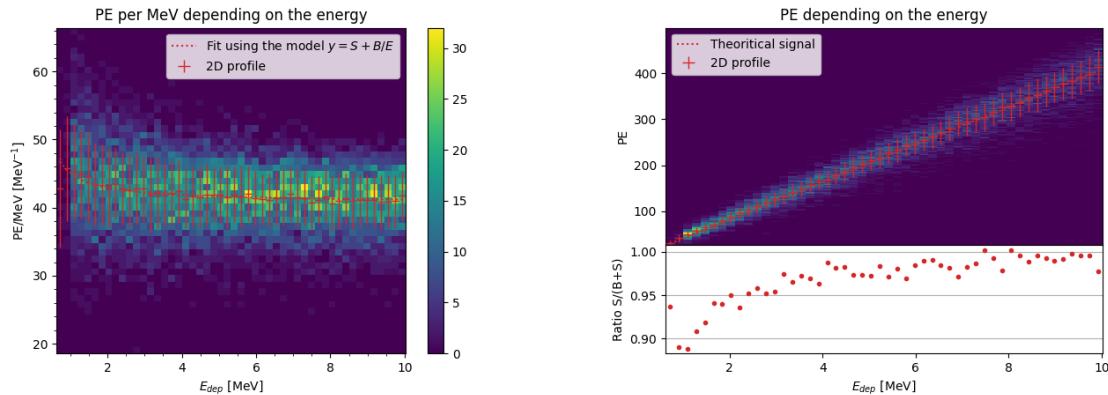


FIGURE 4.6 – Example of a low energy, central event. Here there is no clear signal, the uniformity of the distribution should make it central. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

1448 retrieve.

#### 1449 4.4.1 J21 results

1450 Those results comes from the “gen\_30” model, meaning then 30th model generated using the table  
1451 4.1 or  
1452 — “gen\_30”:  $N_{blocks} = 3$ ,  $N_{channels} = 32$ , FCDNN configuration:  $2048 * 2 + 1024 * 2$ , Loss :=  $E + V$   
1453 The performances of its reconstruction are presented in blue in figure 4.8. Superimposed in black is  
1454 the performances of the classical algorithm from [64].



(A) Distribution of PE/MeV in the J23 Dataset. This distribution is profiled and fitted using equation 4.6

(B) On top: Distribution of PE vs Energy. On bottom: Using the values extracted in 4.7a, we calculate the ration signal over background + signal

FIGURE 4.7

#### 1455 Energy reconstruction

1456 By looking at the figure 4.8a and 4.8b, the CNN has similar performances in its energy resolution.  
 1457 Only at the end of the energy range does the resolution get a little better.

1458 This is explained by looking at the true and reconstructed energy distributions in figure 4.10a. We  
 1459 see that the distributions are similar for energies before 8 MeV but there is an excess of event recon-  
 1460 structed with energies around 9 MeV while a lack of them for 10 MeV. The neural network seems to  
 1461 learn the energy distribution and learn that it exist almost no event with an energy inferior to 1.022  
 1462 MeV and not event with an energy superior to 10 MeV.

1463 The first observation is a physics phenomena: for a positron, its minimum deposited energy is the  
 1464 mass energy coming from its annihilation with an electron 1.022 MeV. There is a few event with  
 1465 energies inferior to 1.022 MeV, in those case the annihilation gammas or even the positron escape the  
 1466 detector. The deposited energy in the LS is thus only a fraction of the energy of the event.

1467 The second observation is indeed true in this dataset but has no physical meaning, it is an arbitrary  
 1468 limit because the physics region of interest is mainly between 1 and 9 MeV of deposited energy  
 1469 (figure 2.2). By learning the energy distribution, the CNN pull event from the border of it to more  
 1470 central value. That's why the energy resolution is better: the events are pulled in a small energy  
 1471 region , thus a small variance but the bias become very high (figure 4.8a).

1472 This behavior also explain the heavy bias at low energy in figure 4.8a. The energy bias of the CNN if  
 1473 fairly constant over the energy range, it is interesting to note that the energy bias depending on the  
 1474 radius is a bit worse than the classical method.

#### 1475 Vertex reconstruction

1476 For the vertex reconstruction we do not study  $x$ ,  $y$  and  $z$  independently but we use  $R$  as a proxy  
 1477 observable. Figure 4.9 shows the error distribution of the different vertex coordinates. We see that  
 1478  $R$  errors and biases are slightly superior to the cartesian coordinates, thus  $R$  is a conservative proxy  
 1479 observable to discuss the subject of vertex reconstruction.

1480 The comparison of radius reconstruction between the classical algorithm and “gen\_30” are presented  
 1481 in the figures 4.8c, 4.8d, 4.8e and 4.8f.

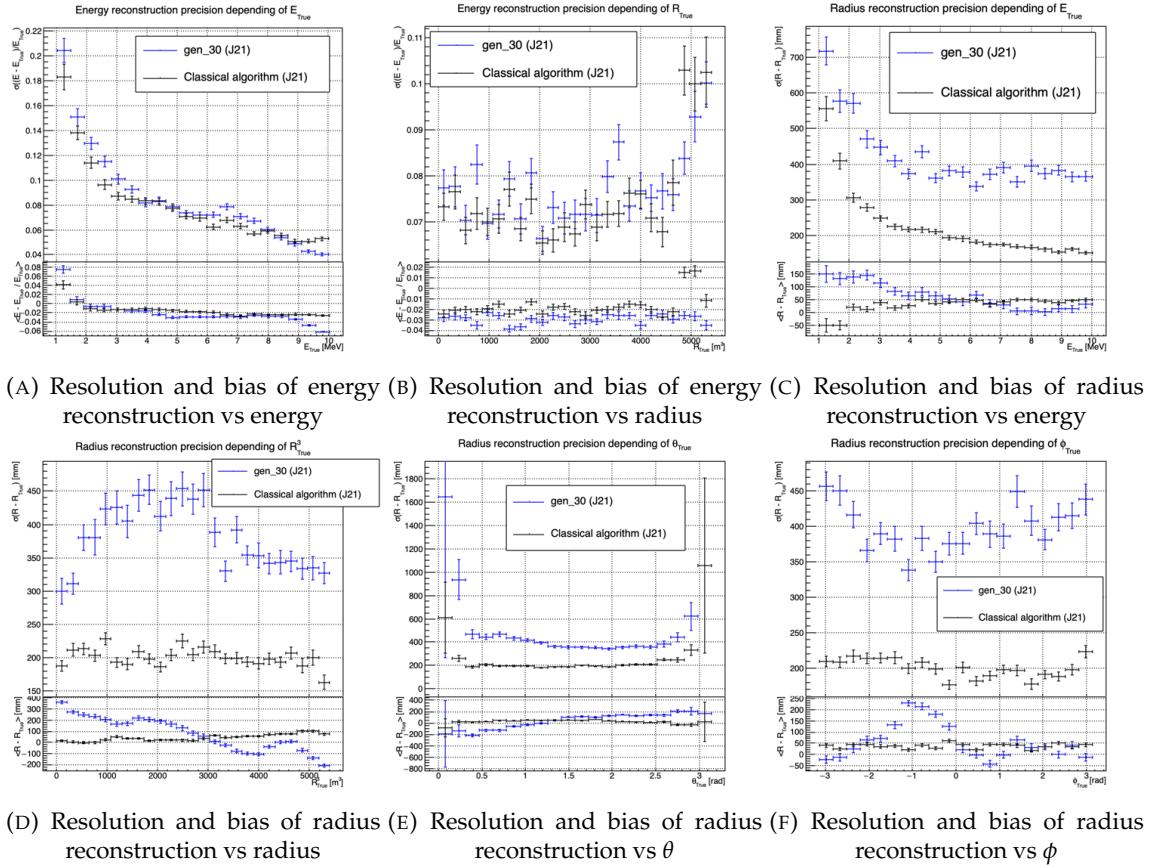


FIGURE 4.8 – Reconstruction performance of the “gen\_30” model on J21 data and its comparison to the performances of the classic algorithm “Classical algorithm” from [64]. The top part of each plot is the resolution and the bottom part is the bias.

Radius reconstruction is worse than the classical algorithms in all configuration. In energy, figure 4.8c, where we see a degradation of almost 20cm over the energy range.

When looking over the true event radius, figure 4.8d, we lose between 30 and 45cm of resolution. The performances are the best for central and radial event.

The precision also worsen when looking at the edge of the image  $\theta \approx 0, \theta \approx 2\pi$  respectively the top and bottom of the image, and when  $\phi \approx -\pi$  and  $\phi \approx \pi$  respectively the left and right side of the image. This is the confirmation that the deformation of the image is problematic for the event reconstruction.

The bias in radius reconstruction is about the same order of magnitude depending of the energy but is of opposite sign. As for the energy, this behavior is studied in more details in section 4.4.2. Over radius,  $\theta$  and  $\phi$  the bias is inconsistent, sometimes event better than the classical reconstruction but can also be much worse than the classical method. This could come from the specialisation of some filters in the convolutional layers for specific part of the detector that would still work “correctly” for other parts but with much less precision.

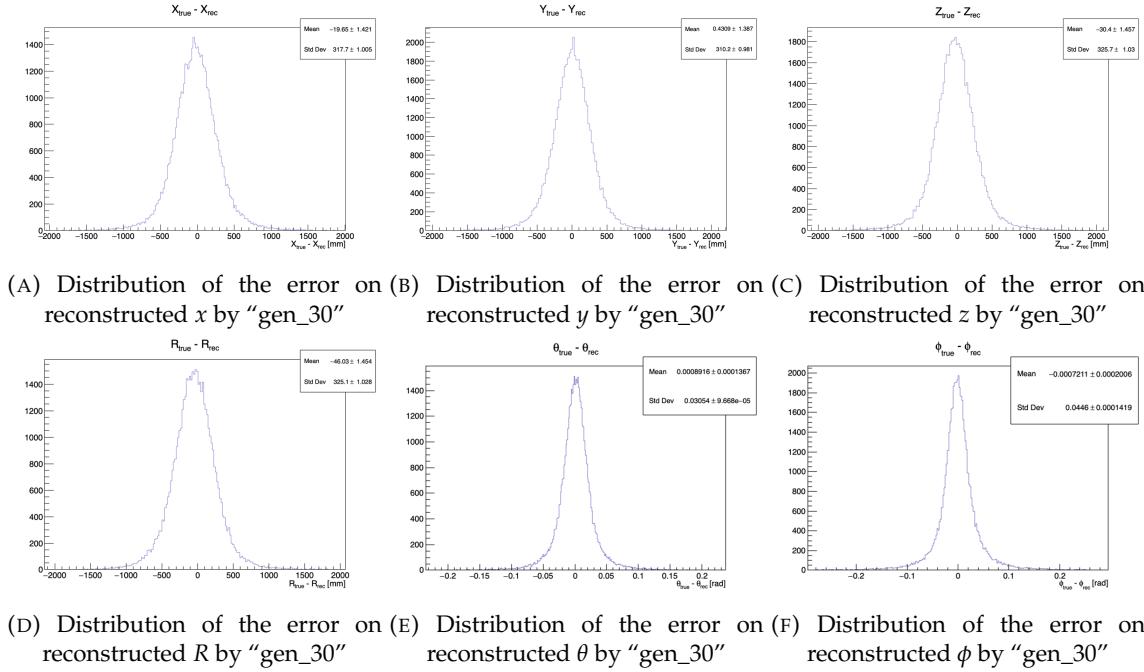


FIGURE 4.9 – Error distribution of the different component of the vertex by "gen\_30". The reconstructed component are  $x$ ,  $y$  and  $z$  but we see similar behavior in the error of  $R$ ,  $\theta$  and  $\phi$ .

#### 1496 4.4.2 J21 Combination of classic and ML estimator

As it has been presented in previous section, there are instances where the reconstructed energy and vertex behaves differently between the neural network and the classic algorithm. For instance, if we look at figure 4.8c, we see that while the CNN tend to overestimate the radius at low energy while the classical algorithm seems to underestimate it. Let's designate the two reconstruction algorithms as estimator of  $X$ , the truth about the event in the phase space  $(E, x, y, z)$ . The CNN and the classical algorithm are respectively designated as  $\theta_N(X)$  and  $\theta_C(X)$ .

$$E[\theta_N] = \mu_N + X; \text{Var}[\theta_N] = \sigma_N^2 \quad (4.7)$$

$$E[\theta_C] = \mu_C + X; \text{Var}[\theta_C] = \sigma_C^2 \quad (4.8)$$

1497 where  $\mu$  is the bias of the estimator and  $\sigma^2$  its variance.

1498 Now if we were to combine the two estimators using a simple mean

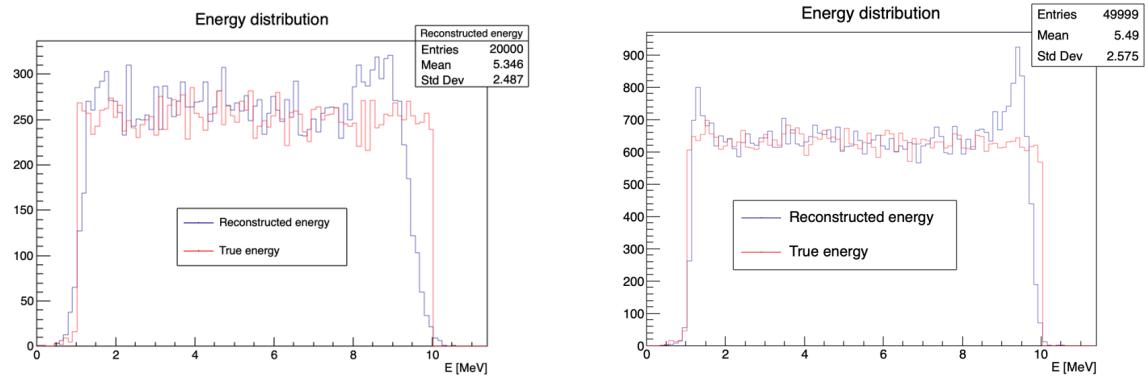
$$\hat{\theta}(X) = \frac{1}{2}(\theta_N(X) + \theta_C(X)) \quad (4.9)$$

then the variance and mean would follow

$$E[\hat{\theta}] = \frac{1}{2}E[\theta_N] + \frac{1}{2}E[\theta_C] \quad (4.10)$$

$$= \frac{1}{2}(\mu_N + X + \mu_C + X) \quad (4.11)$$

$$= \frac{1}{2}(\mu_N + \mu_C) + X \quad (4.12)$$



(A) Distribution of "gen\_30" reconstructed energy and true energy of the analysis dataset (J21)

(B) Distribution of "gen\_42" reconstructed energy and true energy of the analysis dataset (J23)

FIGURE 4.10

$$\text{Var}[\hat{\theta}] = \frac{1}{4}\sigma_N^2 + \frac{1}{4}\sigma_C^2 + 2 \cdot \frac{1}{4} \cdot \sigma_{NC} \quad (4.13)$$

$$= \frac{1}{4}\sigma_N^2 + \frac{1}{4}\sigma_C^2 + \frac{1}{2} \cdot \sigma_{NC} \quad (4.14)$$

$$= \frac{1}{4}\sigma_N^2 + \frac{1}{4}\sigma_C^2 + \frac{1}{2} \cdot \sigma_N\sigma_C\rho_{NC} \quad (4.15)$$

1499 Where  $\sigma_{NC}$  is the covariance between  $\theta_N$  and  $\theta_C$  and  $\rho_{NC}$  their correlation.

1500 We see immediately that if the two estimators are of opposite bias, the bias of the resulting estimator  
1501 is reduced. For the variance, it depends of  $\rho_{NC}$  but in this case if  $\sigma_C^2$  is close to  $\sigma_N^2$  then even for  
1502  $\rho_{NC} \lesssim 1$  then we can gain in resolution.

1503 By generalising the equation 4.9 to

$$\hat{\theta}(X) = \alpha\theta_N + (1 - \alpha)\theta_C; \alpha \in [0, 1] \quad (4.16)$$

1504 we can determine an optimal  $\alpha$  for two combined estimators. The estimators with the smallest  
1505 variance

$$\alpha = \frac{\sigma_C^2 - \sigma_N\sigma_C\rho_{NC}}{\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}} \quad (4.17)$$

1506 and the estimator without bias

$$\alpha = \frac{\mu_C}{\mu_C - \mu_N} \quad (4.18)$$

1507 See annex A for demonstration.

1508 Its pretty clear from the results shown in figure 4.8 that the bias, variances and correlation are not  
1509 constant across the  $(E, R^3)$  phase space. We thus compute those parameters in a grid in  $E$  and  $R^3$  for  
1510 the following results as illustrated in 4.11.

1511 The map we are using are composed of 20 bins for  $R^3$  going from 0 to 5400 m<sup>3</sup> (17.54 m) and 50 bins  
1512 in energy ranging from 1.022 to 10.022 MeV. In the case where we are outside the grid, we use the  
1513 closest cell.

1514 The performance of this weighted mean is presented in figure 4.12. We can see that even when the  
1515 CNN resolution is much worse than the classical algorithm, it can still bring some information thus  
1516 improving the resolution. This comes from the correlation of the reconstruction error to be smaller  
1517 than 1 as presented in figure 4.13. We even see some anticorrelation in the radius reconstruction for

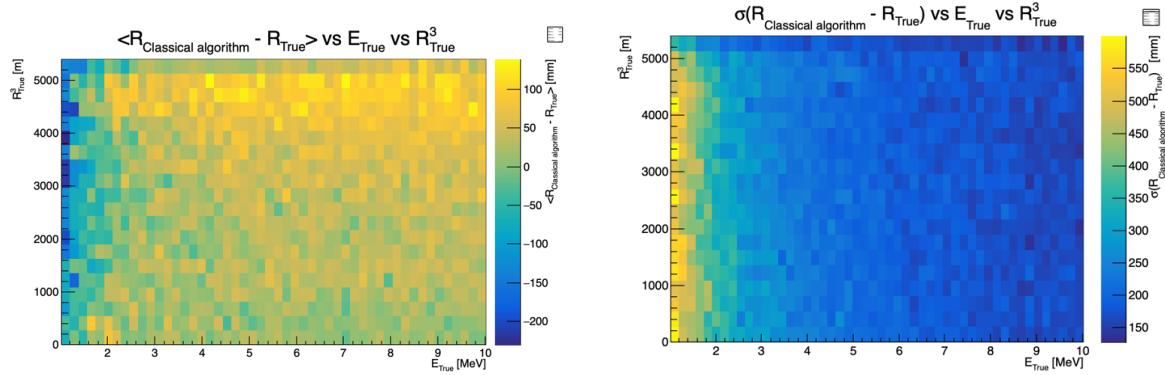


FIGURE 4.11 – Radius bias (on the left) and resolution (on the right) of the classical algorithm in a  $E, R^3$  grid

1518 High radius, high energy, event.

1519 This technique is not suited for realistic reconstruction, we rely too much on the knowledge of  
 1520 the resolution, bias and correlation between the two methods. While this is possible to determine  
 1521 using simulated data or calibration sources, the real data might differ from our model and we  
 1522 would need to really well understand the behavior of the two system. But this is an excellent tool  
 1523 to indicate potential improvements to algorithms and reconstruction methods, showing with this  
 1524 results a potential upper limit to the reconstruction performances.

#### 1525 4.4.3 J23 results

1526 The J21 simulation is fairly old and newer version, such as J23, include refined measurements of the  
 1527 light yield, reflection indices of materials of the detector, structural elements such as the connecting  
 1528 structure and more realistic dark noise. Additionally, the trigger, waveform integration and time  
 1529 window are defined using the algorithms that will ultimately be used by the collaboration to process  
 1530 real physics events.

1531 We retrained the models defined in 4.2.1 on the J23 data and used the same selection procedure. The  
 1532 results from the best architecture, “gen\_42”, are presented in figure 4.14. Following the table 4.1,  
 1533 “gen\_42” is defined as:

1534 — “gen\_42”:  $N_{\text{blocks}} = 3$ ,  $N_{\text{channels}} = 64$ , FCDNN configuration:  $4096 * 2$ , Loss :=  $E + V$

#### 1535 Energy reconstruction

1536 The results of the energy reconstruction are presented in figures 4.14a and 4.14b. Similarly to what  
 1537 we seen for J21, the resolution is close to the one of the classical algorithm with the exception of the  
 1538 start and end of the spectrum. This come from “gen\_42” learning the shape of the distribution and  
 1539 pulling events from the extreme energies, like 1 and 10 MeV, to more common seen energy, like 2 and  
 1540 9 MeV as illustrated in figure 4.10b. The bias disappear with the exception of low and high energy  
 1541 events.

#### 1542 Vertex reconstruction

1543 The vertex reconstruction, presented in figures 4.14c, 4.14d, 4.14e and 4.14f is not yet to the level  
 1544 of the classical reconstruction but the degradation is smaller than for “gen\_32” being at most a

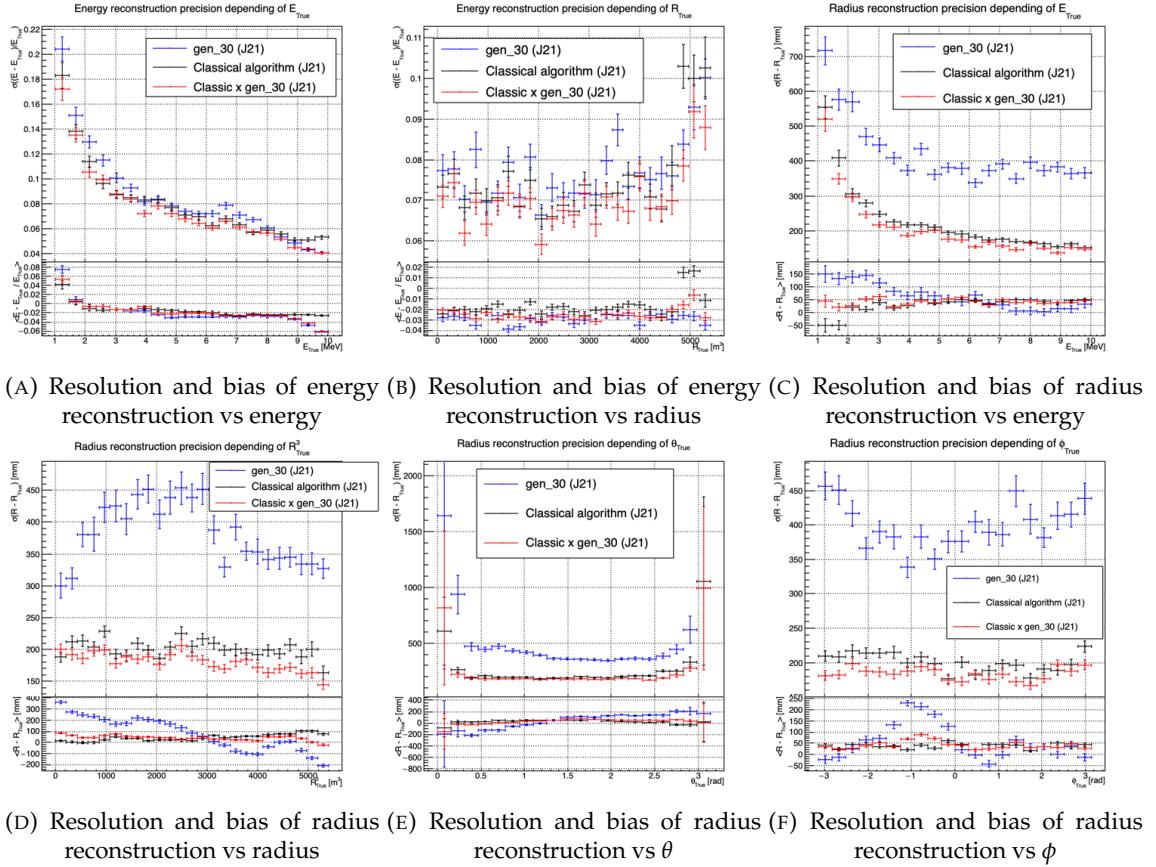


FIGURE 4.12 – Reconstruction performance of the “gen\_30” model on J21, the classic algorithm “Classical algorithm” from [64] and the combination of both using weighted mean. The top part of each plot is the resolution and the bottom part is the bias.

difference of 15cm of resolution and closing to the performance of the classical algorithm in the most favourable condition. “gen\_42” has also very little bias in comparison with the classical method with the exception of the transition to the TR area and at the very edge of the detector.

Unfortunately could not rerun the classical algorithms over the J23 data, as the algorithm was optimised for J21 and was not included and maintained over J23. The combination method need for the two estimators to be run on the same set of event, which was impossible without the classical algorithm being maintained for J23.

Overall the resolution improved over the transition from J21 to J23, effect probably coming from a more complete and rigorous simulation.

## 4.5 Conclusion and prospect

The CNN is a fine tool for event reconstruction in JUNO, and while the reconstruction performances are satisfactory, it show its limitation, the main one concerning the data representation. A lot of training time and resources is consumed going and optimizing over pixel with no physical meaning, the NN needs to optimized itself to take into account edges cases such as event at the edge of the image and deformation of the charge distribution.

Those problems could be circumvented, we could imagine a two part CNN where the first part

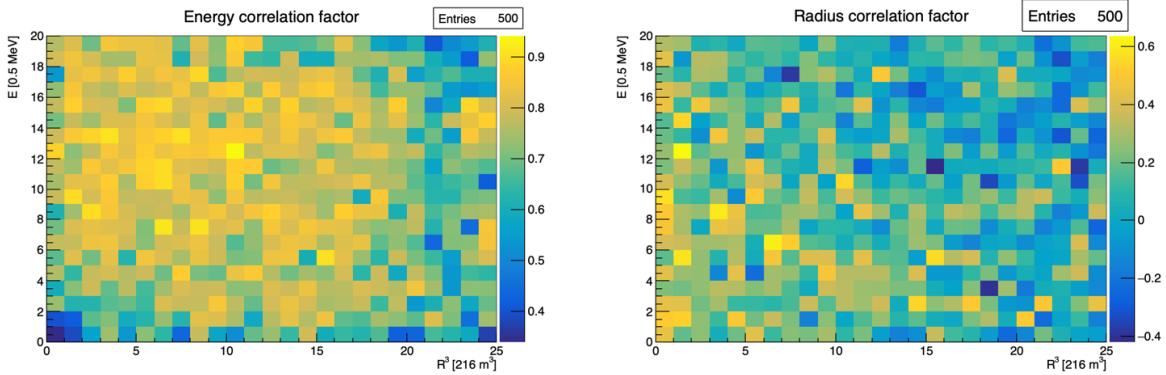


FIGURE 4.13 – Correlation between CNN and classical method reconstruction (on the left) for energy and (on the right) for radius in a  $E, R^3$  grid

reconstruct the  $\theta$  and  $\phi$  spherical coordinates and then rotate the image to locate the event in the center of the image. The second part, from this rotated image, would reconstruct the radius and energy of the event.

To overcome the problematic of the aggregation of PMT time information and the meaning of the time channel in case of no hit, we could transform this channel into a dimension. This would results in an image with multiple charge channels, each one representing the charge sum in a time interval.

In this thesis, we decided to solve those problem by moving away from the 2D image representation, looking into the graph representation and the Graph Neural Network (GNN). This is be the subject of the next chapter.

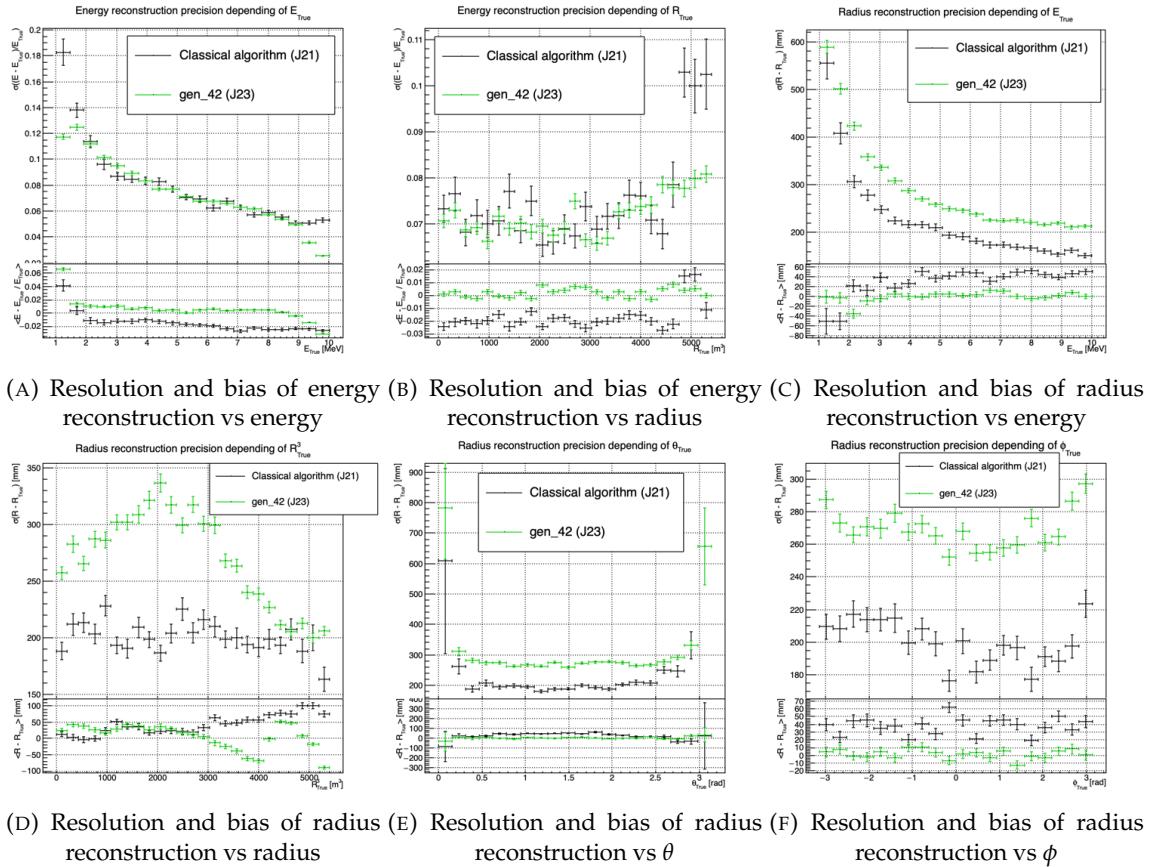


FIGURE 4.14 – Reconstruction performance of the “gen\_42” model on J23 data and its comparison to the performances of the classic algorithm “Classical algorithm” from [64]. The top part of each plot is the resolution and the bottom part is the bias.



<sup>1570</sup> **Chapter 5**

<sup>1571</sup> **Graph representation of JUNO for  
IBD reconstruction**

<sup>1573</sup>

*"The Answer to the Great Question of Life, the Universe and  
Everything is Forty-two"*

*Douglas Adams, The Hitchhiker's Guide to the Galaxy*

<sup>1574</sup>

## Contents

<sup>1575</sup>

<sup>1576</sup>

<sup>1577</sup>

<sup>1578</sup>

<sup>1579</sup>

<sup>1580</sup>

<sup>1581</sup>

<sup>1582</sup>

<sup>1583</sup>

<sup>1584</sup>

<sup>1585</sup>

<sup>1586</sup>

<sup>1587</sup>

|                                      |    |
|--------------------------------------|----|
| <b>5.1 Motivation</b>                | 73 |
| <b>5.2 Data representation</b>       | 74 |
| <b>5.3 Message passing algorithm</b> | 76 |
| <b>5.4 Data</b>                      | 78 |
| <b>5.5 Model</b>                     | 79 |
| <b>5.6 Training</b>                  | 80 |
| <b>5.7 Optimization</b>              | 80 |
| <b>5.8 Results</b>                   | 81 |
| <b>5.9 Conclusion</b>                | 82 |

<sup>1588</sup>

<sup>1589</sup>

<sup>1590</sup>

<sup>1591</sup>

<sup>1592</sup>

<sup>1593</sup>

<sup>1594</sup>

<sup>1595</sup>

We previously showed, in chapter 4, that neural networks are relevant as reconstruction tools in JUNO. Even if they show worse performances, the combination with classical estimators could still bring improvements. We discussed the use of Convolutional Neural Network (CNN) in the previous chapter and their limitations, more specifically the limitation of the image as data representation for the experiment.

In this chapter we propose to use a Graph Neural Network (GNN), a Neural Network specialized to process graph as presented in section 3.3.3, to overcome those limitations.

## 5.1 Motivation

<sup>1596</sup>

<sup>1597</sup>

<sup>1598</sup>

<sup>1599</sup>

<sup>1600</sup>

<sup>1601</sup>

As explained in chapter 2 the JUNO sensors, the Large Photomultipliers (LPMT) and Small Photomultipliers (SPMT), are arranged on a spherical plane. When trying to represent this plane as a 2D image, due to the inherent problem of the projection, some part of the image are distorted and part of the image do not have any physical meaning (see section 4.2.2). A way to represent the data without inducing deformation is the graph, an object composed of a collection of nodes and edges representing the relation between the nodes.

From this graph representation, we can construct a neural network that will process the data while keeping some interesting properties. For example the rotational invariance, i.e. the energy and

1604 radius of the event do change by rotation our referential. For more details see section 3.3.3. Graph  
1605 representation also has the advantage to be able to encode global and higher order informations.

1606 An approach was already proposed in JUNO by Qian et al. [42] where each nodes of the graph are  
1607 like pixels, they represent geometric region of the detector and are connected with their neighbours.  
1608 The LPMT informations are then aggregated on those nodes. The network then process the data  
1609 using the equivalent of convolution but on graph [49].

1610 In this work we want to take a step further in the graph representation by including the SPMT and  
1611 including a maximum of raw informations.

## 1612 5.2 Data representation

1613 In an ideal world we would like to have every PMTs represented as node in the graph, each PMT  
1614 being hit is an informations but the fact that PMTs were not hit is also an important information.  
1615 It's by being aware of the whole of the system that we are able to give meaning to a subpart. As a  
1616 reminder, in the Central Detector (CD), JUNO will posses 17612 LPMTs and 25600 SPMTs for a total  
1617 of 43212 PMTs. This amount of information in itself is still manageable by modern computer if it  
1618 were to be used in a neural network but when defining the relations between the nodes, it become a  
1619 bit more tricky.

1620 Excluding self relation and considering the relation to be undirected, the edge from  $A$  to  $B$  is the  
1621 same from  $B$  to  $A$ , the amount of necessary edges is given by  $\frac{n(n-1)}{2}$  which for 43212 PMTs amount  
1622 for 933'616'866 edges. If we encode an information with double precision (64 bits) in what we call an  
1623 adjacency matrix, each information we want to encode in the relation would consume 4 GB of data.  
1624 When adding the overhead due to gradient computation during training, this would put us over the  
1625 memory capacity of a single V100 gpu card (20 GB of memory). We could use parallel training to  
1626 distribute the training over multiple GPU but we considered that the technical challenge to deploy  
1627 this solution was not worth the trouble.

1628 The option of connecting PMTs node only to their neighbours could be tempting to reduce the num-  
1629 ber of edge, but this solution does not translate well in term of internal representation in memory.  
1630 Edges of sparsely connected nodes can be stored in efficient manner in a sparse matrix but the  
1631 calculation in itself would often results in the concretization of the full matrix in memory, resulting  
1632 in no memory gain during training.

1633 We finally decided of a middle ground where we define three *families* of nodes:

- 1634 — The core of the graph is composed of nodes representing geometric regions of the detector.  
1635 We call those nodes **mesh** nodes. Those mesh nodes are densely connected to each other. We  
1636 keep their number low to gain in memory consumption.
- 1637 — All the fired PMTs, that have been hit, will be represented as nodes. We call those node **fired**.  
1638 Fired nodes are connected to the mesh they geometrically belong.
- 1639 — A final node which will hold global information about the detector and on which we will read  
1640 the interaction vertex and energy. It's designated as the **I/O** node for input/output. This node  
1641 will be connected to every mesh nodes.

1642 Those nodes and their relations are illustrated in figure 5.1a. From this representation, we end up  
1643 with three distinct adjacency adjacency matrix

- 1644 — A  $N_{fired} \times N_{mesh}$  adjacency matrix, representing the relations between fired and mesh. Those  
1645 relations are undirected.
- 1646 — A  $N_{mesh} \times N_{mesh}$  adjacency matrix, representing the relation between meshes. Those relation  
1647 are directed.
- 1648 — A  $N_{mesh} \times 1$  adjacency between the mesh and I/O nodes. Those relations are undirected.

1649 The adjacency matrix representing those relation is illustrated in figure 5.1b.

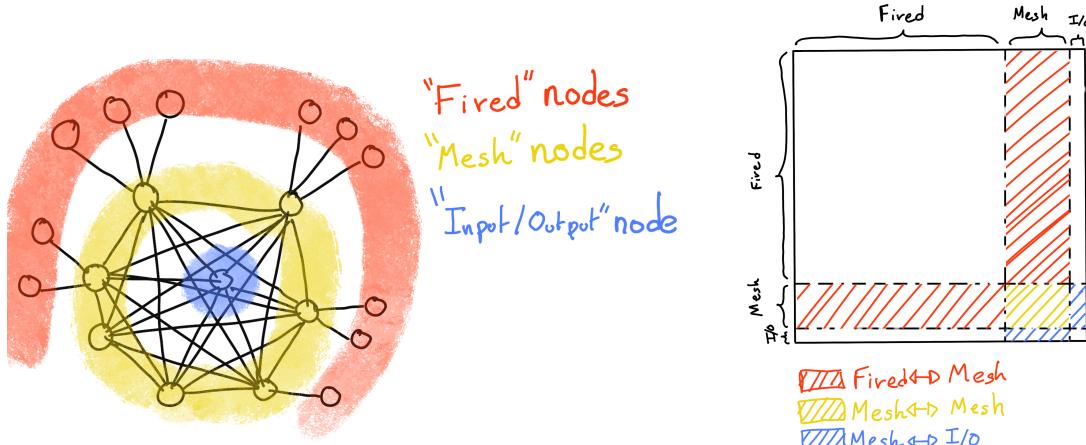


FIGURE 5.1



FIGURE 5.2 – Illustration of the healpix segmentation. On the left: A segmentation of order 0. On the right: A segmentation of order 1

The mesh segmentation is following the Healpix segmentation [74]. This segmentation offer the advantage that almost each mesh have the same number of direct neighbours and it guarantee that each mesh represent the same extent of the detector surface. The segmentation can be infinitely subdivided to provide smaller and smaller pixels. The number of pixel follow the order  $n$  with  $N_{pix} = 12 \cdot 4^n$ . This segmentation is illustrated in figure 5.2. To keep the number of mesh small, we use the segmentation of order 2,  $N_{pix} = 12 \cdot 4^2 = 192$ .

We decided on having the different kind of nodes **mesh (M)**, **fired (F)** and **I/O** have different set of features. The features used in the graph are presented in figure 5.3. Most of the features are low level informations such as the charge or time information but we include some high order features such as

1.  $P_l^h$ : Is the normalized power of the  $l$ th spherical harmonic. For more details about spherical harmonics in JUNO, see annex B.
2.  $\mathbb{A}$  and  $\mathbb{B}$  are informations that represent the likeliness of the interaction vertex to be on the

segment between the center of two meshes.

$$\mathbb{A}_{ij} = (\vec{j} - \vec{i}) \cdot \frac{\vec{l}_1}{D_{ij}} + \vec{i} \quad (5.1)$$

$$\mathbb{B}_{ij} = \frac{Q_i}{Q_2} \left( \frac{l_2}{l_1} \right)^2 \quad (5.2)$$

$$l_1 = \frac{1}{2}(D_{ij} - \Delta t \frac{c}{n}) \quad (5.3)$$

$$l_2 = \frac{1}{2}(D_{ij} + \Delta t \frac{c}{n}) \quad (5.4)$$

where  $\vec{i}$  is the position vector of the mesh  $i$ ,  $D_{ij}$  is the distance between the center of the meshes  $i$  and  $j$ ,  $Q_i$  the sum of charges on the mesh  $i$ ,  $\Delta t = t_i - t_j$  where  $t_i$  the earliest time on the mesh  $i$  and  $n$  the optical index of the LS.  $\mathbb{A}$  is the vertex between center of meshes distance ratio between  $i$  and  $j$  based on the time information. For  $\mathbb{B}$ , the charge ratio evolve with the square of the distance, so the mesh couple with the smallest  $\mathbb{B}$  should be the one with the interaction vertex between its two center.

| Nodes                               |                       |                      | Edges                   |   |                           |
|-------------------------------------|-----------------------|----------------------|-------------------------|---|---------------------------|
| Fire                                | Mesh                  | I/O                  | Fire $\rightarrow$ Mesh | Mesh $\rightarrow$ Mesh (1)   | Mesh $\rightarrow$ I/O    |
| $Q$                                 | $\langle Q_m \rangle$ | $\langle X \rangle$  | $X - X_m$               | $X_{m1} - X_{m2}$   | $\langle X \rangle - X_m$ |
| $t$                                 | $6Q_m$                | $\langle Y \rangle$  | $Y - Y_m$               | $Y_{m1} - Y_{m2}$   | $\langle Y \rangle - Y_m$ |
| $X$                                 | $\min(t_m)$           | $\langle Z \rangle$  | $Z - Z_m$               | $Z_{m1} - Z_{m2}$   | $\langle Z \rangle - Z_m$ |
| $Y$                                 | $\max(t_m)$           | $\Sigma Q$           | $t - \min(t)$           | $\min(t_1) - \min(t_2)$   | $Q_m / \Sigma Q$          |
| $Z$                                 | $6t_m$                | $P_l^h; l \in [0,8]$ | $Q / \Sigma Q_m$        | $\frac{\langle Q_{m1} \rangle - \langle Q_{m2} \rangle}{\langle Q_{m1} \rangle + \langle Q_{m2} \rangle}$ | $\langle t_m \rangle$     |
| <small>LPMT: 1<br/>SPMT: -1</small> |                       | $X_m$                |                         | $D_{m1 \rightarrow m2}^{-1}$  | $\mathbb{A}$              |
|                                     |                       | $Y_m$                |                         |   | $\mathbb{B}$              |
|                                     |                       | $Z_m$                |                         |   |                           |

$Q$  is the charge [nPE]  
 $t$  is the time [ns]  
 $X, Y, Z$  are the coordinates [cm]  
 $Q_m, t_m$  are the set of charge and time in a mesh  
 $X_m, Y_m, Z_m$  the coordinates of the center of the mesh  
 $\langle X \rangle, \langle Y \rangle, \langle Z \rangle$  the position of the charge barycenter.

FIGURE 5.3 – Features held by the nodes and edges in the graph.  $D_{m1 \rightarrow m2}^{-1}$  is the inverse of the distance between two mesh center. The features  $P_l^h$ ,  $\mathbb{A}$  and  $\mathbb{B}$  are detailed in section 5.2

Because our different nodes do not have the same number of features, they live in different spaces. Most library and public algorithms available are designed with node living in the same space in mind, we thus had to develop a custom message passing algorithm.

### 5.3 Message passing algorithm

As introduced in previous section and in figure 5.3, our graphs nodes and edges will have different number of features depending on their nature, meaning that we cannot have a single message passing function. We thus need to define a message passing function for each transition inside or outside a family. Using the notation presented in section 3.3.3

$$n_i^{k+1} = \phi_u(n_i^k, \square_j \phi_m(n_i^k, n_j^k, e_{ij}^k)); n_j \in \mathcal{N}'_i \quad (5.5)$$

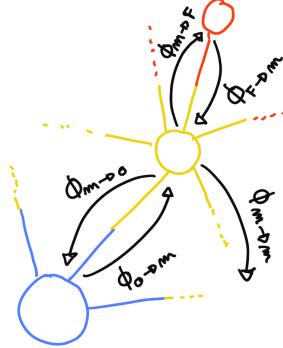


FIGURE 5.4 – Illustration of the different update function needed by our GNN

we need to define

$$\phi_{u;f \rightarrow m} \quad (5.6)$$

$$\phi_{u;m \rightarrow f} \quad (5.7)$$

$$\phi_{u;m \rightarrow m} \quad (5.8)$$

$$\phi_{u;io \rightarrow m} \quad (5.9)$$

$$\phi_{u;io \rightarrow m} \quad (5.10)$$

1676 to update the nodes after each layers as illustrated in figure 5.4. We would also need update function  
 1677 for the edges but for the sake of technical simplicity in this work, we will limit ourself to the nodes  
 1678 update. A wide variety of message passing algorithm exists, with different use cases and goal behind  
 1679 them. To stay generalist and to match to the best the specificity of our architecture, we implement  
 1680 the following algorithm:

$$\phi_u := I_i^{n'} = I_i^n A_{i,e}^i W_n^{e,n'} + I_i^n S_n^{n'} + B^{n'} \quad (5.11)$$

1681 using the Einstein summation notation.  $I_i^n$  is the tensor holding the nodes informations with  $i$   
 1682 the node index and  $n$  the feature index.  $n$  represent the features of the previous layer and  $n'$  the  
 1683 features of this layer.  $A_{i,e}^i$  is the adjacency tensor, discussed in the previous section, representing the  
 1684 connection between the node  $i'$  and the node  $i$ , each connections holding the features indexed by  $e$ .  
 1685 The learnable weights are composed of:

- 1686 — The tensor  $W_n^{e,n'}$  which represent the passage from the previous feature domain  $n$ , the previous  
 1687 layer, to the current domain  $n'$ , this layer, knowing the relation  $e$ .
- 1688 —  $B^{n'}$  which is a learnable bias tensor on the new features  $n'$ .
- 1689 —  $S_n^{n'}$  which can be viewed as a self loop relation where the node update itself based on the  
 1690 previous layer informations.

1691 If a node have neighbours in different families, the different  $I_i^{n'}$  coming from the different  $\phi_u$  are  
 1692 summed.

$$I_i^{n'} = \sum_{\mathcal{N}} \phi_{u,\mathcal{N}} \quad (5.12)$$

1693 where  $\mathcal{N}$  are the neighbouring family and  $\phi_{u,\mathcal{N}}$  the update function between the target node family  
 1694 and the neighbour  $\mathcal{N}$  family.

1695 We thus have a  $S$ ,  $W$  and  $B$  for each of the  $\phi_u$  function we defined above. The *IAW* sum can be seen  
 1696 as the  $\phi_m$  function and  $IS + B$  as the second part of the  $\phi_u$  function. Interestingly, the number on  
 1697 learnable weight in those layer is independent of the number of nodes in each family and depends  
 1698 solely on the number of features on the nodes and the edges.

1699 The expression above only update the node features. We could update the edges, using the results of  
 1700  $\phi_m$  for example, but for technical simplicity we only update the nodes and keep the edges constant.

1701 This operation of message passing is the constituent of our message passing layer, designed in this  
 1702 work as *JWGLayer*. To this layer, we can adjoin an activation function such as *PReLU*

$$I_i^{n'} = PReLU \left( \sum_{\mathcal{N}} I_i^n A_{i',e}^i W_n^{e,n'} + I_i^n S_n^{n'} + B^{n'} \right) \quad (5.13)$$

## 1703 5.4 Data

1704 For this study we will be using a 1M positrons event dataset, uniformly distributed in energy with  
 1705  $E_k \in [0, 9]$  MeV and uniformly distributed in the detector. Those events come from the JUNO official  
 1706 simulation version J23.0.1-rc8.dc1 (released the 7th January 2024). All the event are *calib* level, with  
 1707 simulation of the physics, electronics, digitizations and triggers. 900k events will be used for the  
 1708 training, 50k for validation and loss monitoring and 50k for the results analysis in section 5.8. Each  
 1709 events is between 2k and 12k fired PMTS, resulting in fired nodes being the largest family in our  
 1710 graphs in all circumstances as illustrated in figure 5.5c.

1711 As expected, by comparing the scale between the figure 5.5a and 5.5b we see that the LPMT system  
 1712 is predominant in term of informations in our data. The number of PMT hits grow with energy but  
 1713 do not reach 0 for low energy event due to the dark noise contribution which seems to be around  
 1714 1000 hits per event for the LPMT system (left limit of figure 5.5a) and around 15 hits per event for the  
 1715 SPMT system (left limit of figure 5.5b) which is consistent with the results show in section 4.2.2.

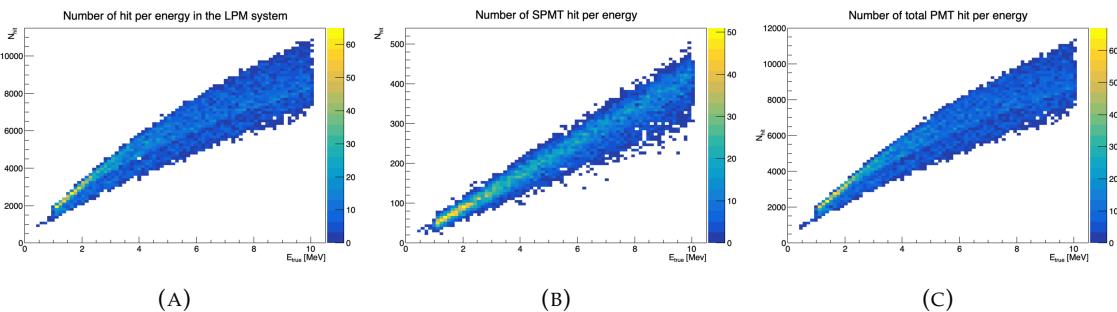


FIGURE 5.5 – Distribution of the number of hits depending on the energy. **On the right:** for the LPMT system. **In the middle :** for the SPMT system. **On the left:** For both system.

1716 The structure seen in the distribution in figure 5.5a comes from the shape of the number of hits  
 1717 depending on the radius as shown in figures 5.6a and 5.6b where the number of hit decrease with  
 1718 radius. It is important to understand that this is not representative of the number of PE per event  
 1719 and the decrease in hits over the radius means that the PE are just more concentrated in a smaller  
 1720 number of PMTs.

1721 No quality cut is applied here, we rely only on the trigger system. It means that event that would not  
 1722 trigger are not present in the dataset but for events that triggered twice, it happens rarely, the two  
 1723 trigger are considered as two separate event.

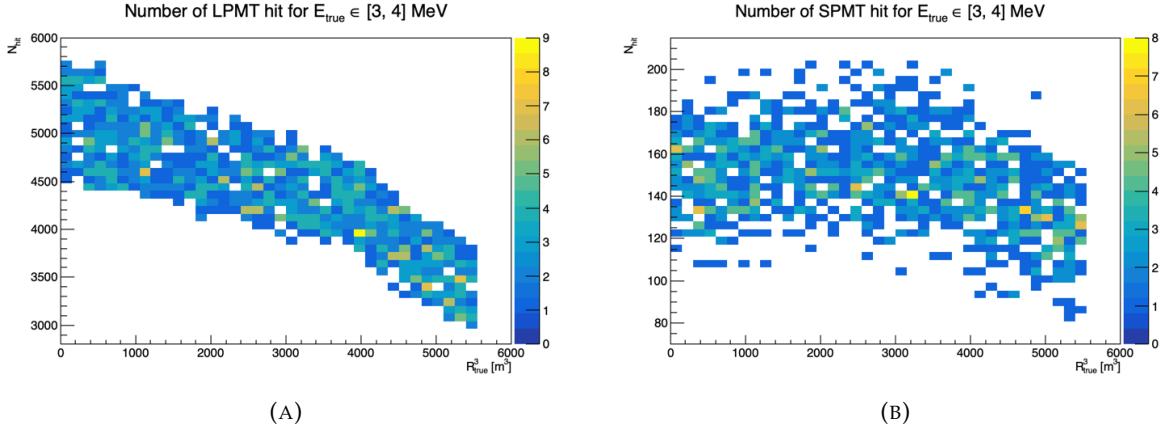


FIGURE 5.6 – Distribution of the number of hits depending on the radius. **On the right:** for the LPMT system. **On the right :** for the SPMT system. To prevent the superposition of structure of different scales we limit ourselves to the energy range  $E_{\text{true}} \in [0, 9]$ .

## 5.5 Model

In this section we'll discuss the different layer composing the final version of the model. As introduced above, each JWGLayer is defined by the number of features on the nodes and edges of the output graph, assuming it takes as input the graph from the precedent layer. For simplicity, when discussing a graph configuration, it will be presented as follow: {  $N_f$ ,  $N_m$ ,  $N_{IO}$ ,  $N_{f \rightarrow m}$ ,  $N_{m \rightarrow m}$ ,  $N_{m \rightarrow f}$  } where

- $N_f$  is the number of feature on the fired nodes.
  - $N_m$  is the number of features on the mesh nodes.
  - $N_{IO}$  is the number of features on the I/O node.
  - $N_{f \rightarrow m}$  is the number of features on the edges between the fired and mesh nodes.
  - $N_{m \rightarrow m}$  is the number of features on the edges between two mesh nodes.
  - $N_{m \rightarrow f}$  is the number of features on the edges between the mesh nodes and the I/O node.
- Because we do not change the number of features on the edges, we can simplify the notation to {  $N_f$ ,  $N_m$ ,  $N_{IO}$  }. As an example, the input graph configuration, following the figure 5.3, is { 6, 8, 13, 5, 8, 5 } or, without the edge features, { 6, 8, 13 }.

The final version of the model, called JWGV8.4.0 is composed of

- An JWGLayer, converting the input graph { 6, 8, 13 } to { 64, 512, 2048 } with a PReLU activation function.
- 3 resnet layers, each of them composed of
  1. 2 JWG layers with a PReLU activation function. They do not change the dimension of the graph
  2. A sum layer that sums the features in the input graph with the one computed from the JWG layers
- A flatten layer that flatten the features of the I/O and mesh nodes in a vector.
- 2 fully connected layers of 2048 neurons with a PReLU activation function.
- 2 fully connected layers of 512 neurons with a PReLU activation function.
- A final, fully connected layer of 4 neurons acting as the output of the network.

A schematic of the model is presented in figure 5.7.

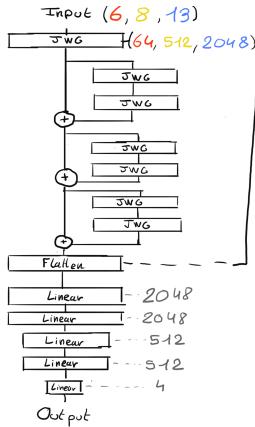


FIGURE 5.7 – Schema of the JWGv8.4.0 architecture, the colored triplet is the graph configuration after each JWG layers

## 1752 5.6 Training

1753 The optimizer used for training is the Adam optimizer and default hyperparameters ( $\beta_1 = 0.9$ ,  
 1754  $\beta_2 = 0.999$  and  $\epsilon = 1e-8$ ) with a learning rate  $\lambda = 1e-8$ . The training last 200 epochs of 800 steps.  
 1755 We use a batch size of 8. The learning rate is constant during the first 20 epochs then exponentially  
 1756 decrease with a rate of 0.99. The model saved is the model with the best validation loss during the  
 1757 training. The validation is computed over a single batch.

## 1758 5.7 Optimization

1759 Due to the extensive training time, up to 90h per training on the more complex architectures, and  
 1760 the heavy memory consumption of the models that would often exceed the 20GB limit of the V100,  
 1761 random search was not a realistic approach to the hyper optimisation. We were able to extend the  
 1762 memory limit to 40GB thanks to a local A100 GPU card available inside the laboratory.

1763 The hyperparameters optimization was thus done “by hand”, by looking at the results of the previous  
 1764 training and tinker hyperparameters that seems to play a role in the training. During this process,  
 1765 the model went into some heavy refactoring. At the start, the message passing algorithm was not  
 1766 the one presented above but each  $\phi_u$  and  $\phi_m$  function were FCDNN. Due to problems of memory  
 1767 consumption and gradient vanishing we pivoted to the message passing algorithm presented above.

1768 Even the features on the graph went under investigation. With the addition of high level observables  
 1769 to the mesh and I/O nodes and edge, there was too much possibility to test everything. We went  
 1770 with the decision to keep the raw observables in the fired and for the higher order observables  
 1771 we tried to take the one that would be difficult for the NN to reconstruct or at least would need  
 1772 multiple layer to reproduce. Basically, because the operation in the JWGLayer are linear operation,  
 1773 any variables dependent on order > 1 of the input would be candidates. This is why we introduce  
 1774 standard deviation,  $A$ ,  $B$  and  $P_l^h$  for example.

1775 Substantial effort went to the data processing process, transforming JUNO files into understandable  
 1776 graphs, before the training. Due to the volatile nature of the graph features during the optimization,  
 1777 the current code do not take preprocessed data and compute the observables, adjacency matrix,  
 1778 etc... on the fly. This data processing is carried out on the CPU, using a worker pool to allow for  
 1779 multiprocess. The raw data are coming from ROOT file produced by the collaboration software,

1780 the Event Data Model (EDM) used internally by the collaboration [75] had to be interfaced to our  
 1781 code, interface maintained through the evolution of the collaboration software. For the harmonic  
 1782 power calculation, we migrated from the Healpix library to Ducc0 [76] for a more fine control of the  
 1783 multithreading.

1784 Over the course of the project, the model went over more than 60 different configurations to end on  
 1785 the one presented in this chapter.

## 1786 5.8 Results

1787 The reconstruction performance of “JWGv8.4” are presented in figure 5.9 and compared to the “Omlil-  
 1788 rec” algorithm, the official IBD reconstruction algorithm in JUNO. Omlilrec is based on the QTMLR  
 1789 reconstruction method that was presented in section 2.6.

1790 We also present the results of the optimal variance combination of the two algorithm labelled as  
 1791 “JWG 8.4 x Omlilrec” where the reconstructed target  $\hat{\theta}_{\text{target}}$  is the weighted sum of the result of the  
 1792 two estimator JWGv8.4  $\theta_J$  and Omlilrec  $\theta_O$ .

$$\hat{\theta} = \alpha\theta_J + (1 - \alpha)\theta_O; \alpha \in [0, 1] \quad (5.14)$$

1793 For more details about the combination and the computation of  $\alpha$ , refer to annex A.2.

1794 One thing that need to be addressed before discussing results is that the Omlilrec algorithm do not  
 1795 reconstruct the deposited energy  $E_{\text{dep}}$  but reconstruct the visible energy  $E_{\text{vis}}$ . The difference between  
 1796 those two different observables comes from the event-wise and channel-wise non-linearity, presented  
 1797 in 2.3. The multiples energy observables are already discussed in section 4.4. For the following  
 1798 results, the systematic bias of Omlilrec that appear due, to the comparison to  $E_{\text{true}}$  instead of  $E_{\text{vis}}$  is  
 1799 corrected using a 5th degree polynomial

$$\frac{E_{\text{true}}}{E_{\text{rec}}} = \sum_{i=0}^5 P_i E_{\text{true}}^i \quad (5.15)$$

1800 The fitted distribution and the corresponding fit is presented in figure 5.8. The value fitted for this  
 1801 correction are presented in table 5.1.

|       |                                |
|-------|--------------------------------|
| $P_0$ | $1.24541 +/- 0.00585121$       |
| $P_1$ | $-0.168079 +/- 0.00716387$     |
| $P_2$ | $0.0489947 +/- 0.00312875$     |
| $P_3$ | $-0.00747111 +/- 0.000622003$  |
| $P_4$ | $0.000570998 +/- 5.7296e-05$   |
| $P_5$ | $-1.72588e-05 +/- 1.98355e-06$ |

TABLE 5.1 – Parameters of the 5th degree polynomial used to correct Omlilrec  
 reconstructed energy.

1802 Overall, energy and radius resolutions are not on par with Omlilrec. We see from the energy de-  
 1803 pendent energy resolution in fig 5.9a that our resolution is a bit more than twice the resolution of  
 1804 Omlilrec and the combination brings no improvements. Same observation for the energy resolution  
 1805 depending on the radius.

1806 The radius resolution, presented in the figures 5.9c, 5.9d, 5.9e and 5.9f is much worse than the  
 1807 Omlilrec one. This comes a bit as a surprise, as the energy reconstruction is dependent on the  
 1808 vertex reconstruction to correct for the non-uniformity and non-linearity effect. This mean that  
 1809 either the GNN could outperform the classical methods if the vertex was correctly reconstructed,

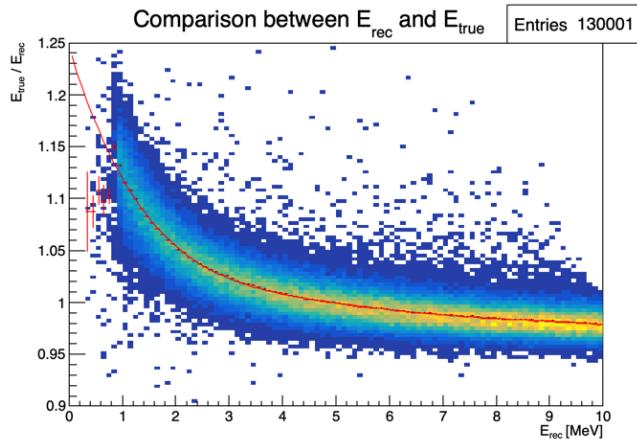


FIGURE 5.8 – Comparison between Omilrec  $E_{rec}$  and the true energy  $E_{true}$ . The profile of the distribution  $E_{true} / E_{rec}$  vs  $E_{rec}$  is fitted with a 5th degree polynomial.

1810 or that somewhere the GNN reconstruct the vertex correctly but has trouble to formulate it in x,y,z  
 1811 coordinates on the latest layer.

1812 The GNN behaviours are close to Omilrec, indicating that the same information is used in the same  
 1813 way by both algorithms, just that the GNN seems to be less fine-tuned than Omilrec. If the precedent  
 1814 reasoning is true, it would mean that by adding more parameters, more layer or a higher pixelisation  
 1815 of the Healpix representation, the GNN could reach Omilrec performances.

## 1816 5.9 Conclusion

1817 In this chapter, I present a proposition for a GNN architecture to reconstruct the energy and position  
 1818 of the prompt signal of an IBD interaction. The GNN is not competitive in terms of resolution with  
 1819 the more classical method Omilrec, which is the state of the art reconstruction method for IBD in the  
 1820 JUNO collaboration, but show encouraging results that could be exploited by going further in the  
 1821 optimisation of the hyper parameters. The message passing algorithm is still pretty naive and could  
 1822 probably be refined for JUNO's need.

1823 Another possible improvement is to find a way to increase the Healpix pixelisation. Through our  
 1824 different work on reconstruction and by looking at the different classical methods, it seems that  
 1825 the time information is crucial for the vertex reconstruction, and thus for the energy reconstruction.  
 1826 While we are keeping every raw informations about the fired PMTs, it is possible that the aggregation  
 1827 on mesh nodes could cause the information loss and it has been noticed that allowing more channels  
 1828 to the hidden layer mesh nodes improve the resolution. This observation can be compared to the  
 1829 convolutional GNN presented section 2.6.3 that has similar performance with the classical method  
 1830 with an order 5 Healpix segmentation resulting in 3072 pixels, comforting the need of a finer pixeli-  
 1831 sation, or more parameters dedicated to aggregation through an increase of channels on the mesh  
 1832 nodes. Both of those improvements require some heavy memory optimisations, distributed training  
 1833 or more powerful hardware to address the memory consumption issue.

1834 A final possible improvement would be to go further in the proximity of raw information. The charge  
 1835 and time used in the PMTs are extracted from a waveform, we could imagine a world where the full  
 1836 PMT waveform in the trigger window would be set of channels on the PMT node.

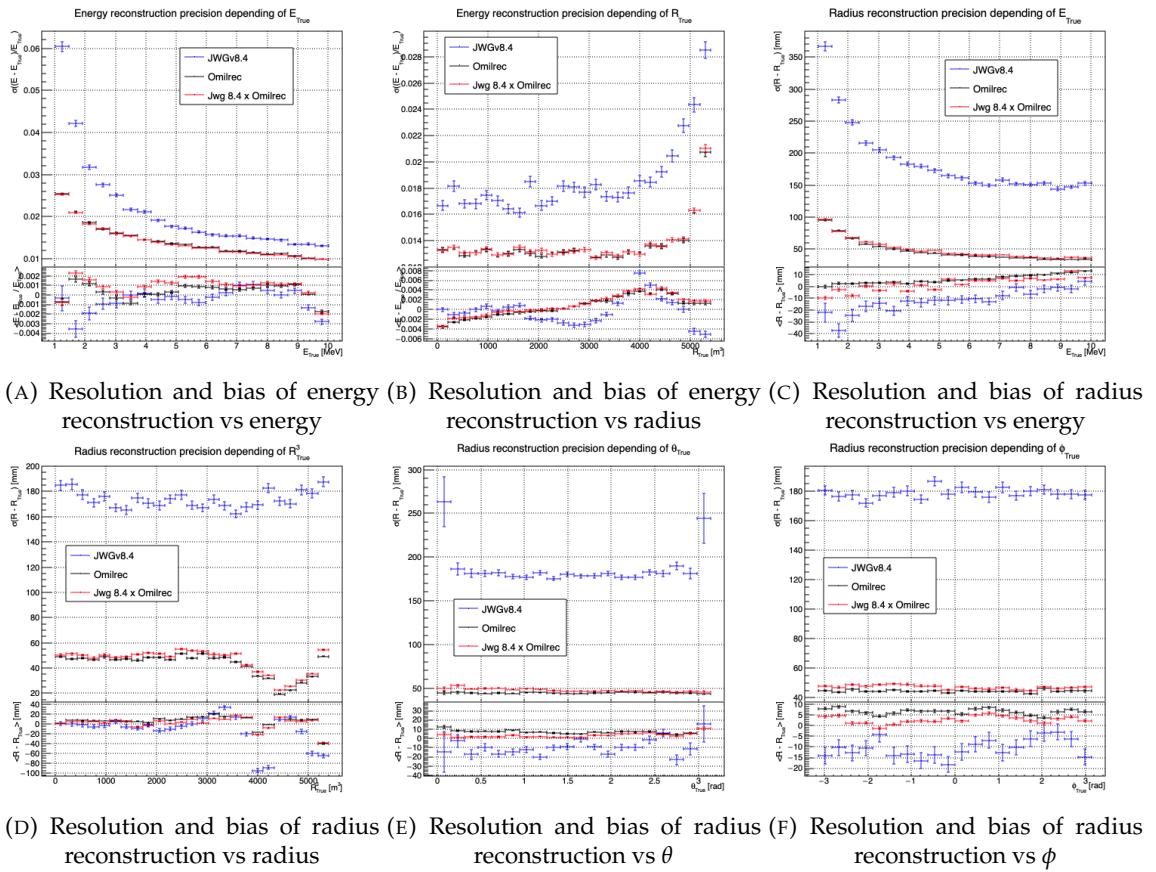


FIGURE 5.9 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.6, JWGv8.4 presented in this chapter and the combination between the two as presented in section 4.4.2. The top part of each plot is the resolution and the bottom part is the bias.



<sup>1837</sup> **Chapter 6**

<sup>1838</sup> **Reliability of machine learning  
methods**

<sup>1840</sup> *"Psychohistory was the quintessence of sociology; it was the science of  
human behavior reduced to mathematical equations. The individual  
human being is unpredictable, but the reactions of human mobs,  
Seldon found, could be treated statistically"*

*Isaac Asimov, Second Foundation*

<sup>1841</sup> **Contents**

---

|   |       |    |
|---|-------|----|
| <sup>1842</sup> <b>6.1 Motivations</b>                      | ..... | 86 |
| <sup>1843</sup> <b>6.2 Method</b>                           | ..... | 86 |
| <sup>1844</sup> <b>6.3 Architecture</b>                     | ..... | 86 |
| <sup>1845</sup> <b>6.3.1 Adversarial Neural Network</b>     | ..... | 86 |
| <sup>1846</sup> <b>6.3.2 Reconstruction Network</b>         | ..... | 87 |
| <sup>1847</sup> <b>6.3.3 Training</b>                       | ..... | 87 |
| <sup>1848</sup> <b>6.4 Results</b>                          | ..... | 87 |
| <sup>1849</sup> <b>6.4.1 Back to identity</b>               | ..... | 88 |
| <sup>1850</sup> <b>6.4.2 Breaking of the reconstruction</b> | ..... | 88 |
| <sup>1851</sup> <b>6.5 Conclusion and prospect</b>          | ..... | 88 |

---

<sup>1853</sup> <sup>1855</sup> As explained in previous chapters, JUNO is a precision experiment where the complete understanding of the effects at hand is crucial. As it will be illustrated in chapter 7, even small invisible biases or uncertainties could lead to the impossibility to run the measurements, or even worse, wrong our mass ordering measurements. While the liquid scintillator technology is well known and straightforward, this is the first time it is deployed to such scale, and for such precision. This novelty brings its fair share of elements, effects or assumption, that, if they were to be overlooked, could cause issue.

<sup>1862</sup> We already shown a large variety of reconstruction algorithms, OMILREC for LPMT reconstruction in section 2.6, numerous machine learning algorithms in section 2.6.3 and our own work in chapters 4 and 5. Those algorithms were compared to each other based on their performance as in [42] but we are the first that looked into the correlation between the reconstruction. The combinations of algorithms shown in chapter 4 and chapter 5 show that some information eludes the algorithms. We used this fact to try to improve our performance but this could also lead the algorithm to being vulnerable to some effect that could affect the detector and wrong the measurements.

<sup>1869</sup> The search for such effect could be done by hand, but the process would be tedious. We propose in this thesis a machine learning method to probe for those effects. In section 6.1, I delve further in the motivations of this work. In section 6.2, I describe the method behind the algorithm. In section 6.3 I detail the architecture of our algorithm and in section 6.4 the results of it. Finally, in section 6.5, I conclude and discuss about the prospect and possible improvements to bring to this work.

## 1874 6.1 Motivations

1875 As introduced above, JUNO needs a very good understanding of the biases and effects affecting its  
 1876 reconstruction as a small bias could wrong the mass ordering measurement. To calibrate those biases  
 1877 and effect, JUNO rely on multiples sources that will be located at various point in the detector. The  
 1878 calibration strategy was already discussed in section 2.3 and show calibrations sources of gammas,  
 1879 neutrons and positrons, with the catch that the positrons will annihilate inside the encapsulation and  
 1880 only the two 511 keV gammas will be seen. All those sources will be located at the center of the  
 1881 detector, impervious to non-uniformity.

1882 A second, natural, source will be used for calibration: The  $^{12}B$  spectrum. The  $^{12}B$  is a cosmogenically  
 1883 produced isotope through the passage of muons inside the LS. The  $^{12}B$  decays via  $\beta^-$  emissions with  
 1884 a Q value of 13.5 MeV with more than 98% of the decay resulting in ground state  $^{12}C$ . The  $^{12}B$  event  
 1885 will be cleanly identified by looking for delayed high energy  $\beta$  events after an energetic muon. Due  
 1886 to its natural causes, the  $^{12}B$  events will be uniformly distributed in the detector. The calibration  
 1887 strategy consist in fitting the energy spectrum of  $^{12}B$  with the results of the simulation to adjust the  
 1888 simulation parameters.

1889 We see that, while the calibration strategy is pretty complete, its missing a few points. First, none  
 1890 of the calibrations sources considered are positrons. While electrons and positrons events should  
 1891 be pretty similar in their interaction with the electronic cloud of the LS atoms, electron events are  
 1892 missing the two annihilations  $\gamma$  and the potential of forming a positronium [77]. The topology of the  
 1893 event thus differ of the order of magnitude of our reconstruction performance, a few nanoseconds  
 1894 for the energy deposit and positronium annihilation against a time transit spread between 3 and 6  
 1895 ns depending on the PMT type [78–80] and the  $\gamma$  will travel distances of the order of magnitude  
 1896 of the typical LPMT resolution of 8 cm (see section 2.6). Moreover, where for calibration sources  
 1897 the localization will be well known, the individual truth of  $^{12}B$  will be unknown. We thus need to  
 1898 compare our model to higher order observables such as energy distribution more than individual  
 1899 comparison.

1900 If there is potential failure point in those considerations, we need to search for them efficiently.

## 1901 6.2 Method

1902 All of the considerations could hide potential unknown or undetected effect that could lead to issue  
 1903 in the mass ordering analysis. But, while we have idea from where the issue could come, the  
 1904 production by hand of event perturbations that would not show in the calibration would be tedious.  
 1905 That's why we propose to use a Neural Network to produce those perturbations if they exists. A  
 1906 schematic of the concept is presented in figure 6.1.

## 1907 6.3 Architecture

- 1908 — Expliquer la problematique dans l'architecture
- 1909 — Ambition de pouvoir etre appliqu  a toutes les methodes, pas que NN
- 1910 — Pb technique: descente de gradient
- 1911 — Pr senter la loss

### 1912 6.3.1 Adversarial Neural Network

- 1913 — Decrire l'architecture de l'ANN

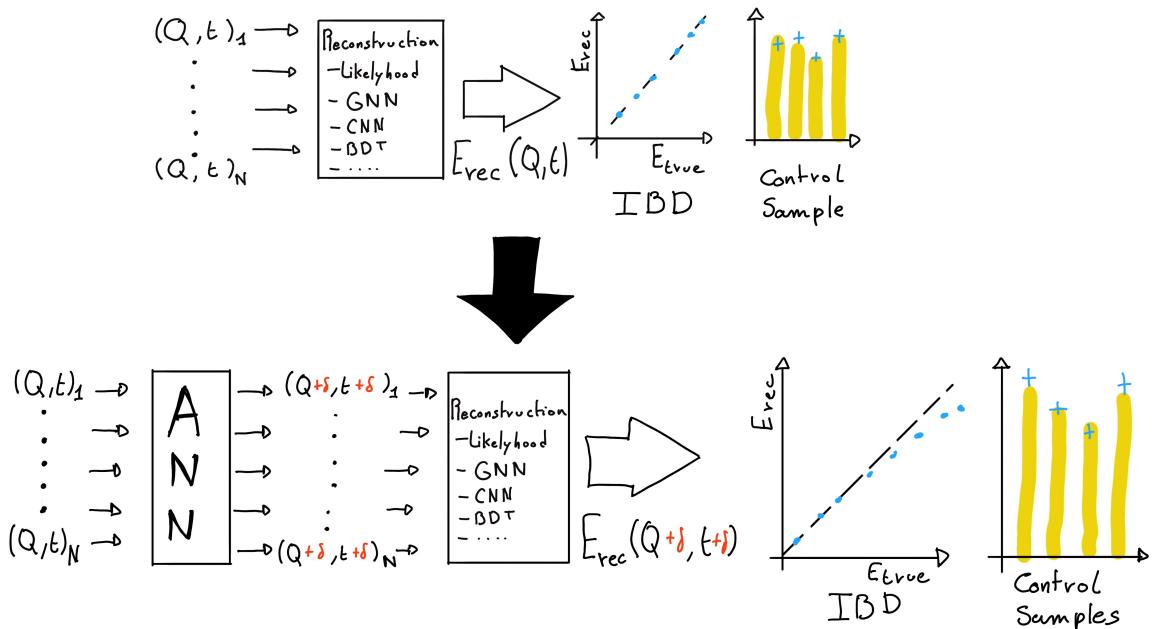


FIGURE 6.1 – Schema of the method to discover vulnerabilities in the reconstruction methods

### 1914 6.3.2 Reconstruction Network

- 1915 — Reseau de Neurone Simple. Deux avantages:
- 1916 — Besoin pour la descente de gradient
- 1917 — Un reseau "simpliste" a plus de chance de présenter des "défauts" que l'ANN pourrait exploiter
- 1918

### 1919 6.3.3 Training

- 1920 — Presentation du dataset
- 1921 — 2 etapes d'entraînement
- 1922 — Retour à l'identité -> que l'ANN ne fasse pas n'importe quoi
- 1923 — Cassage de la reconstruction

### 1924 Hyperparameter optimization

- 1925 — Pour les même raison que l'ANN:
- 1926 — Phase exploratoire, architecture très changeante, random search n'est pas viable
- 1927 — Architecture consomme beaucoup, besoin d'entraîner sur l'A100
- 1928 — Possiblement que de l'optimisation permettrait de faire passer sur V100, mais développement techniques nécessaires.
- 1929

## 1930 6.4 Results

- 1931 — Voir slide Gilles

<sup>1932</sup> **6.4.1 Back to identity**

<sup>1933</sup> **6.4.2 Breaking of the reconstruction**

<sup>1934</sup> **6.5 Conclusion and prospect**

<sup>1935</sup> — Not enough

<sup>1936</sup> — Probably guide the ANN

<sup>1937</sup> **Chapter 7**

<sup>1938</sup> **Joint fit between the SPMT and LPMT  
spectra**

<sup>1940</sup> “We demand rigidly defined areas of doubt and uncertainty!”  
*Douglas Adams, The Hitchhiker’s Guide to the Galaxy*

<sup>1941</sup> **Contents**

---

|   |       |     |
|---|-------|-----|
| <sup>1942</sup> <b>7.1 Motivations</b>                                | ..... | 90  |
| <sup>1944</sup> 7.1.1 Discrepancies between the SPMT and LPMT results | ..... | 90  |
| <sup>1945</sup> 7.1.2 Charge Non-Linearity (QNL)                      | ..... | 91  |
| <sup>1946</sup> <b>7.2 Approach</b>                                   | ..... | 92  |
| <sup>1947</sup> 7.2.1 Data production                                 | ..... | 92  |
| <sup>1948</sup> 7.2.2 Individual fits                                 | ..... | 93  |
| <sup>1949</sup> 7.2.3 Joint fit                                       | ..... | 94  |
| <sup>1950</sup> 7.2.4 Data and theoretical spectrum generation        | ..... | 96  |
| <sup>1951</sup> 7.2.5 Limitations                                     | ..... | 96  |
| <sup>1952</sup> <b>7.3 Fit software</b>                               | ..... | 97  |
| <sup>1953</sup> 7.3.1 IBD generator                                   | ..... | 97  |
| <sup>1954</sup> 7.3.2 Fit   | ..... | 99  |
| <sup>1955</sup> <b>7.4 Technical challenges and development</b>       | ..... | 99  |
| <sup>1956</sup> <b>7.5 Results</b>                                    | ..... | 100 |
| <sup>1957</sup> 7.5.1 Validation                                      | ..... | 100 |
| <sup>1958</sup> 7.5.2 Covariance matrix                               | ..... | 104 |
| <sup>1959</sup> 7.5.3 Statistical tests                               | ..... | 108 |
| <sup>1960</sup> <b>7.6 Conclusion and perspectives</b>                | ..... | 110 |

---

<sup>1963</sup> <sup>1964</sup> JUNO is an experiment of precise measurements, where we try to observe small fluctuation in the energy spectrum and with the goal to achieve sub-percent precision on the oscillation parameters measurement. A precise and complete understanding of the reconstruction and detector effects is thus crucial. The challenge reside in the technology used in the detector, which, while based on well known technology: scintillator observed by PMT, is being deployed on a scale never seen before, in term of scintillator volume and PMT size. Understanding every effects that goes in the detector can become extremely complicated. The ability to compare the results of the same experiment with two systems is thus extremely precious, this is the origin the dual calorimetry with the LPMT and SPMT system.

<sup>1973</sup> The resolution and bias of the reconstruction needs to be extremely well characterized: the target resolution of 3% [50] is unprecedented and is necessary to be able to distinguish between Normal

1975 Ordering (NO) and Inverse Ordering (IO). The non-linearity uncertainty needs to be constrained  
1976 under 1% as exceeding this value, the risk appear to measure the wrong ordering [27].

1977 One of the possible source of non-linearity, which will be used as a reference in this chapter, is the  
1978 charge non-linearity (QNL) that will be discussed in next section. The dual calorimetry can address  
1979 this issue, using calibrations methods and measurements that will be employed to correct it [27].

1980 More generally, comparing the results of the two systems will allow for the detection of potential  
1981 issues on the calibration or reconstruction. This is done in this thesis by comparing directly the  
1982 spectra and oscillation parameters measurements of the two systems.

1983 The study of the independent results of the two system can provide some informations [81] but this  
1984 is missing the important correlation that should be present between the two systems: they see the  
1985 same events, in the same scintillator, they're bound to be correlated. We explore in this chapter a  
1986 preliminary study of the impact of those correlations via multiple methods and the impact of QNL  
1987 at various degrees.

1988 In the next section we will discuss the motivations behind this study. In section 7.2, I present the  
1989 approaches and assumptions in this study. In section 7.3, I present the fit framework used, and then,  
1990 in section 7.4 the technical improvement brought and the difficulties faced during the development.  
1991 To end this chapter I present the results in 7.5 and discuss the conclusions and perspectives in 7.6.

## 1992 7.1 Motivations

### 1993 7.1.1 Discrepancies between the SPMT and LPMT results

1994 As discussed in the introduction of this chapter, the SPMT and LPMT systems will observe the same  
1995 events. This mean that, after calibration, if the two system show significant differences in their results  
1996 this is the signal of potential overlook of an effect or problem. Being able to detect such differences  
1997 is thus crucial, as discussed above, even the smallest deviation from our model could lead to the  
1998 impossibility to measure the Mass Ordering (MO) or even worse, wrong our measurement.

1999 The two systems are expected to have the same sensitivity to the oscillation parameters  $\theta_{12}$  and  $\Delta m_{21}^2$   
2000 [11]. We will thus rely on the measurement of those two parameters to detect potential discrepancies.

2001 We could just look at the value and compare them to the estimated independent error of the two  
2002 system, but we believe and will demonstrate in this chapter that the independent study of the two  
2003 system is missing a lot of informations, and that, by taking into account the statistic and systematic  
2004 correlations between the two systems, we can produce much more powerful statistical tests.

2005 Our work in this chapter is to develop such tools. The first step is, of course, to verify that in the  
2006 case of no discrepancies, the results are coherent with the independent analysis. This will give us the  
2007 distribution of those statistical test in absence of discrepancies. When we will have real data, we will  
2008 be able to compare it to those distributions to compute a p-value characterizing the absence of those  
2009 potential discrepancies.

2010 To evaluate the power of our methods, we need to simulate a concrete difference between the two  
2011 spectra. We have decided to study a plausible effect, the Charge Non-Linearity (QNL) that is detailed  
2012 next section. But the goal of those tools is to be discrepancy agnostic, as those discrepancies could  
2013 come from a variety of source (calibration issue, insufficient simulation tuning, etc...)

### 2014 7.1.2 Charge Non-Linearity (QNL)

2015 The CD energy response is subject to two kinds of non-linearity, the first one is the LS response  
 2016 non-linearity, where the LS photo-production is not linear with the deposited energy as illustrated  
 2017 in figure 2.12a. The second one is the LPMT response non-linearity where the charge read from the  
 2018 LPMT is not linear with respect to the number of collected Photo-Electrons (PE) (see section 2.3).

2019 The LS non-linearity comes from physic sources. Particle interactions in the LS will produce mainly  
 2020 scintillation light, as discussed in section 2.2.2, but will also produce some Cherenkov light (< 10%  
 2021 of the collected light). Both mechanisms possess intrinsic non-linearity, for the Cherenkov emission  
 2022 it depends on the velocity of charged particle velocity while the scintillation photon-yield follows a  
 2023 so-called Birk's law with a "quenching" effect depending on the energy and type of particle [16]. This  
 2024 results in am event-wise QNL.

2025 The LPMT response non-linearity can come from sheer saturation when subject to a high photon rate  
 2026 inducing a gain non-linearity or come from readout effects such as electronic noise, overshoot, the  
 2027 integration time window and even the waveform algorithm. All of these effects result in a channel-  
 2028 wise QNL.

2029 Precedent studies [27] suggest a model to emulate the non-linearity response that will be used in this  
 2030 work. We define the channel wise non-linearity that would be applied to each LPMT readout

$$\frac{Q_{rec}}{Q_{true}} = \frac{-\gamma_{qnl}}{9} Q_{true} + \frac{\gamma_{qnl} + 9}{9} \quad (7.1)$$

2031 where  $Q_{rec}$  is the reconstructed number of PE by the PMT,  $Q_{true}$  is true number of PE that hit the  
 2032 PMT, and  $\gamma_{qnl}$  is a factor representing the amplitude of the non-linearity.

2033 We also define an event-wise non-linearity characterized by

$$\frac{E_{vis}}{E_{true}} = \frac{-\alpha_{qnl}}{9} E_{true} + \frac{\alpha_{qnl} + 9}{9} \quad (7.2)$$

2034 where  $E_{vis}$  is the visible energy that is collected by the detector and  $E_{true}$  is the true deposited energy.  
 2035 An example of the effect of such event-wise QNL is presented in figure 7.1.

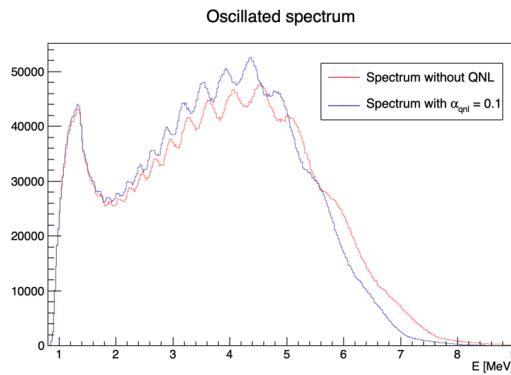
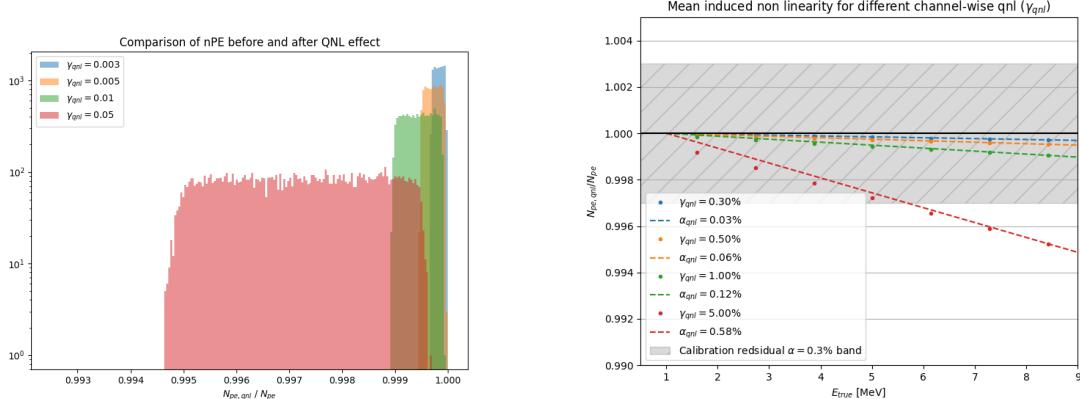


FIGURE 7.1 – Two oscillated spectra of  $1e7$  event expected in JUNO. In red the spectrum without supplementary QNL. In blue the same spectrum but where an event-wise QNL  $\alpha_{qnl} = 10\%$  is introduced.

2036 Using 1M events from the JUNO official simulation J23.0.1-rc8.dcl (released on 7th January 2024), we  
 2037 simulated events up to the photon collection in LPMTs and introduced an additional channel-wise  
 2038 QNL by using the equation 7.1 to modify the number of collected photons.



(A) Distribution of ratio of collected nPE after the additional QNL over the number of nPE that would be collected for different  $\gamma_{qnl}$ . We select event with an interaction radius  $R < 4\text{m}$  to not be affected by the non-uniformity.

(B) Ratio of collected nPE after the additional QNL over the number of nPE that would be collected at different energies. We select event with an interaction radius  $R < 4\text{m}$  to not be affected by the non-uniformity. The dots represent the mean of the distributions in figure 7.2a and the dashed line are the equivalent event-wise non-linearity from eq 7.2. The hatched zone is the residual non-linearity expected after calibration [29].

FIGURE 7.2

2039 In figure 7.2a we show the distribution of the ratio  $\frac{Q_{rec}}{Q_{true}}$  for central events ( $R < 4\text{m}$ ) and different  
 2040 values of  $\gamma_{qnl}$ . In figure 7.2a, we show the mean of this distribution as a function of the energy. We  
 2041 also present the effective  $\alpha_{qnl}$  for each value of  $\gamma_{qnl}$ . We observe that using the event-wise QNL is  
 2042 equivalent to the mean behavior of using channel-wise QNL.

2043 When using channel-wise non-linearity, we need to simulate a number of PE per LPMT, the process  
 2044 can be quite tedious if we want a realistic simulation. So in this study we are only using event-wise  
 2045 non-linearity to make the process simpler. This event-wise non-linearity will be characterized by  $\alpha_{qnl}$   
 2046 in this work.

## 2047 7.2 Approach

2048 In this section, we detail the testing procedure for each of our tools.

### 2049 7.2.1 Data production

#### 2050 IBD spectra

2051 The first step involves generating the data on which our tools will be tested. In this study we  
 2052 use Monte-Carlo toys. For each toy we generate a  $\bar{\nu}_e$  energy spectrum from the Taishan, Yangjiang  
 2053 and Dayabay nuclear power plants, the reactors used as source for the NMO analysis. The reactors  
 2054 parameters comes from JUNO official database, which shared among all physics analysis, the JUNO  
 2055 common inputs. This provides the initial spectra for the LPMT and SPMT systems. We then incorporate  
 2056 physic effects such as the LS non-linearity etc... (more details in section 7.3.1). Finally, we apply

2057 the reconstruction resolution for each system to their respective spectra, resulting in the final LPMT  
2058 and SPMT spectra.

2059 We will study the effect of exposure on our methods at different threshold: 100 days, 1 year, 2 year  
2060 and finally 6 years which is the nominal data taking period for the NMO analysis.

2061 These spectra are generated for different QNL,  $\alpha_{qnl} = 0$  (no spectrum distortion) and for  $\alpha_{qnl} \in$   
2062  $\{0.01, 0.005, 0.003, 0.002, 0.001\}$ . As a reminder, the calibration guarantees a residual event-wise non-  
2063 linearity of  $\alpha_{qnl} \leq 0.003$  [29].

The first test does not require any fitting, we are just comparing the LPMT and SPMT spectra using the expected statistical correlation matrix in the case  $\alpha_{qnl} = 0$ . For details about the generation of this correlation matrix, refer to section 7.5.2. This test is the spectrum  $\chi^2$  or  $\chi^2_{spe}$ . In this test we compute a  $\chi^2$  representing the compatibility between the LPMT and SPMT spectra:

$$\Delta_i = h_{L,i} - h_{S,i} \quad (7.3)$$

$$U = AVA^T \quad (7.4)$$

$$\chi^2_{spe} = \vec{\Delta}^T U^{-1} \vec{\Delta} \quad (7.5)$$

2064 Where  $h_{L,i}$  and  $h_{S,i}$  are the contents of the  $i$ th bin of the LPMT and SPMT spectra respectively.  $V$  is  
2065 the covariance matrix of the LPMT + SPMT spectra.  $A$  is a transformation matrix defined as:

$$A_{ij} = \frac{\partial \Delta_i}{\partial h_j} = \frac{\partial (h_{L,i} - h_{S,i})}{\partial h_j} \quad (7.6)$$

2066 Thus,  $A_{ij} = 1$  if  $i = j$ , and  $A_{ij} = -1$  if  $j$  is the SPMT bin corresponding to the  $i$  LPMT bin.

2067 This  $\chi^2_{spe}$  is minimal when the statistic between the bins of the LPMT and SPMT spectra follow the  
2068 covariance matrix  $V$ . By looking at the distribution of this  $\chi^2_{spe}$  when  $\alpha_{qnl} = 0$  we can produce  
2069 p-values for the values found when  $\alpha_{qnl} \neq 0$ .

## 2070 Background spectra

2071 The JUNO common inputs provide only LPMT background spectra. These background spectra are  
2072 already smeared by the LPMT resolution and thus need to be regenerated to be smeared to account  
2073 for the SPMT resolution. Fortunately the SPMT resolution is greater than that of the LPMT, allowing  
2074 us to apply additional smearing to the spectrum using

$$S(E) = L(E) \star \frac{1}{\sqrt{|\Delta\sigma^2|}\sqrt{2\pi}} e^{-\frac{E^2}{2|\Delta\sigma^2|}}; |\Delta\sigma^2| = \sigma_L^2 - \sigma_S^2 \quad (7.7)$$

2075 Where  $S(E)$  is the SPMT spectrum,  $L(E)$  the LPMT spectrum,  $\sigma_L$  and  $\sigma_S$  the LPMT and SPMT resolution  
2076 respectively. This formula is valid under the assumption that the LPMT and SPMT smearing are  
2077 gaussian and that the LPMT and SPMT have the same bias. Those two assumptions are valid in the  
2078 context of the IBD spectrum production as detailed in section 7.3.1. The demonstration of equation  
2079 7.7 can be found in annex C.

### 2080 7.2.2 Individual fits

Each of the spectra, LPMT and SPMT, are then fitted individually with and without the presence of QNL over multiples toys. The results allow us to compute the correlation between the oscillations parameters measured by both of the systems when there is no QNL allowing us to compute a  $\chi^2$

representing the compatibility between the measurements of the systems. Because the SPMT system is not sensible to the oscillation parameters  $\Delta m_{31}^2$  and  $\theta_{13}$ , the test is only done on the oscillation parameters  $\theta_{12}$  and  $\Delta m_{21}^2$ . We can thus produce the individual chi square  $\chi_{ind}^2$

$$\Delta_\lambda = \lambda_L - \lambda_S \quad (7.8)$$

$$\vec{\Delta} = [\Delta_{\theta_{12}} \Delta_{\Delta m_{21}^2}] \quad (7.9)$$

$$U = A V A^T \quad (7.10)$$

$$\chi_{ind}^2 = \vec{\Delta}^T U^{-1} \vec{\Delta} \quad (7.11)$$

where  $\lambda_L$  and  $\lambda_S$  are the measured parameters by the LPMT and SPMT systems respectively. The different  $\lambda$  considered are  $\theta_{12}$  and  $\Delta m_{21}^2$ .  $V$  here is the  $4 \times 4$  covariance matrix between the parameters  $\theta_{12,L}, \Delta m_{21,L}^2, \theta_{12,S}$  and  $\Delta m_{21,S}^2$ .  $A$  is the transformation matrix that allow us to compute the covariance matrix de  $\vec{\Delta}$  from  $V$  following

$$A_{ij} = \frac{\partial \Delta_i}{\partial j}; i \in \{\theta_{12}, \Delta m_{21}^2\}; j \in \{\theta_{12,L}, \Delta m_{21,L}^2, \theta_{12,S}, \Delta m_{21,S}^2\} \quad (7.12)$$

Same as described above, by comparing the distribution of this  $\chi_{ind}^2$  when  $\alpha_{qnl} = 0$  and  $\alpha_{qnl} \neq 0$  we can compute the power of this test in term of p-values.

### 7.2.3 Joint fit

#### Standard joint fit

The final step is to produce a joint fit between the two spectra. In this case we adjust our model, the oscillated spectrum, over two spectra at the same time. We minimize a  $\chi_{joint}^2$  defined over the two spectra, the LPMT and SPMT one

$$\Delta_i = D_i - T_i \quad (7.13)$$

$$\chi_{joint}^2 = \vec{\Delta}^T V^{-1} \vec{\Delta} \quad (7.14)$$

where  $D_i$  is the content of the  $i$ th bin measured, from the data, and  $T_i$  is the theoretical number of event in this bin.  $V$  is the covariance matrix of our spectrum.

$T$  is the fitted function and depend on multiple parameters

- The oscillation parameters  $\theta_{12}, \Delta m_{21}^2, \theta_{13}$  and  $\Delta m_{31}^2$ . Those parameters can be free, have a pull term or be fixed during the fit.
- We take into account in the data production the matter effect and parametrize it by the parameter  $\rho$ , the effective rock density between the reactors and the experiment. Same as the oscillation parameters, this parameter can be free, pulled or fixed.
- The exposure of the considered data which is just a normalization factor in front of the theoretical spectrum. This parameter is fixed at the start of the fit.

In the standard joint fit, the free parameters are  $\sin^2(2\theta_{12}), \Delta m_{21}^2$  and  $\Delta m_{31}^2$ .  $\sin^2(2\theta_{13})$  is fixed to the PDG nominal value. For simplicity, we refer to  $\sin^2(2\theta_{12})$  and  $\sin^2(2\theta_{13})$  as  $\theta_{12}$  and  $\theta_{13}$  respectively.

Both of the LPMT and SPMT systems are sensitive to  $\theta_{12}$  and  $\Delta m_{21}^2$ , thus these parameters are totally free and start at the PDG nominal value. Only the LPMT system is sensitive to  $\Delta m_{31}^2$ , we let it free so we can observe the effect of the deformation on it while the solar parameters  $\theta_{12}, \Delta m_{21}^2$  are constrained by the SPMT system. To prevent  $\Delta m_{31}^2$  to take absurd value, we add a pull term using the PDG nominal value and errors. The PDG nominal values used in this study can be found in table

| $\sin^2(2\theta_{12})$    | $\Delta m_{21}^2$                           | $\Delta m_{31}^2$                              | $\sin^2(2\theta_{13})$ |
|---------------------------|---|--|------------------------|
| $0.851^{+0.020}_{-0.018}$ | $7.53 \pm 0.18 \times 10^{-5} \text{ eV}^2$ | $2.5283 \pm 0.034 \times 10^{-3} \text{ eV}^2$ | $0.08523 \pm 0.00268$  |

TABLE 7.1 – Nominal PDG2020 value [16]. All value are reported assuming Normal Ordering.

2106 7.1.

$$\chi_{joint}^2 = \vec{\Delta}^T V^{-1} \vec{\Delta} + \frac{\Delta m_{31}^2 - \Delta m_{31,PDG}^2}{\sigma_{31,PDG}} \quad (7.15)$$

2107  $\theta_{13}$  is the parameter on which we are least accurate. It's fixed to nominal value to prevent degeneracy  
2108 (table 7.1).

2109 The covariance matrix is produced from a correlation matrix  $C$

$$V_{ij} = \sigma_i \sigma_j C_{ij} \quad (7.16)$$

2110 where  $\sigma_i$  is the uncertainty on the number of event in the  $i$ th bin. We consider in this study that the  
2111 content of each bin follow a Poisson statistic, thus the uncertainty is  $\sigma_i = \sqrt{N_i}$  where  $N_i$  is the content  
2112 of the  $i$ th bin. The bin content used for the uncertainty can come from two sources: the data and the  
2113 theoretical spectra  $\sigma_i = \sqrt{D_i}$  (Pearson test) and  $\sigma_i = \sqrt{T_i}$  (Neyman test). Precedent studies have  
2114 show that both Pearson and Neyman tests show bias at low statistic, we thus use the Pearson V test  
2115 where

$$\chi_{joint}^2 = \vec{\Delta}^T V^{-1} \vec{\Delta} + \frac{\Delta m_{31}^2 - \Delta m_{31,PDG}^2}{\sigma_{31,PDG}} + \ln|V| \quad (7.17)$$

2116 and the covariance matrix  $V$  is computed using the data spectrum for the uncertainty.

2117 The estimation of the covariance is crucial in this study as the strength of this test rely on the sys-  
2118 tematic and statistical correlations between the LPMT and SPMT spectrum. The generation methods  
2119 and results of this matrix is detailed in section 7.5.2.

## 2120 Delta joint fit

2121 Using the same structure we define a second joint fit, the Delta joint fit where, in addition to every-  
2122 thing that was discussed above, we add two other parameters  $\delta\theta_{12}$  and  $\delta\Delta m_{21}^2$  and split the theoretical  
2123  $T(\theta_{12}, \Delta m_{21}^2, \dots)$  spectrum in two

$$T_{LPMT} \equiv T(\theta_{12} + \delta\theta_{12}, \Delta m_{21}^2 + \delta\Delta m_{21}^2, \dots) \\ T_{SPMT} \equiv T(\theta_{12}, \Delta m_{21}^2, \dots) \quad (7.18)$$

2124 If the there is no additional distortion between the LPMT and the SPMT spectra, the fit should  
2125 converge to  $\delta\theta_{12} = \delta\Delta m_{21}^2 = 0$ . By observing the dispersion of those parameters we can define  
2126 the probability  $P(\alpha_{qnl} = 0 | (\delta\theta_{12}, \delta\Delta m_{21}^2))$  and use the median value of  $(\delta\theta_{12}, \delta\Delta m_{21}^2)$  when  $\alpha_{qnl} \neq 0$   
2127 to define a p-value.

2128 The last test we explore in this thesis is to fit the same spectrum with the Standard Joint fit, that  
2129 we consider as the hypothesis without distortion  $H_0$ , and the Delta Joint fit, designated as the  $H_1$   
2130 hypothesis. By looking at the dispersion of  $\chi_{joint,H_0}^2 - \chi_{joint,H_1}^2$  we can extract a sensitivity to potential  
2131 distortion.

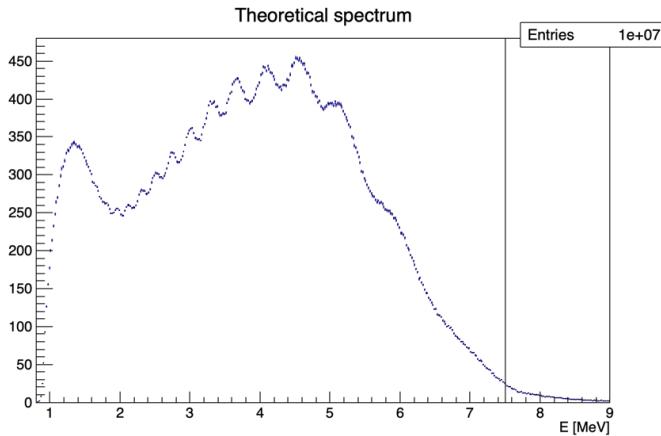


FIGURE 7.3 – Theoretical LPMT spectrum at nominal oscillation values binned using 410 bins from 0.8 to 9 MeV. It is rescaled to 6 years statistic. The black line represent the 335 bin cut

#### 2132 7.2.4 Data and theoretical spectrum generation

2133 To implement the joint fit, we have technically two data spectra and two theoretical spectra. The data  
 2134 in this study are produced using an IBD generator *IBD gen*, see section 7.3.1. The theoretical spectrum  
 2135 are produced the same way as data spectrum but with much higher statistics,  $10^7$  events to compare  
 2136 with the  $\approx 10^5$  events for 6 years statistic. The two spectrum, that we get as a collection of events,  
 2137 are binned in two histograms from 0.8 to 9 MeV of reconstructed energy with bins of 0.02 MeV each,  
 2138 resulting in 410 bins per spectrum. An illustration of the theoretical spectrum can be found in figure  
 2139 7.3. The low number of events in the tail of the spectrum can cause instability due to the low statistic,  
 2140 we thus cut the spectrum at 7.5 MeV / 335 bins for the fit.

2141 All the IBD spectra presented and used in this study are produced assuming Normal Ordering using  
 2142 the PDG nominal value [16] for the oscillation parameters. Those values are reported in table 7.1.

#### 2143 7.2.5 Limitations

2144 In this work we are only working considering the statistical errors. We can ignore systematic effects,  
 2145 such as effects that would affect the neutrino spectrum or the background spectrum, as they are  
 2146 entirely correlated between the two systems. The details of those systematic effects can be found in  
 2147 [11].

2148 Most of our results assume decorrelated detection effects between the SPMT and LPMT systems.  
 2149 Their respective reconstruction effects are simulated using simple gaussian drawing on the resolution,  
 2150 independently from the event position. This approach was used in previous sensitivity and  
 2151 precision studies [11, 82]. The potential effect of those reconstruction effects and a first attempt to  
 2152 take them into account are explored in section 7.5.2.

2153 Even if the goal of this work is to propose deformation agnostic tools, the QNL we use in this study is  
 2154 simplistic as we consider event-wise, position uniform deformation. We show in figure 7.2a and 7.2b  
 2155 that event-wise QNL is equivalent to the mean behaviour of channel-wise QNL but a more complete  
 2156 study would simulate channel-wise deformation for each event.

## 2157 7.3 Fit software

2158 In this section, I describe the ft framework that was used in this study. The software is composed  
 2159 of two parts as illustrated in figure 7.4: A standalone part composed of ROOT [83] macros, and the  
 2160 Avenue framework.

2161 The Avenue framework is responsible for the spectrum and configuration reading, transforming  
 2162 the raw collection of events into spectra, managing the physics effect such as the oscillation and  
 2163 computing and minimizing the  $\chi^2$  with the help of the RooFit library. The macros are invoking, if  
 2164 necessary, the Avenue framework and are the entry point for fitting, generating the necessary inputs  
 2165 quantity such as the spectra and correlation matrix, analysing the fit results and managing jobs for  
 2166 distributed computing.

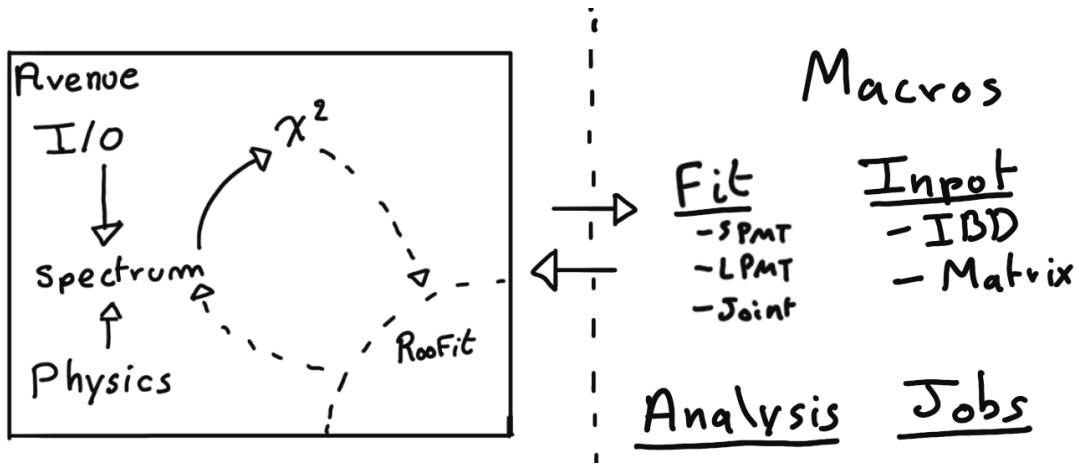


FIGURE 7.4 – Schematic description of the fit framework

2167 In this section we will focus on the IBD generator in section 7.3.1 and the fit macro in itself in section  
 2168 7.3.2.

### 2169 7.3.1 IBD generator

2170 The IBD generator is a standalone generator used to produce oscillated and non oscillated spectra  
 2171 as the one seen by the JUNO experiment. It takes as inputs physics parameters and a collection  
 2172 of histograms, values and function provided by JUNO to its analysis groups, referred as the JUNO  
 2173 common inputs.

2174 Options allow to enable or disable effects such as non-uniformity and non-linearity. It finally take as  
 2175 an argument the number of events to generate  $N_{evt}$ . Optionally, we generate an effective number of  
 2176 events  $N$  by drawing in a Poisson distribution of mean  $N_{evt}$ .

2177 Then for each event we

- 2178 1. Choose randomly, following the reactor power fraction, the source reactor of the neutrino.
- 2179 2. Generate a random interaction position in the detector following a uniform distribution over  
 2180 the detector volume.
- 2181 3. Draw a random neutrino energy  $E_\nu$  from the expected neutrino emission spectrum of every  
 2182 reactor. This spectrum is computed by:
  - 2183 (a) Computing the power spectrum of each isotopes  $^{235}\text{U}$ ,  $^{238}\text{U}$ ,  $^{239}\text{Pu}$ ,  $^{241}\text{Pu}$  using the Huber-  
 2184 Mueller model [5, 8].

- 2185 (b) Summing the contribution of each isotopes following the respective fission fraction [0.58,  
 2186 0.07, 0.30, 0.05] as reported in [84].  
 2187 (c) The power of each reactor is then adjusted by their distances from the detector, the detector  
 2188 efficiency and their mean duty cycle (11 of 12 month).  
 2189 (d) The total spectrum is then finally adjusted by taking into account the correction of the Day  
 2190 Bay bump [85], adjustment due to spent nuclear fuel and due to the non-equilibrium.  
 2191 4. (Optional) Compute the survival probability due to oscillation at nominal oscillation param-  
 2192 eters value. If the neutrino does not survive, the event is rejected and the algorithm restart  
 2193 from step (1).  
 2194 5. Compute the emitted positron energy  $E_{pos}$  from the mass difference. If the neutrino does not  
 2195 have enough energy reject the event and start from step (1).  
 2196 6. Compute the deposited energy  $E_{dep}$  by incrementing  $E_{pos}$  by 511 keV to account for the positron  
 2197 annihilation. We do not consider cases where some of the energy leak outside of the detector  
 2198 (positron or annihilation gammas escaping the CD).  
 2199 7. Correct the deposited energy with the expected event-wise non-linearity from [29] to obtain  
 2200 the visible energy  $E_{vis}$ .  
 2201 8. (Optional) Add a custom non-linearity as described in section 7.1.2. This non linearity is  
 2202 characterized by  $\alpha_{qnl}$  to obtain  $E_\alpha$ .  
 2203 9. Finally, using the expected resolution of the LPMT and SPMT systems, provided in the JUNO  
 2204 common inputs, we draw from a gaussian characterized by those resolution the reconstructed  
 2205 energy  $E_{rec}$  or  $E_{lpmt}$  and  $E_{spmt}$  for each systems. The resolutions are provided as ABC param-  
 2206 eters using

$$\frac{\sigma E_{vis}}{E_{vis}} = \sqrt{\left(\frac{A}{\sqrt{E_{vis}}}\right)^2 + B^2 + \left(\frac{C}{E_{vis}}\right)^2} \quad (7.19)$$

2207 where A is the term driven by the Poisson statistics of the total number of detected photoelec-  
 2208 trons, C is dominated by the PMT dark noise, and B is dominated by the detector's spatial  
 2209 non-uniformity. The relative and absolute resolutions of the LPMT and SPMT systems are  
 2210 illustrated in figure 7.5.

2211 The events are stored as n-tuples and are not yet binned at the end of the generator.

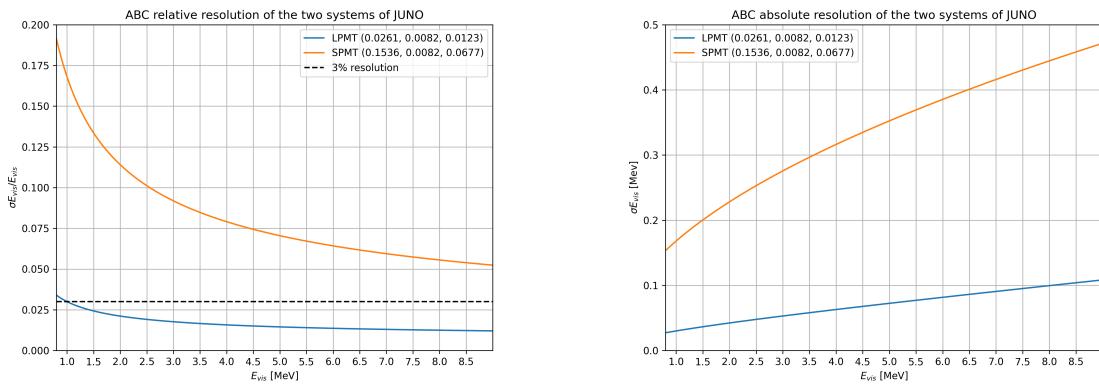


FIGURE 7.5 – Relative (On the left) and absolute (On the right) resolutions of the LPMT and SPMT systems used in this study. The number in parenthesis are the parameter A, B and C respectively for each systems.

---

### 2212 7.3.2 Fit

2213 The fit macro is the core of this fitting procedure. This macro is responsible for loading the fit  
 2214 configuration and setup the Avenue framework. Using Avenue, it will setup the data files, theoretical  
 2215 spectrum, choose the binning,  $\chi^2$ , etc... It also have the possibility to generate toys on the fly based  
 2216 on the theoretical spectrum. Given this theoretical spectrum we can randomize the bin content either  
 2217 by:

- 2218 1. Drawing the bin content in a Poisson distribution with the bin content as parameter.
- 2219 2. Drawing the bin content in a Gaussian distribution with the bin content as mean and variance.  
     The bin content is then rounded to the nearest integer.
- 2220 3. Drawing the bin difference following a given covariance matrix using the Choleski decomposi-  
     tion. This matrix is at least the statistical covariance matrix but can also contain systematic  
     uncertainties.

$$V = LL^T \quad (7.20)$$

$$\mathbf{R} \sim \mathcal{N}(0, 1) \quad (7.21)$$

$$\tilde{\mathbf{h}} = \lceil \mathbf{h} + L\mathbf{R} \rceil \quad (7.22)$$

$$(7.23)$$

2221 where  $V$  is covariance matrix used to produce the fluctuations,  $\mathbf{R}$  is drawn in a multinomial  
 2222 distribution of mean 0 and variance 1,  $\mathbf{h}$  the bin content of the theoretical spectrum and  $\tilde{\mathbf{h}}$  the  
 2223 bin content of the generated toy.

2224 The first two methods allow for the fast production of independent toys while the third allow for  
 2225 the production of statistical and systematical dependent toys. Unfortunately, none of those methods  
 2226 are fitted to produce toy with a QNL different from the theoretical spectrum. The uncertainty on the  
 2227 reconstructed energy  $\sigma E_{rec}$  being dependent on  $E_{vis}/E_\alpha$  makes that we would need to deconvolute  
 2228 the reconstruction effect from the theoretical spectrum. It is much easier to just produce those toys  
 2229 from the IBD generator.

## 2230 7.4 Technical challenges and development

2231 The fit framework Avenue was already partially developed with multispectra fitting in mind but  
 2232 a lot technical development was necessary to allow for a joint fit. The first step was to migrate  
 2233 the framework from ROOT5 (last release in March 2018) to ROOT6 (v6.26.06 released in July 2022)  
 2234 to ensure compatibility with the data coming from the JUNO collaboration, and benefiting of the  
 2235 improvement and corrections that came with ROOT6. This allow us to upgrade the C++ standard  
 2236 from C++11 to C++17. A substantial effort has been done to modernize the code, generalizing the  
 2237 functions and methods via templating to help readability and using smart pointer to prevent possible  
 2238 memory leaks.

2239 The Avenue framework had to be adapted, notably on the chi-square calculation and spectrum gen-  
 2240 eration to correctly take into account the correlation between the SPMT and LPMT spectra. The delta  
 2241 joint fit requiring two more parameters over a spectrum twice as large as before with LPMT takes  
 2242 much more time, around 15h for 6 years exposure, than the single LPMT fit. Thus the framework  
 2243 and the fit macro had to be updated for distributed computing. Notably the aggregation of fit results  
 2244 can now be done in a single file instead of managing a file per fit. In case of numerous toy, the hard  
 2245 drive access time could lead to long analysis time.

2246 While the IBD generator was already able to generate LPMT and SPMT spectrum, it was not designed  
 2247 for generating correlated spectrum. As detailed in section 7.3.1, up to the reconstruction effect, the

2248 two spectrum need to share the same generation else the two spectrum would be decorrelated and it  
 2249 would be like we would run two different experiment.

## 2250 7.5 Results

### 2251 7.5.1 Validation

2252 The first step is to confirm that the updated fit framework is able to reproduce existing results and  
 2253 that the joint fit behave as expected, meaning

- 2254 — Without QNL, the individual (*LPMT* and *SPMT*) fit converge to the parameters nominal  
 values and their errors are similar to the ones reported in existing analysis such as [11].
- 2255 — The standard joint fit with an independent covariance matrix (*Indep Standard joint*), meaning  
 2256 that the covariance between the LPMT and SPMT spectra is 0, believe to have twice as much  
 2257 informations, and thus believe to have a grater precision than the individual fits.
- 2258 — The standard joint (*Standard joint*) fit with a correlated covariance matrix has errors similar to  
 2259 the LPMT individual fit as the LPMT drive the precision on  $\theta_{13}$  and  $\Delta m_{31}^2$  and that the LPMT  
 2260 as SPMT are expected to have close precision on  $\theta_{12}$  and  $\Delta m_{21}^2$ .
- 2261 — The delta joint (*Delta joint*) fit with covariance matrix have the same resolution as the standard  
 2262 joint fit. The supplementary parameter  $\delta\theta_{12}$  and  $\delta\Delta m_{21}^2$  should not bring supplementary  
 2263 precision.

2264 The italicized name are the name used in the results reports to identify each fit. We also look into the  
 2265 *Indep Delta joint*, which is the Delta Joint fit but the covariance between the LPMT and SPMT spectra  
 2266 is 0, and the *Weighted* results where

$$\frac{1}{\sigma_{\text{Weighted}}^2} = \frac{1}{\sigma_{\text{LPMT}}^2} + \frac{1}{\sigma_{\text{SPMT}}^2} \quad (7.24)$$

2267 We expect the weighted resolution to be similar to the *Indep Standard joint* as, in both of those test, we  
 2268 do not consider the correlation between the SPMT and LPMT results.

### 2270 Asimov studies

2271 We ran Asimov studies on the tests presented above on the updated framework, the results are  
 2272 reported in table 7.2. All those test are ran considering statistics error only, 6 years exposure with  
 2273 all backgrounds, Pearson  $\chi^2$  (covariance is estimated using data spectrum) and  $\theta_{13}$  fixed to nominal  
 2274 value. For the *SPMT* fit  $\Delta m_{31}^2$  is fixed at nominal value as the SPMT system is net expected to be  
 2275 sensitive to this parameter.

2276 In every cases presented above, the fit converges to the parameters nominal value thus only the  
 2277 errors are presented.

2278 We observe, as expected, that  $\sigma_{\text{Weighted}} \approx \sigma_{\text{Indep Standard joint}}$  with the exception of  $\sigma\theta_{12}$ . This could  
 2279 from the slight difference in statistic between the SPMT and LPMT spectra. Indeed, due to a larger  
 2280 smearing in energy resolution, events that would be inside the spectrum range [0.8, 7.5] MeV are  
 2281 smeared outside it. This deficit is partially compensated by event outside the spectrum coming back  
 2282 in it but we expect very few event outside the spectrum in comparison to event at the edges of it.  
 2283 Thus the event deficit is not totally compensated.  $\theta_{12}$  being mainly driven by the amplitude of the  
 2284 spectrum (see illustration 2.2), that's why we think this the origin of the difference.

2285 The second observation is that  $\sigma_{\text{Standard joint}} \approx \sigma_{\text{LPMT}}$ . Once the covariance matrix between the  
 2286 LPMT and SPMT is correctly introduced, the fit “understand” that it does not have supplementary  
 2287 information and the LPMT system, which have the best precision, dominate the resolution.

|   | $\Delta m_{21}^2$ error | $\delta \Delta m_{21}^2$ error | $\theta_{12}$ error | $\delta \theta_{12}$ error | $\Delta m_{31}^2$ error | $\chi^2$    |
|---|-------------------------|--------------------------------|---------------------|----------------------------|-------------------------|-------------|
| LPMT  | 1.29936e-07             |                                | 1.33852e-03         |                            | 4.39399e-06             | 3.23088e-18 |
| SPMT  | 1.38297e-07             |                                | 1.38653e-03         |                            |                         | 2.87502e-18 |
| Indep Standard joint                          | 9.48731e-08             |                                | 9.86765e-04         |                            | 4.39212e-06             | 6.10592e-18 |
| Standard joint                                | 1.29723e-07             |                                | 1.18342e-03         |                            | 4.39287e-06             | 3.38055e-18 |
| Weighted                                      | 9.46966e-08             |                                | 9.63002e-04         |                            |                         |             |
| Delta joint                                   | 1.35780e-07             | 3.43529e-08                    | 1.38236e-03         | 1.46865e-04                | 4.39309e-06             | 3.38055e-18 |
| Indep Delta joint                             | 1.38297e-07             | 1.89391e-07                    | 1.38653e-03         | 1.87830e-03                | 4.39241e-06             | 6.10592e-18 |
| Fixed $\Delta m_{21}^2$ and $\Delta m_{31}^2$ |                         |                                |                     |                            |                         |             |
| Indep Standard joint                          |                         |                                | 9.33082e-04         |                            |                         | 4.82955e-26 |
| LPMT  |                         |                                | 1.27032e-03         |                            |                         | 2.58849e-26 |
| SPMT  |                         |                                | 1.31070e-03         |                            |                         | 2.24106e-26 |
| Weighted                                      |                         |                                | 9.12193e-04         |                            |                         |             |
| Fixed $\Delta m_{31}^2$ and $\theta_{12}$     |                         |                                |                     |                            |                         |             |
| Indep Standard joint                          | 8.97117e-08             |                                |                     |                            |                         | 6.10617e-18 |
| SPMT  | 1.30734e-07             |                                |                     |                            |                         | 2.87522e-18 |
| LPMT  | 1.23319e-07             |                                |                     |                            |                         | 3.23095e-18 |
| Weighted                                      | 8.97066e-08             |                                |                     |                            |                         |             |

TABLE 7.2 – Results of the Asimov studies on the updated framework. All results are Asimov fit, considering 6 years exposure,  $\theta_{13}$  is fixed to nominal value,  $\chi^2$  is pearson meaning that he error is estimated using the data spectrum

Finally for the *Delta* fit, the error on  $\delta\theta_{12}$  and  $\delta\Delta m_{21}^2$  are of the same order of magnitude than the errors on  $\theta_{12}$  and  $\Delta m_{21}^2$  in the absence of the covariance matrix. As the LPMT and SPMT spectra are not connected through the covariance matrix, the delta parameters are unconstrained thus the similar errors. Once the covariance matrix is introduced, the delta are much more constrained and show errors of an order of magnitude smaller than the error on their respective parameters.

Overall, the asimov studies are satisfactory. The joint fit behave as expected and the errors on the delta parameters are significantly smaller than the error on their respective parameters, indicating great potential if they converge to value too far from 0.

### Toy studies

Once we validated that the asimov study is yielding coherent results, we study the behaviour of toy studies. The above asimov study was using the Pearson  $\chi^2$  (Eq. 7.13) without pull parameter. We show in figure 7.6 the effect of using a simple Pearson  $\chi^2$ . We see that  $\sin^2(2\theta_{12})$  (reported as  $\theta_{12}$  for simplicity) is biased of about  $0.5\sigma$  and  $\Delta m_{21}^2$  biased of about  $0.1\sigma$ . When introducing the PearsonV  $\chi^2$  (Eq. 7.17) the bias disappear as reported in figure 7.7.

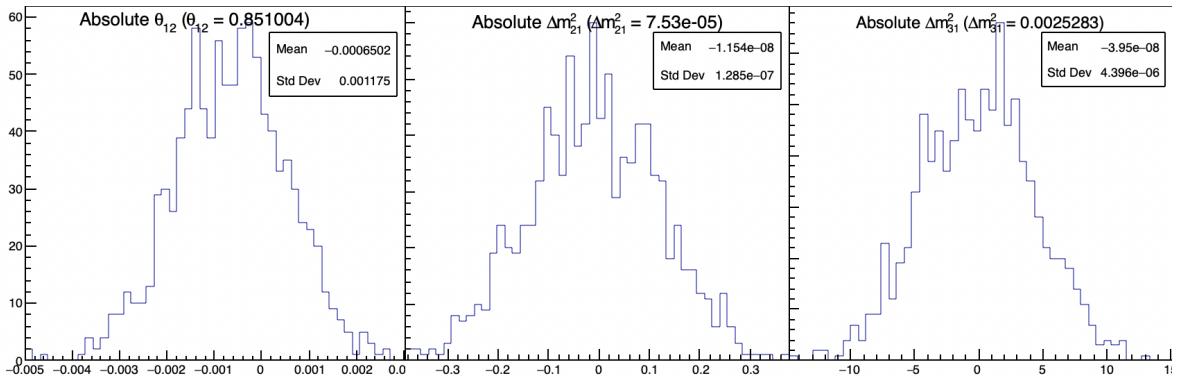


FIGURE 7.6 – Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, Pearson  $\chi^2$ ,  $\theta_{13}$  fixed.

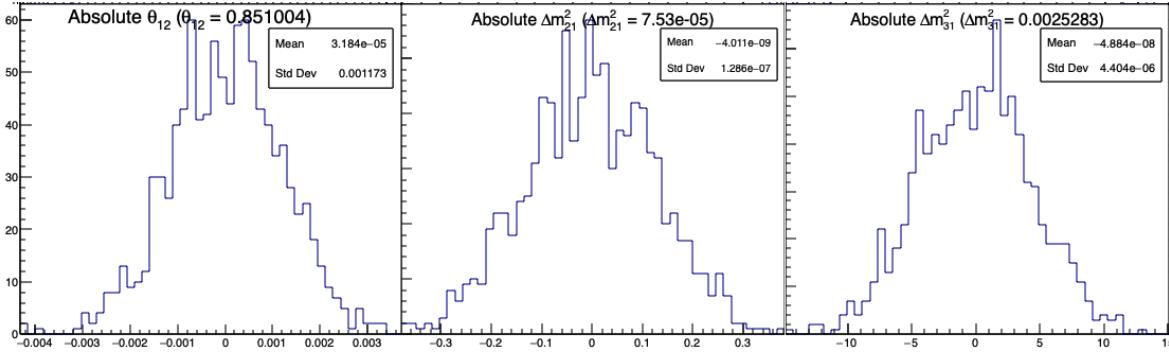


FIGURE 7.7 – Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, PearsonV  $\chi^2$ ,  $\theta_{13}$  fixed.

When the supplementary parameters are introduced in the Delta Joint fit, the fit is stable as shown in the results figure 7.8. The resolutions on the oscillation parameters are slightly worse in the Delta joint fit due to the supplementary freedom. As seen in the asimov studies, the resolution of the  $\delta$  parameters is an order of magnitude smaller than their respective parameters, indicating that they can be powerful tools to detect discrepancies between the SPMT and LPMT spectra.

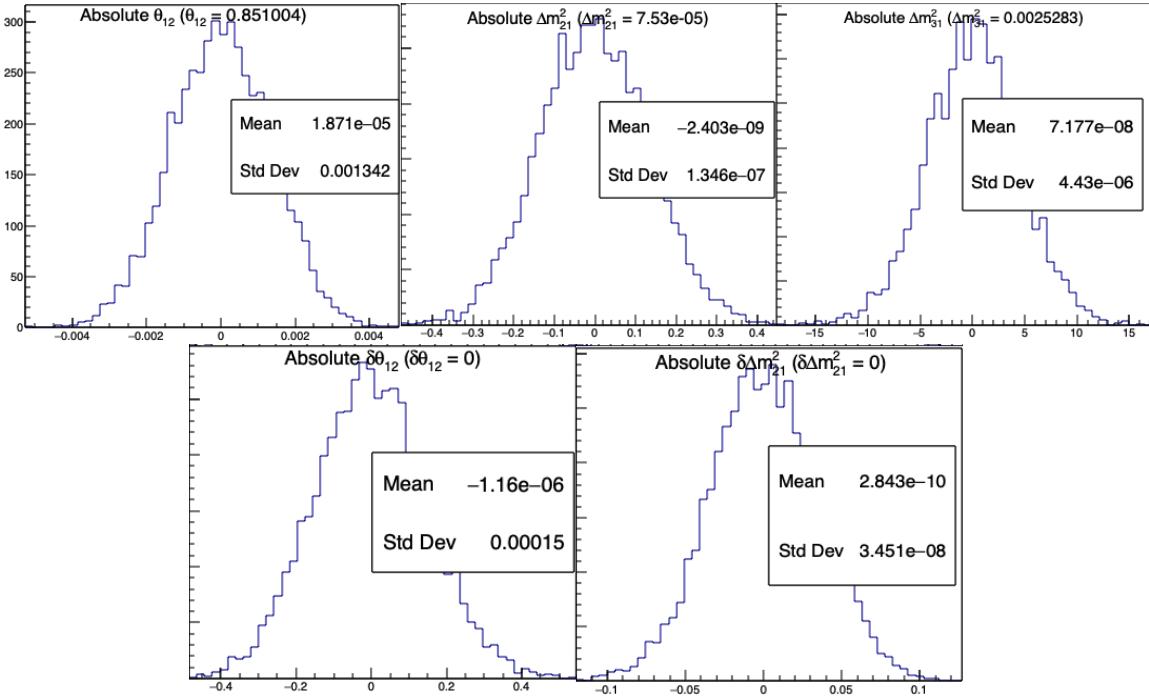


FIGURE 7.8 – Distribution of BFP - nominal value for 5000 toy Delta joint fit. 6 years exposure, all background, PearsonV  $\chi^2$ ,  $\theta_{13}$  fixed.

### Effect of supplementary QNL on the LPMT spectrum

Now that we know that the framework and joint fit behave correctly on unbiased data, we test the effect of introducing the QNL, as presented in Eq. 7.2, in the LPMT spectrum. To test the effect, we consider a QNL  $\alpha_{qnl} = 1\%$ . For reference, this is about three time the expected residual QNL after

calibration ( $\alpha_{qnl} = 0.3\%$  [29]). The background had to be removed as JUNO provide them already smeared, thus the introduction of supplementary QNL is not trivial, the resolution being dependent of  $E_{vis}$  which is affected by the QNL. We use a covariance matrix assuming no QNL. The effect of this QNL on the spectrum is illustrated in figure 7.9. In table 7.3 we report the results of the different scenarios.

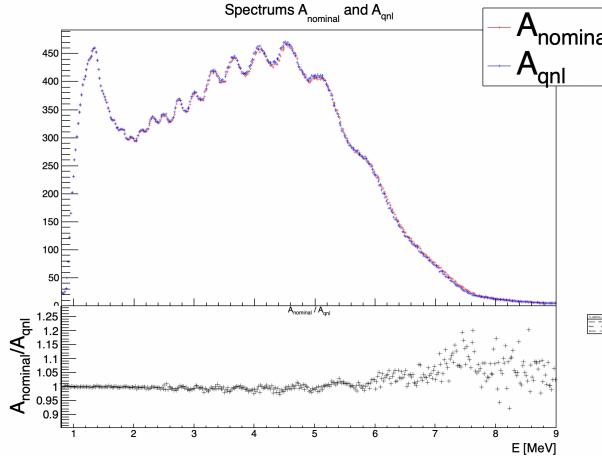


FIGURE 7.9 – **Top:** Theoretical spectrum without QNL (in red) and with  $\alpha_{qnl} = 1\%$  (in blue). **Bottom:** Ratio between the theoretical spectrum with and without QNL.

| Mean (std dev) | $\theta_{12} [10^{-3}]$ | $\Delta m_{21}^2 [10^{-7}\text{eV}^2]$ | $\Delta m_{31}^2 [10^{-6}\text{eV}^2]$ | $\delta\theta_{12} [10^{-3}]$ | $\delta\Delta m_{21}^2 [10^{-7}\text{eV}^2]$ |
|----------------|-------------------------|--|--|-------------------------------|--|
| LPMT           | -1.569 (1.171)          | -0.957 (0.989)                         | -8.235 (3.898)                         | Irrelevant                    | Irrelevant                                   |
| SPMT           | -0.164 (1.191)          | -0.603 (1.054)                         | Not sensitive                          | Irrelevant                    | Irrelevant                                   |
| Indep Standard | -0.880 (1.174)          | -0.786 (1.004)                         | -8.195 (3.900)                         | Irrelevant                    | Irrelevant                                   |
| Standard       | -8.106 (1.423)          | -2.483 (1.018)                         | -6.649 (4.008)                         | Irrelevant                    | Irrelevant                                   |
| Indep Delta    | -0.169 (1.190)          | -0.598 (1.054)                         | -8.234 (3.899)                         | -1.397 (0.259)                | -0.361 (0.366)                               |
| Delta          | -0.163 (1.183)          | -1.532 (1.036)                         | -8.193 (3.934)                         | -1.441 (0.193)                | 0.654 (0.303)                                |

TABLE 7.3 – Results of the different fit scenarios on QNL distorted data  $\alpha_{qnl} = 1\%$ . The mean value are reported subtracted from their nominal value. For SPMT  $\Delta m_{31}^2$  is fixed at nominal value. The  $\chi^2$  is PearsonV. The correlation matrix used to fit assume no QNL in the spectrum.

The results in table 7.3 are subtracted from their nominal value, themselves reported in table 7.1. We clearly see the bias induced by  $\alpha_{qnl} = 1\%$  when comparing the SPMT and LPMT results. The Indep Standard is, as expected, the mean value between the SPMT and LPMT: the fit having no informations about the correlation between the spectrum think it have two uncorrelated experiments thus report an in between value. When introducing the relationship between the LPMT and SPMT spectra in the Standard fit, the joint fit cannot find a clean minima, it thus converge to a completely incorrect value.

Introducing the  $\delta$  without the correlation in Delta Indep remove the bias and converge to the SPMT minima, the  $\delta$  absorbing the deformation of the LPMT spectra.

Finally, with the  $\delta$  and the covariance matrix,  $\theta_{12}$  is unbiased,  $\delta\theta_{12}$  absorbing the deformation.  $\delta\Delta m_{21}^2$  is still heavily biased, even more than LPMT only, for the same reason than the Standard fit: the correlation make it difficult to converge to the nominal value.

Overall  $\Delta m_{31}^2$  bias is unchanged as the SPMT spectrum bring no information about the parameter. The  $\delta$  are significant, naively up to  $7.46\sigma$  for  $\delta\theta_{12}$  in the Delta fit.

### 2330 7.5.2 Covariance matrix

2331 The covariance matrix between the LPMT and SPMT spectra is at the heart of this study as it  
 2332 was already mentioned in section 7.2 and demonstrated in section 7.5.1. In this section we discuss  
 2333 the different approaches taken to estimate it. In this work we will mainly discuss the statistical  
 2334 covariance matrix between the two spectra, how the number of event in a LPMT bin influence the  
 2335 number of bin in the SPMT spectrum due to the resolution. We will still discuss the reconstruction  
 2336 effects, mostly due to non-uniformity, in on reconstruction correlation.

2337 **Analytical method**

2338 The first method discussed is the analytical method where we propagate the resolution of the LPMT  
 2339 and SPMT spectra over a non-smeared spectrum. Following the approach used in the IBD generation  
 2340 in section 7.3.1, we consider the system resolution  $\sigma(E)$  to be only dependent in energy. We do not  
 2341 consider the position of the event.

2342 The first step is to compute the statistical uncertainty of the input spectrum while taking into account  
 2343 the smearing, considering no uncertainty on the smearing. For this, using the notation of  
 2344 section 39.2.5 *Propagation of errors* of PDG2020 [16] and considering an extended spectrum of 820 bins  
 2345 following the binning scheme introduced in 7.2.4, the first 410 for the LPMT and the last 410, we  
 2346 consider

2347 —  $\theta = (\theta_0, \dots, \theta_n)$ ;  $n = 820$  the content of the spectrum bins.

2348 —  $\eta(\theta) = (\eta_0(\theta), \dots, \eta_m(\theta))$ ;  $m = 820$  the set of smearing functions representing the PMT resolutions.

2349  $\eta_m$  can thus be defined as

$$\eta_i = \sum_j^n G(i, \sigma(E_i))(j) \theta_j \quad (7.25)$$

2350 where  $G(i, \sigma(E_i))(j)$  is the smearing function defined as

$$G(i, \sigma(E_i))(j) = \int_{\lfloor E_i \rfloor}^{\lceil E_i \rceil} \frac{1}{\sigma(E_i)\sqrt{2\pi}} e^{-\frac{(E_i-E)^2}{2\sigma(E_i)^2}} dE \quad (7.26)$$

2352 where  $E_i$  is the mean energy in the bin  $i$  and  $\lfloor E_i \rfloor$  and  $\lceil E_i \rceil$  are the lower and higher energy bound of  
 2353 the  $i$ th bin respectively.

2354 We can then construct the transfer matrix  $A$  as

$$A_{ij} = \frac{\partial \eta_i}{\partial \theta_j} = G(i, \sigma(E_i))(j) \quad (7.27)$$

2355 and then compute the first part of our covariance matrix

$$U = A V A^T \quad (7.28)$$

2356 where  $V$  is the uncorrelated covariance matrix simply defined, under the assumption of poissonian  
 2357 statistic for the bin content,

$$V_{ij} = \sqrt{\theta_i \theta_j} \quad (7.29)$$

2358 Now we just need to consider the uncertainty on the smearing  $\sigma \eta_i$ , considering no uncertainty on  
 2359 the unsmeared spectrum. From Eq. 7.25, the  $G(i, j) \equiv G(i, \sigma(E_i))(j)$  are considered independents  
 2360 from each other  $\forall i, j$ . This mean that this covariance matrix is diagonal, we only need  $\sigma G(i, j)$ . We  
 2361 can derive this term from two equation:

- 2362 — The term  $G(i, j)\theta_j$  represent the number of event smeared from the bin  $j$  that end up in the bin  
 2363  $i$ . This is a number, we thus assume poissonian statistic so that  $\sigma[G(i, j)\theta_j] = \sqrt{G(i, j)\theta_j}$ .  
 2364 — Using basic error propagation we can say that  $\sigma^2[G(i, j)\theta_j] = \theta_j^2\sigma^2G(i, j) + G(i, j)^2\sigma^2\theta_j$ .

Using  $\sigma\theta_j = \sqrt{\theta_j}$  we derive

$$G(i, j)\theta_j = \sigma^2[G(i, j)\theta_j] = \theta_j^2\sigma^2G(i, j) + G(i, j)^2\theta_j \quad (7.30)$$

$$\Rightarrow \sigma^2G(i, j) = \frac{G(i, j)\theta_j - G(i, j)^2\theta_j}{\theta_j^2} \quad (7.31)$$

$$= \frac{(1 - G(i, j))G(i, j)}{\theta_j} \quad (7.32)$$

2365 By summing the two covariance matrix, we can extract a correlation matrix presented in figure 7.10.  
 2366 The correlation between the SPMT and LPMT spectra is greater at the start of the spectrum, where  
 2367 the absolute smearing is the smallest, up to 5% correlation, and diffuse as the bins are further from  
 2368 each other and the absolute resolution grow.

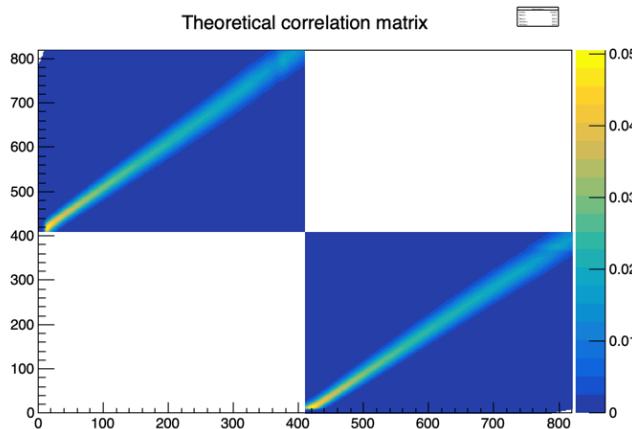


FIGURE 7.10 – Theoretical correlation matrix between the LPMT spectrum (bins 0-409) and the SPMT spectrum (410-819). The diagonal has been set to 0 (it was 1) for readability purpose.

### 2369 Empiric method

2370 The second method is the empiric way where we generate toys and just compute the empirical  
 2371 correlation between the bin contents.

$$\text{Corr}(\theta_i, \theta_j) = \frac{\mathbb{E}[\theta_i\theta_j] - \mathbb{E}[\theta_i]\mathbb{E}[\theta_j]}{\sigma\theta_i\sigma\theta_j} \quad (7.33)$$

2372 We thus generate  $10^7$  event using the IBD generator presented in section 7.3.1, then produce spectra  
 2373 from this finite set of events, meaning we must choose a number  $N$  of toy each composed of  $M$  event  
 2374 in order to have the best estimate.

2375 Due to the nature of our estimator, the estimated correlation coefficient is subject to statistical fluctuation  
 2376 as any estimator. There is no definite formula to compute the standard deviation of the correlation coefficient as suggested in this study [86] but all cited formula depend solely on the  
 2377

number of samples, in our case the number of toy  $N$ , and the correlation coefficient. This indicate that maximizing the number of toy is the right decision, even if each toy posses only one sole event.

To study this rather counter intuitive observation (How can a spectrum with only one event can be representative of the experiment ?), I present in figure 7.11 the upper left corner of the estimated correlation matrix for different configurations of  $N$  and  $M$  in the limit of  $10^7$  total event. We see in figure 7.11a that if the toy number  $N$  is too low, the statistical noise make the correlation pattern almost completely disappear, in figure 7.11b we see clearly the same correlation patter as in the theoretical matrix in figure 7.10. On the final matrix in figure 7.11c the pattern is clearly visible, but we see a shade of anti-correlation around the spectrum that was not present in the theoretical correlation matrix.

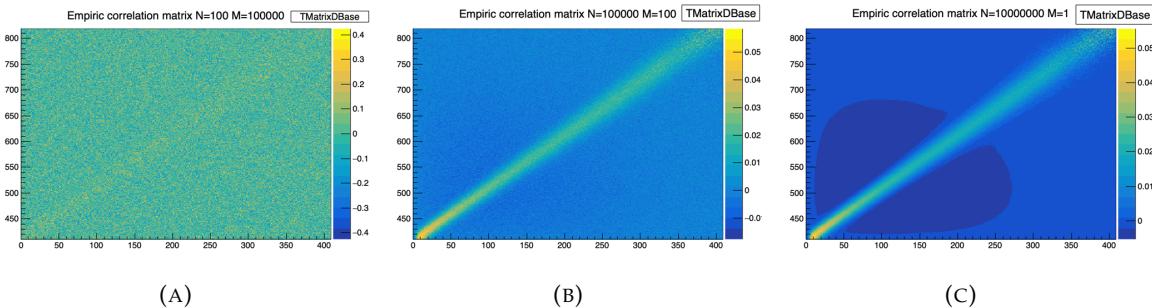


FIGURE 7.11 – Upper left corner of the estimated correlation matrix between the LPMT and SPMT spectrum for different configuration of  $N$  toy with different number of  $M$  events per toy

The difference between the element of the theoretical and the empiric correlation matrices are presented in figure 7.12a. We that the difference between the two is very small with a bias of  $1.8 \cdot 10^{-3}$  and a standard deviation of  $1.9 \cdot 10^{-3}$  while the interesting correlation are of the order  $10^{-2}$ . As presented in figure 7.12b, the most extreme differences comes from the low end of the spectrum.

This low energy difference could be explained as the theoretical does not take into account event that would be smeared from outside the spectrum.  $E < 0.8$ , MeV back inside the spectrum thus missing on the potential correlations.

The second major difference between the empirical and theoretical correlation matrices is the anti-correlation of magnitude  $\approx -5 \cdot 10^{-3}$  around the spectrum. In the theoretical correlation matrix, we assume that  $G(i, j)$  is uncorrelated from  $G(i, k)$  but this is not true in the case of a finite dataset.  $G(i, j)$  represent the number of events that migrate from the bin  $i$  to  $j$ , in the case of a finite number of event to distribute between the bins, the number of event that can be distributed in the bin  $k$  is constrained by the number of event distributed in the bin  $j$  leading to the anti-correlation between this two bins.

These empirical correlation matrices still pose an issue: These matrices needs to be invertible for  $\chi^2$  calculation. The framework use the Cholesky decomposition [87] for this, requiring the correlation matrices to be positive definite, which is not guarantee using this empirical methods. Due to this issue, the theoretical matrix is used in the studies presented in this thesis.

#### Empirical correlation matrix from fully simulated event

The last study on the correlation matrix between the LPMT and SPMT spectrum consists in simulating and reconstructing full events in the official JUNO simulation framework and computing an empirical matrix based on those events.

The core of the idea is that the LPMT and SPMT reconstruction errors is bound to be correlated due to systematic effects. The first and most obvious one, for example, is energy escaping from the central

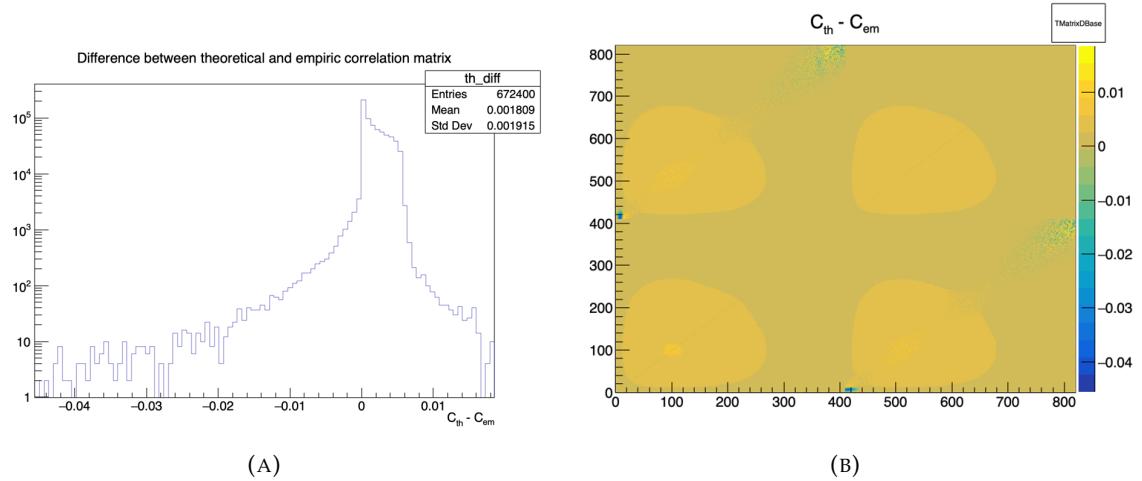


FIGURE 7.12 – Difference between the element of the theoretical and empiric correlation matrix

detector. If the positron, or one of the two annihilation gamma, escape from the detector, less energy is deposited thus both of the systems will reconstruct a lower energy than was actually deposited. On a more subtle scale, the randomness in the production of scintillation photons is common for the two systems, if the liquid scintillator produces fewer scintillation photons for an event, both systems are likely to underestimate the energy.

We study those effects by computing from a dataset of IBD events, uniformly distributed in the CD, the correlation between the reconstruction errors on the energy

$$\text{Corr}(E_{lpmt} - E_{dep}, E_{spmt} - E_{dep}) \quad (7.34)$$

where  $E_{lpmt}$  and  $E_{spmt}$  are the reconstructed energies from both systems and  $E_{dep}$  is the deposited energy in the detector.

With this observable, the bias difference between the two reconstructions at fixed  $R$  and  $E$  is irrelevant. However, since we compute the correlation in  $E$  and  $R^3$  bins, we need to account for the potential spurious relationship between the errors and their respective biases. If the bias is small relative to the resolution, it can be ignored; but if the bias variation is on the same order of magnitude as the error, it may introduce false correlations. For this reason, based on the CNN results shown in figure 4.8, we restrict our analysis to the  $1 < E_{dep} < 9$  MeV range.

The results of those correlations are presented in figure 7.13 for the single energy and radius dependency and figure 7.14 for the dual energy and radius dependency.

We see correlation increase with respect to the energy which can be attributed to the signal over dark noise ratio. As more PMTs hits come from the signal, the reconstruction becomes more signal related. Regarding the  $R^3$  distribution, we see almost no dependency until the total reflection area. After this point the correlation rises as the event are exposed to the optical effect of the total reflection area.

By looking at figure 7.14, we can see that the rising in correlation with respect to the energy is mostly due to the radius dependency.

The exploitation of those correlations in the fit and the data production, without generating and reconstructing full spectra from SNIPER, is a bit more complicated. As seen in section 7.3.1, we characterize the resolution of both systems by the ABC parameters. The correlation shown here take into account all of the ABC terms, as they are the complete correlation between the two systems, but the generation and the modeling this correlation needs to be very well understood as, as seen before,

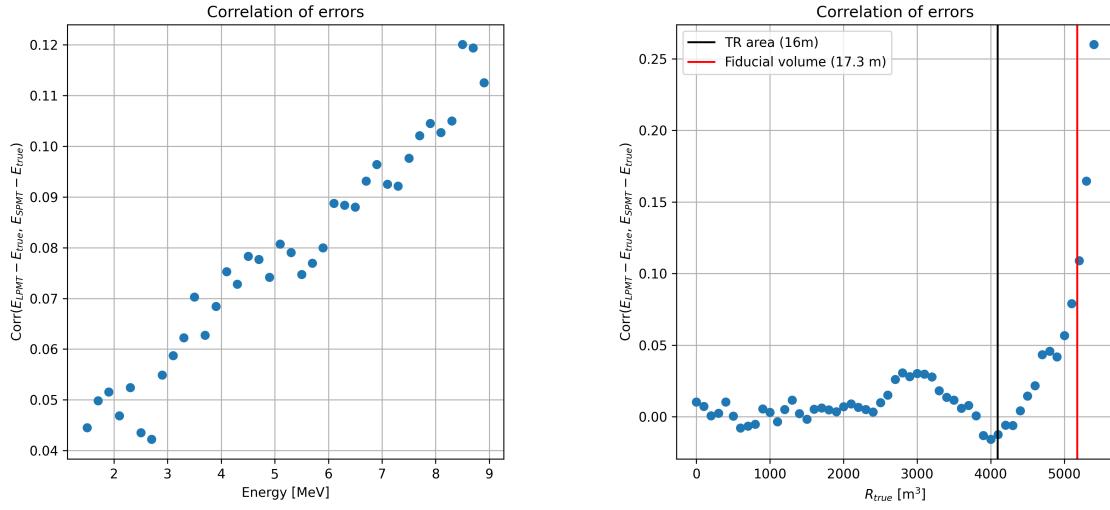


FIGURE 7.13 – Correlation on the reconstruction error between the LPMT and SPMT system as a function of (On the left) the energy, (On the right) the radius. The SPMT reconstruction comes from the NN presented in chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.6. To prevent effect due to the CNN bad reconstruction, we select the event with  $1 < E_{dep} < 9$  MeV.

2439 the mass ordering and parameters measurements are very sensitive to even small correlations.  
 2440 We consider the binned approach that we used here, knowing that the CNN reconstruction was  
 2441 deemed efficient but flawed, to be insufficient for the complete study of those effects on the fit.

### 2442 7.5.3 Statistical tests

2443 In this part, I present the results of the statistical tests presented in section 7.2.

#### 2444 Test $\chi^2_{spe}$

2445 The  $\chi^2_{spe}$  is a chi-square representing the compatibility between the LPMT ans SPMT spectra under  
 2446 constraints of the correlation matrix between the two.

$$2447 \quad \chi^2_{spe} = \Delta h V_{spe} \Delta h^T; \Delta h = \{(h_0^L - h_0^S), \dots, (h_n^L - h_n^S)\} \quad (7.35)$$

2448 where  $h_i^L$  and  $h_i^S$  are the contents of the  $i$ th bins of the LPMT and SPMT spectra. For details about the  
 calculation of  $V_{spe}$ , see section 7.2.

2449 The results for different exposures can be found in figure 7.15. To give an idea of the significance of  
 2450 this test, we provide the median p-value for each test  $\alpha_{qnl} \neq 0$ . As expected, the power of this test  
 2451 rises as the exposure does. We see significant discrimination at 6 years for  $\alpha_{qnl} \geq 0.3\%$  where the  
 2452 p-value for  $\alpha_{qnl} = 3\%$  is  $0.005 \pm 0.0022$ .

2453 This test relies solely on the estimated covariance matrix between the two spectra, requiring no  
 2454 fitting. As a result, it is a very lightweight test that can still provide valuable indications of potential  
 2455 unknown distortions between the two spectra.

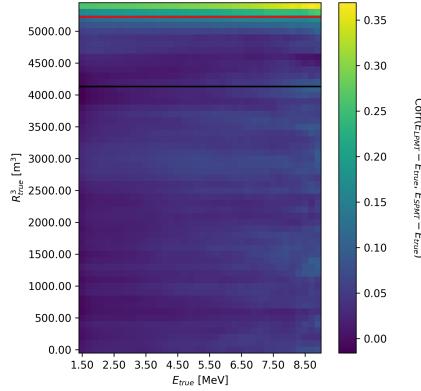


FIGURE 7.14 – Correlation on the reconstruction error between the LPMT and SPMT system as a function of the energy and the radius. The SPMT reconstruction comes from the NN presented in chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.6. To prevent effect due to the CNN bad reconstruction, we select the event with  $1 < E_{dep} < 9$  MeV.

2456 **Test  $\chi_{ind}^2$**

2457 The  $\chi_{ind}^2$  is the chi-square that represent the agreement between the measured oscillation parameters  
 2458  $\theta_{12}$  and  $\Delta m_{21}^2$ . This test is defined as

$$\chi_{ind}^2 = \Delta\lambda V_{ind} \Delta\lambda^T; \Delta\lambda = \{\theta_{12}^L - \theta_{12}^S, (\Delta m_{21}^2)^L - (\Delta m_{21}^2)^S\} \quad (7.36)$$

2459 where  $\theta_{12}^L$  and  $(\Delta m_{21}^2)^L$  are the oscillation parameters measured by the LPMT system. Same for  $\theta_{12}^S$   
 2460 and  $(\Delta m_{21}^2)^S$  for the SPMT system. We use  $V_{ind}$  computed for  $\alpha_{qnl} = 0$ . For more details about the  
 2461 calculation of  $V_{ind}$  see section 7.2.

2462 The results are presented in figure 7.16. This test does not require any joint fit or covariance matrix  
 2463 estimation between the two spectrum, it just need the estimated covariance matrix between the four  
 2464 parameters. We see that the p-value are much less significant than the other tests, this is because this  
 2465 test possess much less information about the relation between the LPMT and SPMT systems.

2466 This test is the most straightforward as it require only the fit of the two spectra and the estimation  
 2467 of the parameters covariances, but is also the less powerful with a p value for  $\alpha_{qnl} = 0.3\%$  of  $0.09 \pm$   
 2468  $0.009$ .

2469  **$\delta$  parameters significance**

2470 This test involves observing the values of the  $\delta$  parameters in the Delta Joint fit and comparing them  
 2471 tho their dispersion in the case where  $\alpha_{qnl} = 0$ . The results are shown in figures 7.17 and 7.18.

2472 We can see that the  $\delta\Delta m_{21}^2$  has a very small discriminative power (figure 7.18) even at 6 years  
 2473 exposure with a p-value of  $0.34 \pm 0.01$  for  $\alpha_{qnl} = 0.3\%$ . On the other hand  $\delta\theta_{12}$  (figure 7.17) has  
 2474 much more discriminative power with a p-value for  $\alpha_{qnl} = 0.3\%$  of  $0.025 \pm 0.005$ . This test with a  
 2475 single joint fit seems to be still less powerful than the  $\chi_{spe}^2$ . This can be explained as this method  
 2476 only get information through the oscillation parameters  $\theta_{12}$  and  $\Delta m_{21}^2$  missing potential informations  
 2477 contained in  $\Delta m_{31}^2$ .

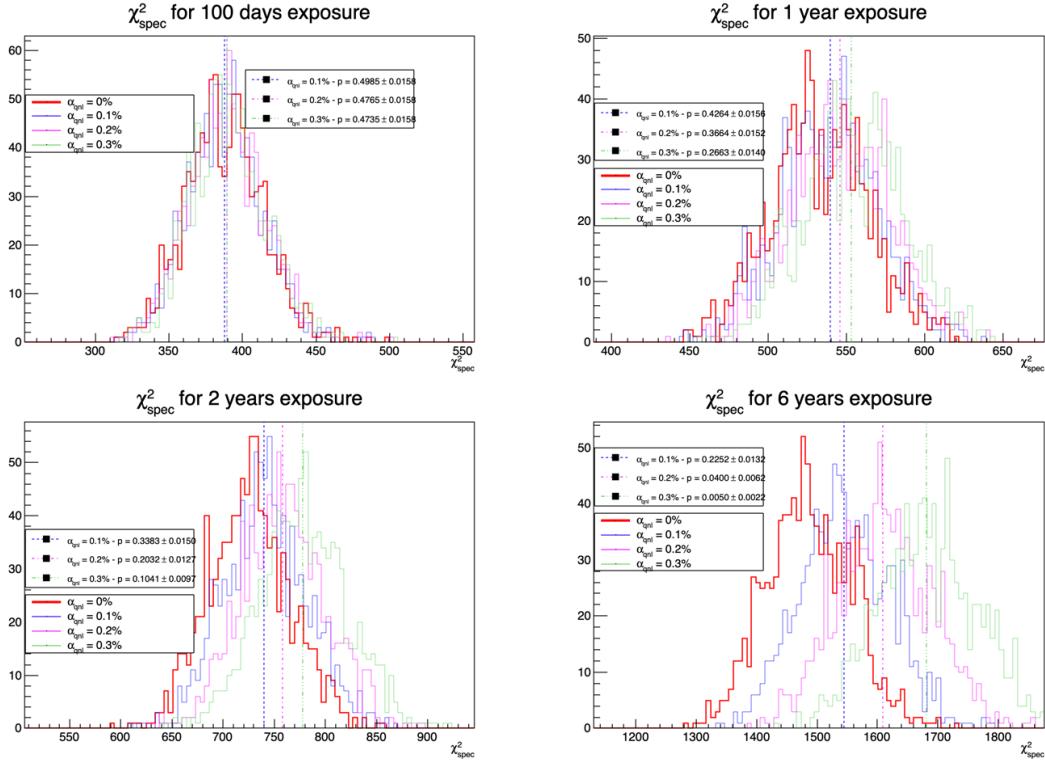


FIGURE 7.15 – Distribution of the  $\chi^2_{\text{spe}}$  for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the  $\alpha_{qnl} = 0$  distribution that are greater than those medians.

#### 2478 Hypothesis test

2479 In this last test we consider the two fit Standard Joint and Delta Joint as two hypothesis. The first  
 2480 one, Standard Joint, is the  $H_0$  hypothesis: we do not need supplementary parameters to describe the  
 2481 energy spectrum. The second one, Delta Joint, is the  $H_1$  hypothesis: we do need those supplementary  
 2482  $\delta$  parameters to, if not correctly, approach the energy spectrum. If the  $\delta$  parameter are unnecessary  
 2483 the  $\chi^2_{H_0}$  should be close to  $\chi^2_{H_1}$ . On the other hand, if one spectrum is distorted, then those parameters  
 2484 are relevant and  $\chi^2_{H_1} < \chi^2_{H_0}$ . For this test we thus observe the  $\chi^2_{H_0} - \chi^2_{H_1}$  distributions for different  
 2485 exposures and  $\alpha_{qnl}$ . The results are presented in figure 7.19.

2486 This test is the most complex, requiring two fit and the covariance matrix between the LPMT and  
 2487 SPMT spectra. The results are good, close to the  $\chi^2_{\text{spe}}$ , one with a p-value at 6 years for  $\alpha_{qnl} = 0.3\%$  of  
 2488  $0.01 \pm 0.003$ .

2489 As explained in section 7.2.4, the spectra used for the fit are cut at 335 bins / 7.5 MeV to prevent  
 2490 instability, while in  $\chi^2_{\text{spe}}$  we use full 410 bins spectra. The  $\chi^2_{\text{spe}}$  thus has more informations that the  
 2491 hypothesis test leading to this difference in power.

## 2492 7.6 Conclusion and perspectives

2493 In this chapter, we present the development of a fit framework that allows us to fit multiple spectra  
 2494 simultaneously. We also introduce a set of tools that enable us to detect potential distortions in one of  
 2495 the two spectra. As an illustration of the capability of these tools, we use supplementary event-wise

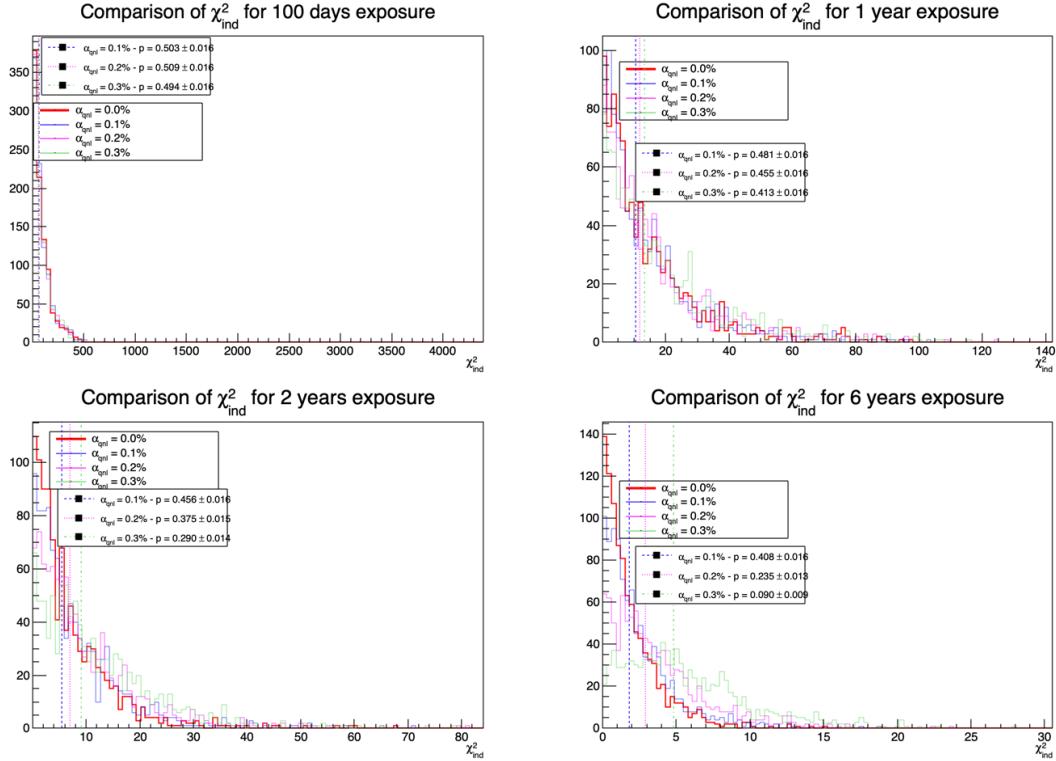


FIGURE 7.16 – Distribution of the  $\chi^2_{\text{Ind}}$  for 1000 toys for different exposures. The dashed lines represent the median of the distributions and the p-value are the percentage of the  $\alpha_{qnl} = 0$  distribution that are greater than those medians.

non-linearity and compare it to the potential residual event-wise non-linearity after calibration. Our results show that after 6 years of data collection, we can reject the median residual distortion with a p-value of 0.5% under the conditions outlined in this chapter.

Additionally, this study is preliminary, as the background was neglected in the distortion test, and no systematic uncertainties were considered. The supplementary non-linearity was introduced event-wise but should be applied channel-wise to account for the detector's non-uniformity. The correlation matrix between the LPMT and SPMT spectra should also be further analyzed, as indicated by the discrepancies between the theoretical and empirical correlation matrices. We should also further investigate the effect of non-uniformity on the correlation matrix.

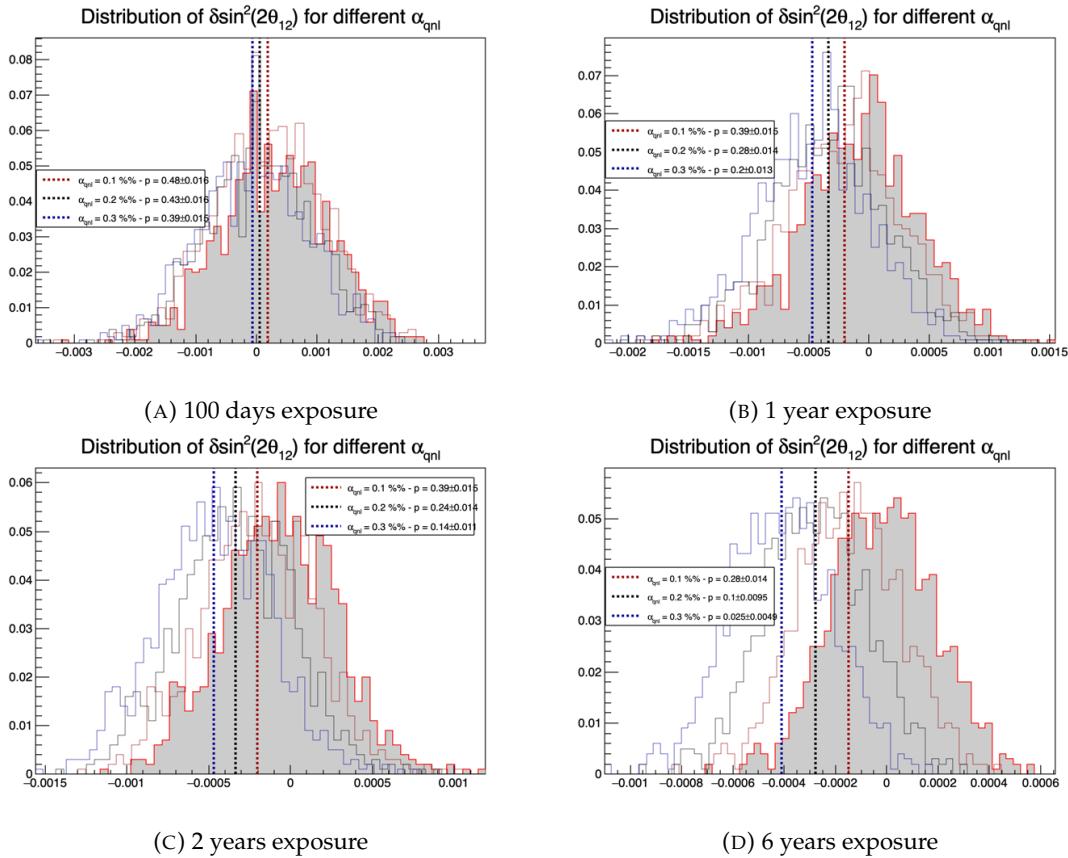


FIGURE 7.17 – Distribution of the  $\delta \sin^2(2\theta_{12})$  for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the  $\alpha_{qnl} = 0$  distribution that are greater than those medians.

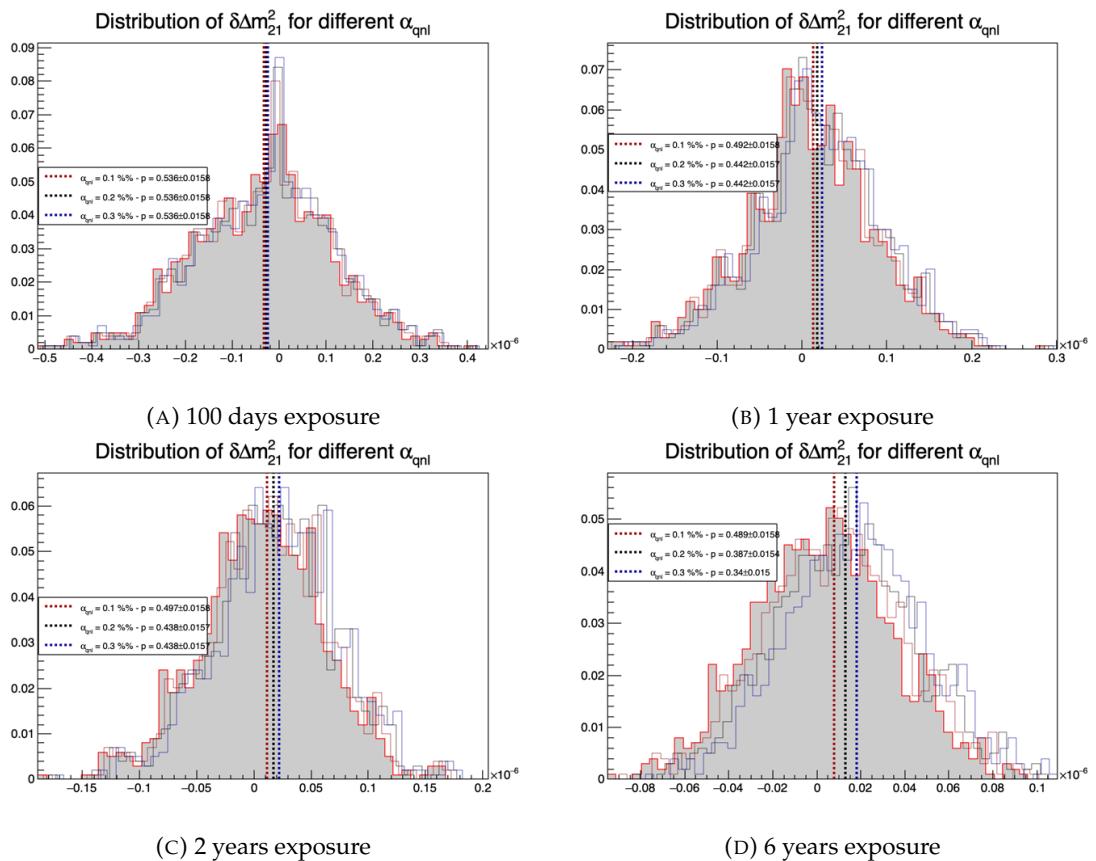


FIGURE 7.18 – Distribution of the  $\delta\Delta m_{21}^2$  for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the  $\alpha_{qnl} = 0$  distribution that are greater than those medians.

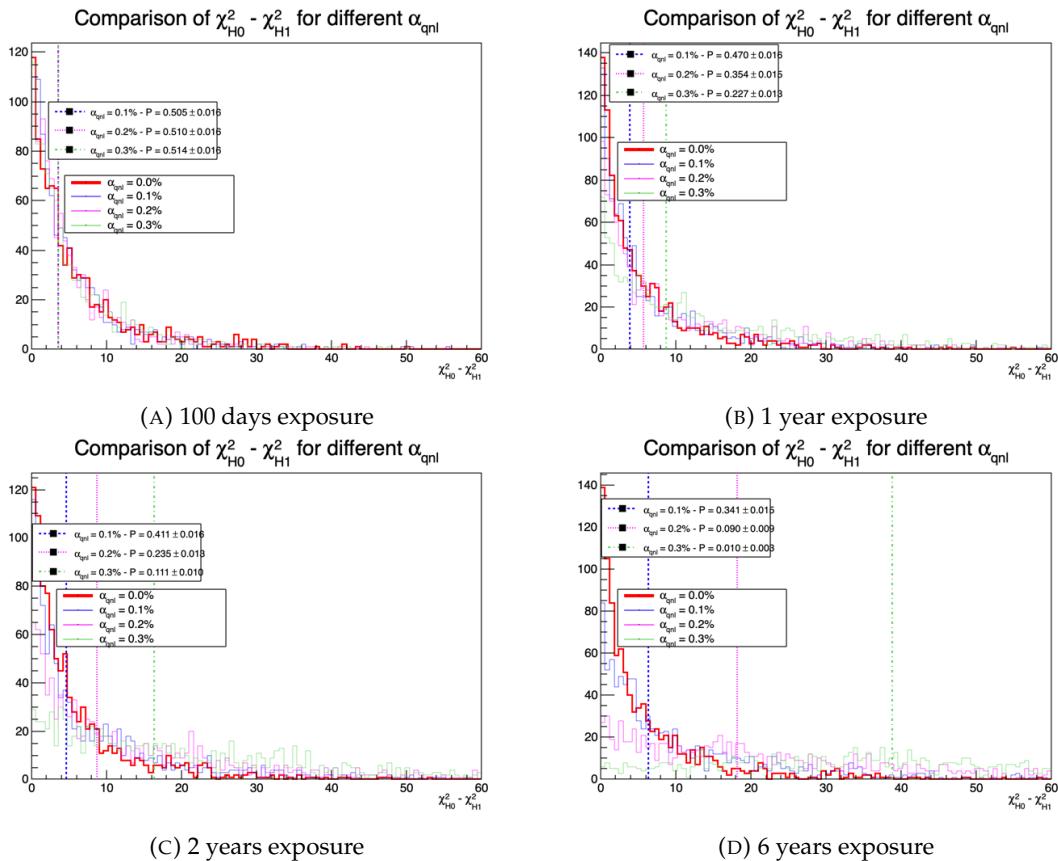


FIGURE 7.19 – Distribution of  $\chi^2_{H_0} - \chi^2_{H_1}$  for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the  $\alpha_{qnl} = 0$  distribution that are greater than those medians.

2505 **Chapter 8**

2506 **Conclusion**



<sup>2507</sup> **Appendix A**

<sup>2508</sup> **Calculation of optimal  $\alpha$  for estimator  
combination**

<sup>2509</sup>

<sup>2510</sup> This annex the details of the determination of the optimal  $\alpha$  for estimator combination presented in  
<sup>2511</sup> section 4.4.2.

<sup>2512</sup> As a reminder, the combined estimator  $\hat{\theta}$  of  $X$  is defined as

$$\hat{\theta}(X) = \alpha\theta_N + (1 - \alpha)\theta_C; \alpha \in [0; 1] \quad (\text{A.1})$$

<sup>2513</sup> where  $\theta_N$  and  $\theta_C$  are both estimator of  $X$ .

<sup>2514</sup> **A.1 Unbiased estimator**

For the unbiased estimator, it is straight-forward. We search  $\alpha$  such as  $E[\hat{\theta}] = X$

$$E[\hat{\theta}] = E[\alpha\theta_N + (1 - \alpha)\theta_C] \quad (\text{A.2})$$

$$= E[\alpha\theta_N] + E[(1 - \alpha)\theta_C] \quad (\text{A.3})$$

$$= \alpha E[\theta_N] + (1 - \alpha)E[\theta_C] \quad (\text{A.4})$$

$$= \alpha(\mu_N + X) + (1 - \alpha)(\mu_C + X) \quad (\text{A.5})$$

$$X = \alpha\mu_N + \mu_C - \alpha\mu_C + X \quad (\text{A.6})$$

$$0 = \alpha(\mu_N - \mu_C) + \mu_C \quad (\text{A.7})$$

$$(A.8)$$

$$\Rightarrow \alpha = \frac{\mu_C}{\mu_C - \mu_N} \quad (\text{A.9})$$

<sup>2515</sup> **A.2 Optimal variance estimator**

The  $\alpha$  for this estimator is a bit more tricky. By expanding the variance we get

$$\text{Var}[\hat{\theta}] = \text{Var}[\alpha\theta_N + (1 - \alpha)\theta_C] \quad (\text{A.10})$$

$$= \text{Var}[\alpha\theta_N] + \text{Var}[(1 - \alpha)\theta_C] + \text{Cov}[\alpha(1 - \alpha)\theta_N\theta_C] \quad (\text{A.11})$$

$$= \alpha^2\sigma_N^2 + (1 - \alpha)^2\sigma_C^2 + 2\alpha(1 - \alpha)\sigma_N\sigma_C\rho_{NC} \quad (\text{A.12})$$

<sup>2516</sup> where, as a reminder,  $\rho_{NC}$  is the correlation factor between  $\theta_C$  and  $\theta_N$ .

Now we try to find the minima of  $\text{Var}[\hat{\theta}]$  with respect to  $\alpha$ . For this we evaluate the derivative

$$\frac{d}{d\alpha} \text{Var}[\hat{\theta}] = 2\alpha\sigma_N^2 - 2(1-\alpha)\sigma_C^2 + 2\sigma_N\sigma_C\rho_{NC}(1-2\alpha) \quad (\text{A.13})$$

$$= 2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) - 2\sigma_C^2 + 2\sigma_N\sigma_C\rho_{NC} \quad (\text{A.14})$$

then find the minima and maxima of this derivative by evaluating

$$\frac{d}{d\alpha} \text{Var}[\hat{\theta}] = 0 \quad (\text{A.15})$$

$$2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) - 2\sigma_C^2 + 2\sigma_N\sigma_C\rho_{NC} = 0 \quad (\text{A.16})$$

$$2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) = 2\sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC} \quad (\text{A.17})$$

$$\alpha = \frac{\sigma_C^2 - \sigma_N\sigma_C\rho_{NC}}{\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}} \quad (\text{A.18})$$

2517 This equation shows only one solution which is a minima. From Eq. A.18 arise two singularities:

- 2518 —  $\sigma_N = \sigma_C = 0$ . This is not a problem because as physicists we never measure with an absolute  
2519 precision, neither us or our detectors are perfect.
- 2520 —  $\sigma_N = \sigma_C$  and  $\rho_{CN} = 1$ . In this case  $\theta_C$  and  $\theta_N$  are the same estimator in term of variance thus  
2521 any value for  $\alpha$  yield the same result: an estimator with the same variance as the original ones.

2522 **Appendix B**

2523 **Charge spherical harmonics analysis**

2524 When looking at JUNO events we can clearly see some pattern in the charge repartition based on  
2525 the event radius as illustrated in figure B.4. When dealing with identifying features and pattern on a  
2526 spherical plane, the astrophysics community have been using, with success, the spherical harmonic  
2527 decomposition. The principle is similar to a frequency analysis via Fourier transform. It comes to  
2528 saying that a function  $f(r, \theta, \phi)$ , here our charge repartition of the spherical plane constructed by our  
2529 PMTs, can be expressed

$$f(r, \theta, \phi) = \sum_{l=0}^{\infty} \sum_{m=-l}^l a_l^m r^l Y_l^m(\theta, \phi) \quad (\text{B.1})$$

2530 where  $a_l^m$  are constants complex factor,  $Y_l^m(\theta, \phi) = Ne^{im\phi} P_l^m(\cos \theta)$  are the spherical harmonics of  
2531 degree  $l$  and order  $m$  and  $P_l^m$  their associated Legendre Polynomials. Those harmonics are illustrated  
2532 in figure B.1. By reducing the problem to the unit sphere  $r = 1$ , we get rid of the term  $r^l$ . The Healpix  
2533 library [74] offer function to efficiently find the  $a_l^m$  factor from a given Healpix map.

2534 For the above decomposition, we will define the *Power* of an harmonic as

$$S_{ff}(l) = \frac{1}{2l+1} \sum_{m=-l}^l |a_l^m|^2 \quad (\text{B.2})$$

2535 and the *Relative Power* as:

$$P_l^h = \frac{S_{ff}(l)}{\sum_l S_{ff}(l)} \quad (\text{B.3})$$

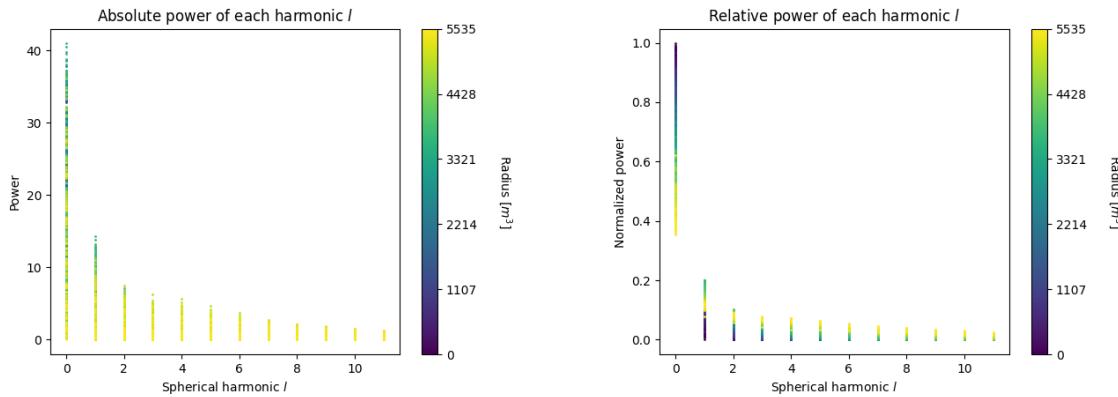
2536 For this study we will use 10k positron events with  $E_{kin} \in [0; 9]$  MeV uniformly distributed in the  
2537 CD from the JUNO official simulation version J23.0.1-rc8.dc1 (released the 7th January 2024). All the  
2538 event are *calib* level, with simulation of the physics, electronics, digitizations and triggers. We first  
2539 take a sub-set of 1k events and look at the power and relative power distribution depending on the  
2540 radius and harmonic degree  $l$ . The results are shown in figure B.2. While don't see any pattern in  
2541 absolute power, it is pretty clear that there is a correlation between the relative power of  $l = 0$  and  
2542 the radius of the event.

2543 When applying the same study but dependent on the energy, no clear correlation appear. The results  
2544 for the  $l = 0$  harmonic are presented in the figure B.5. Thus, in this study we will focus on the radial  
2545 dependency of the relative power of each harmonic.

2546 In figures B.6 and B.7 are presented the distribution of the relative power of each harmonic for  $l \in$   
2547  $[0, 11]$ . The relation between the radius and the relative power become even more clear, especially  
2548 for the first harmonics  $l \in [0, 4]$ . After that for  $l > 4$  their relative power is close to 0 for central event,  
2549 thus loosing power. It also interesting to note the change of behavior in the TR area, clearly visible  
2550 for  $l = 1$  and  $l = 2$ .

| $l:$ | $P_\ell^m(\cos \theta) \cos(m\varphi)$ | $P_\ell^{ m }(\cos \theta) \sin( m \varphi)$ |
|------|--|--|
| 0 s  |  |  |
| 1 p  |  |  |
| 2 d  |  |  |
| 3 f  |  |  |
| 4 g  |  |  |
| 5 h  |  |  |
| 6 i  |  |  |
| $m:$ | 6 5 4 3 2 1 0                          | -1 -2 -3 -4 -5 -6                            |

FIGURE B.1 – Illustration of the real part of the spherical harmonics

FIGURE B.2 – Scatter plot of the absolute and relative power, respectively on the left and right plot, of each harmonic degree  $l$ . The color indicate the radius of the event.

As an erzats of reconstruction algorithm, we fit each of those distribution with a 9th degree polynomial which give us the relation

$$F(R^3) \longmapsto P_l^h \quad (\text{B.4})$$

We do it this way because some of the distribution have multiple solution for a given relative power, for example  $l = 1$ , while each radius give only one power. We now just need to find

$$F^{-1}(P_l^h) \longmapsto R^3 \quad (\text{B.5})$$

Inverting a 9th degree polynomial is hard, if not impossible. The presence of multiple roots for the same power complexify the task even more. To circumvent this problem, we reconstruct the radius by locating the minima of  $(F(R^3) - \hat{P}_l^h)^2$  where  $\hat{P}_l^h$  is the measured power fraction.

To distinguish between multiple possible minima, we use as a starting point the radius given by the procedure on  $l = 0$  that, by looking at the fit in figure B.6, should only present one minima. For  $l > 0$  we also impose bound on the possible reconstructed  $R^3$  as  $R^3 \in [R_0^3 - 100, R_0^3 + 100]$  where  $R_0^3$  is the reconstructed  $R^3$  by the harmonic  $l = 0$ .

2562 The minimization algorithm used are the Bent algorithm for  $l = 0$  and the Bounded algorithm for  
 2563  $l > 0$  provided by the Scipy library [88]. We then do the mean of the reconstructed radius from  
 2564 the different harmonics. The reconstruction results are shown in figure B.3. The performance seems  
 2565 correct but we see heavy fluctuation in the bias. To really be used as a reconstruction algorithm, the  
 2566 method needs to be refined as discussed in the next section.

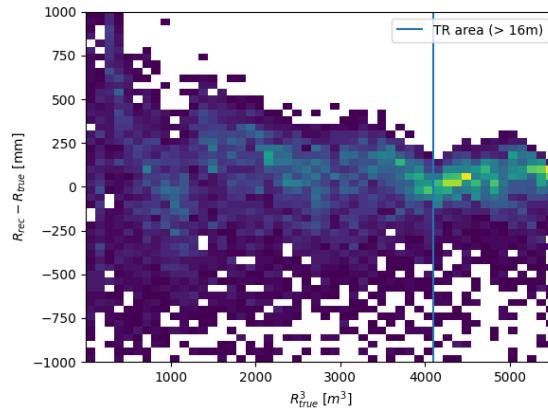


FIGURE B.3 – Error on the reconstructed radius vs the true radius by the harmonic method

## Conclusion

2568 We have clearly shown in this analysis the relevance the of relative harmonic power for radius  
 2569 reconstruction, and provided an erzats of a reconstruction algorithm. We will not delve further in  
 2570 this thesis but if we wanted to refine this algorithm multiple paths can be explored:

- 2571 — No energy signature in the harmonics: This is surprising that there is no correlation between  
 2572 the energy and the amplitude of the harmonics. We know that the energy is heavily correlated  
 2573 with the total number of photoelectrons collected, it would be unintuitive that we see no  
 2574 relation.
- 2575 — Localization of the event: We shown here the relation between the relative power of the har-  
 2576 monic and the radius but don't get any information about the  $\theta$  and  $\phi$  spherical coordinates.  
 2577 This information is probably hidden in the individual power of each order  $m$  of the degree  $l$ .  
 2578 This intuition comes from the figure B.1 where in the higher degree  $l$  we see that the order  $m$   
 2579 are oriented. Intuitively, the order should be able to indicate a direction where the signal is  
 2580 more powerful.
- 2581 — Combination of the degree power: Here we combined the radius reconstructed by the dif-  
 2582 ferent degree via a simple mean but we shown in section 4.4.2 and annex A that this is note  
 2583 the optimal way to combine estimator. A more refined algorithm probably exist to take into  
 2584 account the predicting power of each order.

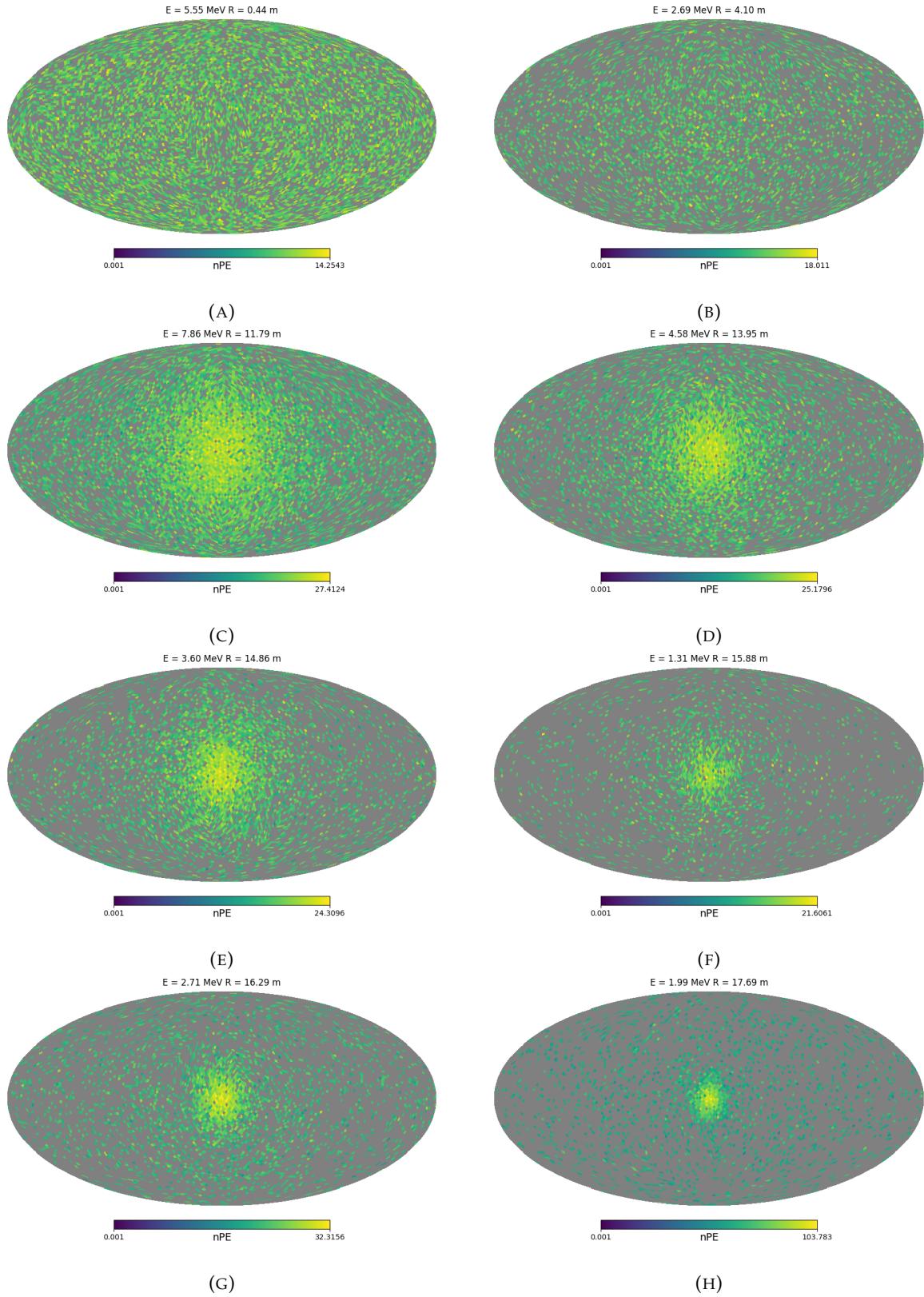


FIGURE B.4 – Charge repartition in JUNO as seen by the Healpix segmentation. Those are Healpix map of order 5 (i.e. 12288 pixels). The color represent the summed charge of the PMTs in each pixels. The color scale is logarithmic. The view have been centered to prevent event deformations.

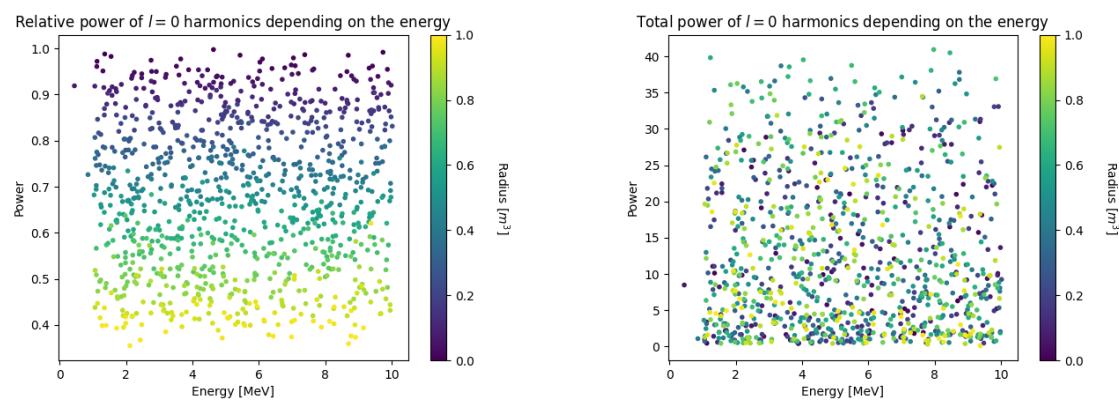


FIGURE B.5 – Scatter plot of the absolute and relative power, respectively on the left and right plot, of the  $l = 0$  harmonic. The color indicate the radius of the event.

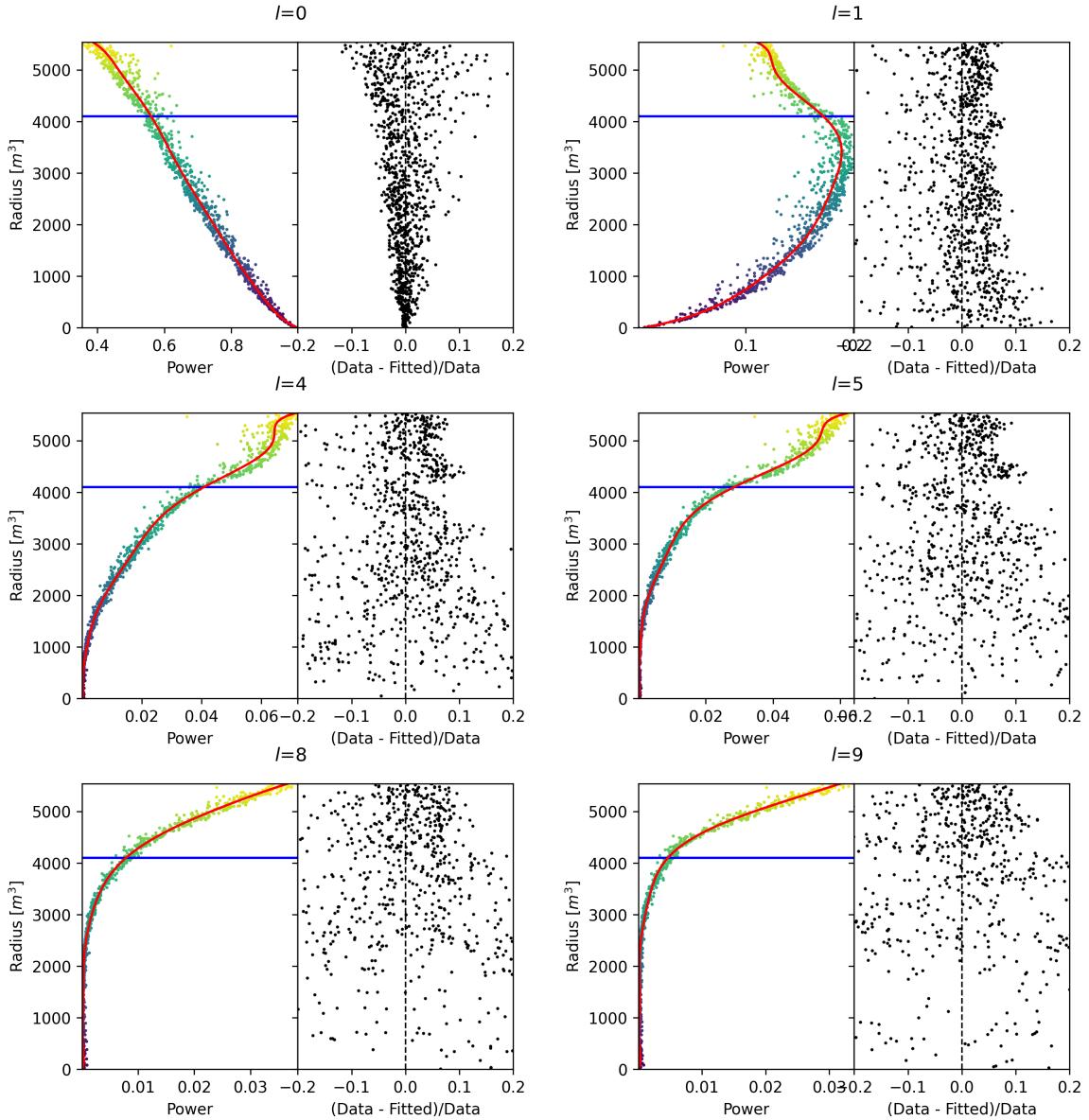


FIGURE B.6 – Plot of the distribution of the relative power of each harmonic dependent on  $R^3$  (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. **Part 1**

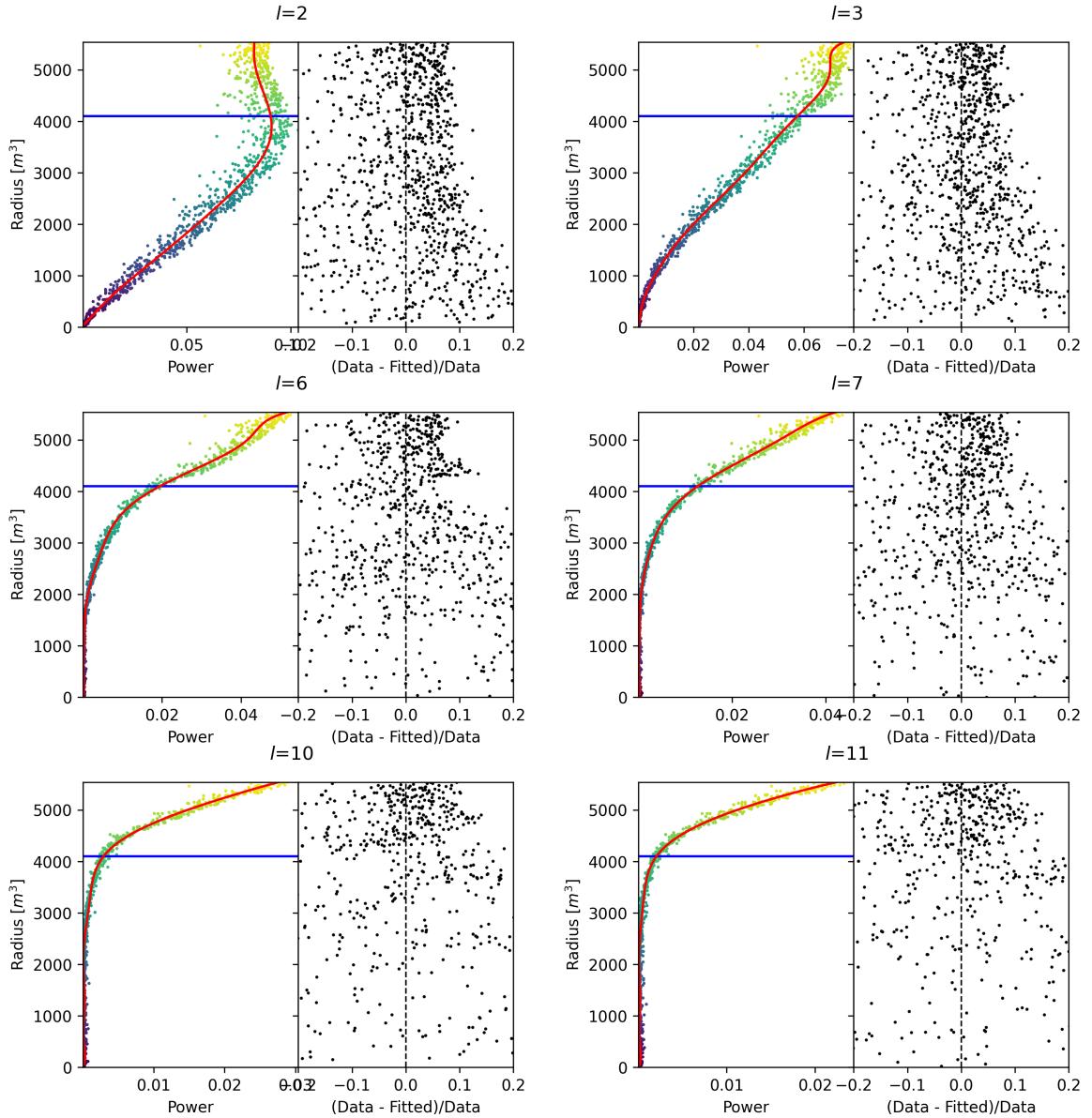


FIGURE B.7 – Plot of the distribution of the relative power of each harmonic dependent on  $R^3$  (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. **Part 2**



<sup>2585</sup> **Appendix C**

<sup>2586</sup> **Additional spectrum smearing**

<sup>2587</sup> In this section we demonstrate that a spectrum  $S$  smeared by a gaussian  $G$  parametrized by its  
<sup>2588</sup> variance  $\sigma_1^2$  can be smeared by a gaussian parametrized by the variance  $\sigma_2^2$  from the smeared spectrum  $K(E, \sigma_1) = S(E) \star G(E, \sigma_1)$  under the condition that  $\sigma_2^2 > \sigma_1^2$ .

Let  $K'(E, \sigma_2) = S(E) \star G(E, \sigma_2)$  the target spectrum we can expand

$$K'(E, \sigma_2) = S(E) \star G(E, \sigma_1) \star G^{-1}(E, \sigma_1) \star G(E, \sigma_2) \quad (\text{C.1})$$

$$= K(E, \sigma_1) \star G^{-1}(E, \sigma_1) \star G(E, \sigma_2) \quad (\text{C.2})$$

<sup>2590</sup> where  $G^{-1}(E, \sigma_1)$  is defined as  $G(E, \sigma_1) \star G^{-1}(E, \sigma_1) = \delta(E)$ .

By moving into Fourier space we can express

$$G(E, \sigma_1) \star G^{-1}(E, \sigma_1) = \delta(E) \quad (\text{C.3})$$

$$F[G(E, \sigma_1)](\nu) \times F[G^{-1}(E, \sigma_1)](\nu) = 1 \quad (\text{C.4})$$

<sup>2591</sup> with  $F[G(E, \sigma_1)](\nu)$  the fourier transform of  $G$

$$F[G(E, \sigma_1)](\nu) = e^{-\frac{\sigma_1^2(2\pi)^2}{2}\nu^2} \quad (\text{C.5})$$

we have

$$F[G^{-1}(E, \sigma_1)](\nu) = (F[G(E, \sigma_1)](\nu))^{-1} = (e^{-\frac{\sigma_1^2(2\pi)^2}{2}\nu^2})^{-1} \quad (\text{C.6})$$

$$= e^{\frac{\sigma_1^2(2\pi)^2}{2}\nu^2} \quad (\text{C.7})$$

Thus we express

$$F[G^{-1}(E, \sigma_1) \star G(E, \sigma_2)] = e^{\frac{\sigma_1^2(2\pi)^2}{2}\nu^2} \times e^{-\frac{\sigma_2^2(2\pi)^2}{2}\nu^2} \quad (\text{C.8})$$

$$= e^{\frac{(2\pi)^2}{2}(\sigma_1^2 - \sigma_2^2)\nu^2} \quad (\text{C.9})$$

$$= e^{\frac{(2\pi)^2}{2}\Delta\sigma^2\nu^2}; \Delta\sigma^2 = (\sigma_1^2 - \sigma_2^2) \quad (\text{C.10})$$

<sup>2592</sup> We see that  $F^{-1}[F[G^{-1}(E, \sigma_1) \star G(E, \sigma_2)]]$  is solvable if  $\Delta\sigma^2 = (\sigma_1^2 - \sigma_2^2) < 0 \Rightarrow \sigma_2 > \sigma_1$ . In that case

$$G^{-1}(E, \sigma_1) \star G(E, \sigma_2) = \frac{1}{\sqrt{|\Delta\sigma^2|}\sqrt{2\pi}} e^{-\frac{E^2}{2|\Delta\sigma^2|}} \quad (\text{C.11})$$



# List of Tables

|      |            |   |     |
|------|------------|---|-----|
| 2595 | <b>2.1</b> | Characteristics of the nuclear power plants observed by JUNO. . . . .   | 14  |
| 2596 | <b>2.2</b> | A summary of precision levels for the oscillation parameters. The reference value (PDG<br>2020 [16]) is compared with 100 days, 6 years and 20 years of JUNO data taking. . . . .   | 15  |
| 2597 | <b>2.3</b> | Detectable neutrino signal in JUNO and the expected signal rates and major back-<br>ground sources . . . . .  | 16  |
| 2599 | <b>2.4</b> | List of sources and their process considered for the energy scale calibration . . . . .   | 24  |
| 2600 | <b>2.5</b> | Calibration program of the JUNO experiment . . . . .  | 26  |
| 2601 | <b>2.6</b> | Features used by the BDT for vertex reconstruction . . . . .  | 36  |
| 2602 | <b>2.7</b> | Features used by the BDTE algorithm. <i>pe</i> and <i>ht</i> reference the charge and hit-time<br>distribution respectively and the percentages are the quantiles of those distributions.<br><i>cht</i> and <i>cc</i> reference the barycenters of hit time and charge respectively . . . . .                           | 37  |
| 2603 | <b>4.1</b> | Sets of hyperparameters values considered in this study . . . . .   | 58  |
| 2604 | <b>5.1</b> | Parameters of the 5th degree polynomial used to correct Omilrec reconstructed energy. .   | 81  |
| 2605 | <b>7.1</b> | Nominal PDG2020 value [16]. All value are reported assuming Normal Ordering. . . . .  | 95  |
| 2606 | <b>7.2</b> | Results of the Asimov studies on the updated framework. All results are Asimov fit,<br>considering 6 years exposure, $\theta_{13}$ is fixed to nominal value, $\chi^2$ is pearson meaning that<br>the error is estimated using the data spectrum . . . . .  | 101 |
| 2607 | <b>7.3</b> | Results of the different fit scenarios on QNL distorted data $\alpha_{qnl} = 1\%$ . The mean value<br>are reported subtracted from their nominal value. For SPMT $\Delta m_{31}^2$ is fixed at nominal<br>value. The $\chi^2$ is PearsonV. The correlation matrix used to fit assume no QNL in the<br>spectrum. . . . . | 103 |



# List of Figures

|                |   |    |
|----------------|---|----|
| 2617      2.1  | <b>On the left:</b> Location of the JUNO experiment and its reactor sources in southern china. <b>On the right:</b> Aerial view of the experimental site . . . . .  | 12 |
| 2618      2.2  | Expected number of neutrinos event per MeV in JUNO after 6 years of data taking. The black curve shows the flux if there was no oscillation. The light gray curve shows the oscillation if only the solar terms are taken in account ( $\theta_{12}$ , $\Delta m_{21}^2$ ). The blue and red curve shows the spectrum in the case of, respectively, NO and IO. The dependency of the oscillation to the different parameters are schematized by the double sided arrows. We can see the NMO sensitivity by looking at the fine phase shift between the red and the blue curve. . . . .                            | 13 |
| 2619      2.3  | Expected visible energy spectrum measured with the LPMT system with (grey) and without (black) backgrounds. The background amount for about 7% of the IBD candidate and are mostly localized below 3 MeV [11] . . . . .   | 15 |
| 2620      2.4  | a      Schematics view of the JUNO detector. . . . .  | 18 |
| 2621      2.4  | b      Top down view of the JUNO detector under construction . . . . .  | 18 |
| 2622      2.5  | Schematics of an IBD interaction in the central detector of JUNO . . . . .  | 18 |
| 2623      2.6  | Schematics of the supporting node for the acrylic vessel . . . . .  | 19 |
| 2624      2.7  | <b>On the left:</b> Quantum efficiency (QE) and emission spectrum of the LAB and the bis-MSB [20]. <b>On the right:</b> Sensitivity of the Hamamatsu LPMT depending on the wavelength of the incident photons [22]. . . . .   | 20 |
| 2625      2.8  | Schematic of a PMT . . . . .  | 21 |
| 2626      2.9  | The LPMT electronics scheme. It is composed of two part, the <i>wet</i> electronics on the left, located underwater and the <i>dry</i> electronics on the right. They are connected by Ethernet cable for data transmission and a dedicated low impedance cable for power distribution . . . . .  | 22 |
| 2627      2.10 | Schematic of the JUNO SPMT electronic system ( <b>left</b> ), and exploded view of the main component of the UWB ( <b>right</b> ) . . . . .   | 22 |
| 2628      2.11 | The JUNO top tracker . . . . .  | 23 |
| 2629      2.12 | Fitted and simulated non linearity of gamma, electron sources and from the $^{12}\text{B}$ spectrum. Black points are simulated data. Red curves are the best fits. Figures taken from [29]. . . . .  | 24 |
| 2630      2.12 | a      Gamma non-linearity . . . . .  | 24 |
| 2631      2.12 | b      Boron spectrum . . . . .   | 24 |
| 2632      2.12 | c      Electron non-linearity . . . . .   | 24 |
| 2633      2.13 | Overview of the calibration system . . . . .  | 25 |
| 2634      2.14 | Event-level instrumental non-linearity, defined as the ratio of the total measured LPMT charge to the true charge for events uniformly distributed in the detector. The solid red line represents event-level non-linearity without the channel-level correction, with position non-uniformity obtained at 1 MeV applied, in an extreme hypothetical scenario of 50% non-linearity over 100 PEs for the LPMTs. The dashed blue line represents that after the channel-level correction. The gray band shows the residual uncertainty of 0.3%, after the channel-level correction. Figure taken from [29]. . . . . | 27 |
| 2635      2.15 | . . . . .   | 28 |

|      |      |   |    |
|------|------|---|----|
| 2660 | a    | Schematic of the TAO satellite detector . . . . .   | 28 |
| 2661 | b    | Schematic of the OSIRIS satellite detector . . . . .  | 28 |
| 2662 | 2.16 | . . . . .   | 29 |
| 2663 | a    | Illustration of the different optical photons reflection scenarios. 1 is the reflection of the photon at the interface LS-acrylic or acrylic-water. 2 is the transmission of the photons through the interfaces. 3 is the conduction of the photon in the acrylic. . . . .  | 29 |
| 2664 | b    | Heatmap of $R_{rec}$ and $R_{rec} - R_{true}$ as a function of $R_{true}$ for 4MeV prompt signals uniformly distributed in the detector calculated by the charge based algorithm  | 29 |
| 2665 | 2.17 | . . . . .   | 30 |
| 2666 | a    | $\Delta t$ distribution at different iterations step $j$ . . . . .  | 30 |
| 2667 | b    | Heatmap of $R_{rec}$ and $R_{rec} - R_{true}$ as a function of $R_{true}$ for 4MeV prompt signals uniformly distributed in the detector calculated by the time based algorithm . . . . .  | 30 |
| 2668 | 2.18 | Bias of the reconstructed radius R (left), $\theta$ (middle) and $\phi$ (right) for multiple energies by the time likelihood algorithm . . . . .  | 31 |
| 2669 | 2.19 | On the left: Resolution of the reconstructed R as a function of the energy in the TR area ( $R^3 > 4000\text{m}^3 \equiv R > 16\text{m}$ ) by the charge and time likelihood algorithms. On the right: Bias of the reconstructed R in the TR area for different energies by the charge likelihood algorithm . . . . .   | 32 |
| 2670 | 2.20 | Radial resolution of the different vertex reconstruction algorithms as a function of the energy . . . . .   | 33 |
| 2671 | 2.21 | . . . . .   | 33 |
| 2672 | a    | Spherical coordinate system used in JUNO for reconstruction . . . . .   | 33 |
| 2673 | b    | Definition of the variables used in the energy reconstruction . . . . .   | 33 |
| 2674 | 2.22 | . . . . .   | 35 |
| 2675 | a    | Radial resolutions of the likelihood-based algorithm TMLE, QMLE and QTMLE . . . . .   | 35 |
| 2676 | b    | Energy resolution of QMLE and QTMLE using different vertex resolutions . . . . .  | 35 |
| 2677 | 2.23 | Projection of the LPMTs in JUNO on a 2D plane. (a) Show the distribution of all PMTs and (b) and (c) are example of what the charge and time channel looks like respectively . . . . .  | 37 |
| 2678 | 2.24 | Radial (left) and energy (right) resolutions of different ML algorithms. The results presented here are from [42]. DNN is a deep neural network, BDT is a BDT, ResNet-J and VGG-J are CNN and GNN-J is a GNN. . . . .   | 38 |
| 2679 | 3.1  | Example of a BDT that determine if the given object is a duck . . . . .   | 42 |
| 2680 | 3.2  | Schema of a simple neural network . . . . .   | 43 |
| 2681 | 3.3  | Illustration of the training lifecycle . . . . .  | 45 |
| 2682 | 3.4  | . . . . .   | 46 |
| 2683 | a    | Illustration of SGD falling into a local minima . . . . .   | 46 |
| 2684 | b    | Illustration of the Adam momentum allowing it to overcome local minima . . . . .  | 46 |
| 2685 | 3.5  | Illustration of the SGD optimizer. In blue is the value of the loss function, orange, green and red are the path taken by the optimized parameter during the training for different LR. . . . .   | 47 |
| 2686 | a    | Illustration of the SGD optimizer on one parameter $\theta$ on the MAE Loss. We see here that it has trouble reaching the minima due to the gradient being constant. . . . .  | 47 |
| 2687 | b    | Illustration of the SGD optimizer on one parameter $\theta$ on the MSE Loss. We see two different behavior: A smooth one (orange and red) when the LR is small enough and a more chaotic one when the LR is too high. . . . .   | 47 |
| 2688 | 3.6  | . . . . .   | 48 |
| 2689 | a    | Illustration of overtraining. The task at hand is to determine depending on two input variable $x$ and $y$ if the data belong to the dataset $A$ or the dataset $B$ . The expected boundary between the two dataset is represented in grey. A possible boundary learnt by overtraining is represented in brown. . . . . | 48 |
| 2690 | b    | Illustration of a very simple NN . . . . .  | 48 |

|      |      |   |    |
|------|------|---|----|
| 2712 | 3.7  | Illustration of the ResNet framework . . . . .  | 49 |
| 2713 | 3.8  | Illustration of the gradient explosion. Here it can be solved with a lower learning rate<br>but its not always the case. . . . .  | 49 |
| 2714 |      |   |    |
| 2715 | 3.9  | . . . . .   | 50 |
| 2716 | a    | Schema of a FCDNN . . . . .   | 50 |
| 2717 | b    | Illustration of a composition of ReLU “approximating” a function. (1) No ReLU<br>is taking effect (2) One ReLU is activating (3) Another ReLU is activating . . . . .   | 50 |
| 2718 |      |   |    |
| 2719 | 3.10 | Illustration of the effect of a convolution filter. Here we apply a filter with the aim<br>do detect left edges. We see in the resulting image that the left edges of the duck are<br>bright yellow where the right edges are dark blue indicating the contour of the object.<br>The convolution was calculated using [57]. . . . . | 51 |
| 2720 |      |   |    |
| 2721 | 3.11 | . . . . .   | 52 |
| 2722 | a    | Example of images in the MNIST dataset . . . . .  | 52 |
| 2723 | b    | Schema of the CNN used in Pytorch example to process the MNIST dataset . . . . .  | 52 |
| 2724 |      |   |    |
| 2725 | 3.12 | Illustration of a graph and its tensor representation. . . . .  | 53 |
| 2726 |      |   |    |
| 2727 | 3.13 | Illustration of the message passing algorithm. The detailed explanation can be found<br>in section 3.3.3 . . . . .  | 53 |
| 2728 |      |   |    |
| 2729 | 4.1  | Graphic representation of the VGG-16 architecture, presenting the different kind of<br>layer composing the architecture. . . . .  | 57 |
| 2730 |      |   |    |
| 2731 | 4.2  | Repartition of SPMTs in the image projection. The color scale is the number of SPMTs<br>per pixel . . . . .   | 61 |
| 2732 |      |   |    |
| 2733 | 4.3  | Example of a high energy, radial event. We see a concentration of the charge on the<br>bottom right of the image, clear indication of a high radius event. <b>On the left:</b> the<br>charge channel. The color is the charge in each pixel in NPE equivalent. <b>On the right:</b><br>The time channel in nanoseconds. . . . .     | 62 |
| 2734 |      |   |    |
| 2735 | 4.4  | Example of a low energy, radial event. The signal here is way less explicit, we can<br>kind of guess that the event is located in the top middle of the image. <b>On the left:</b> the<br>charge channel. The color is the charge in each pixel in NPE equivalent. <b>On the right:</b><br>The time channel in nanoseconds. . . . . | 62 |
| 2736 |      |   |    |
| 2737 | 4.5  | Example of a high energy, central event. In this image we can see a lot of signal but<br>uniformly spread, this is indicative of a central event. <b>On the left:</b> the charge channel.<br>The color is the charge in each pixel in NPE equivalent. <b>On the right:</b> The time channel<br>in nanoseconds. . . . .              | 63 |
| 2738 |      |   |    |
| 2739 | 4.6  | Example of a low energy, central event. Here there is no clear signal, the uniformity<br>of the distribution should make it central. <b>On the left:</b> the charge channel. The color<br>is the charge in each pixel in NPE equivalent. <b>On the right:</b> The time channel in<br>nanoseconds. . . . .                           | 63 |
| 2740 |      |   |    |
| 2741 | 4.7  | . . . . .   | 64 |
| 2742 | a    | Distribution of PE/MeV in the J23 Dataset. This distribution is profiled and<br>fitted using equation 4.6 . . . . .   | 64 |
| 2743 | b    | <b>On top:</b> Distribution of PE vs Energy. <b>On bottom:</b> Using the values extracted<br>in 4.7a, we calculate the ration signal over background + signal . . . . .   | 64 |
| 2744 |      |   |    |
| 2745 | 4.8  | Reconstruction performance of the “gen_30” model on J21 data and it’s comparison<br>to the performances of the classic algorithm “Classical algorithm” from [64]. The top<br>part of each plot is the resolution and the bottom part is the bias. . . . .   | 65 |
| 2746 |      |   |    |
| 2747 | a    | Resolution and bias of energy reconstruction vs energy . . . . .  | 65 |
| 2748 | b    | Resolution and bias of energy reconstruction vs radius . . . . .  | 65 |
| 2749 | c    | Resolution and bias of radius reconstruction vs energy . . . . .  | 65 |
| 2750 | d    | Resolution and bias of radius reconstruction vs radius . . . . .  | 65 |
| 2751 | e    | Resolution and bias of radius reconstruction vs $\theta$ . . . . .  | 65 |
| 2752 | f    | Resolution and bias of radius reconstruction vs $\phi$ . . . . .  | 65 |
| 2753 |      |   |    |

|      |      |   |    |
|------|------|---|----|
| 2763 | 4.9  | Error distribution of the different component of the vertex by “gen_30”. The reconstructed component are $x$ , $y$ and $z$ but we see similar behavior in the error of $R$ , $\theta$ and $\phi$ . . . . .  | 66 |
| 2764 | a    | Distribution of the error on reconstructed $x$ by “gen_30” . . . . .  | 66 |
| 2765 | b    | Distribution of the error on reconstructed $y$ by “gen_30” . . . . .  | 66 |
| 2766 | c    | Distribution of the error on reconstructed $z$ by “gen_30” . . . . .  | 66 |
| 2767 | d    | Distribution of the error on reconstructed $R$ by “gen_30” . . . . .  | 66 |
| 2768 | e    | Distribution of the error on reconstructed $\theta$ by “gen_30” . . . . .   | 66 |
| 2769 | f    | Distribution of the error on reconstructed $\phi$ by “gen_30” . . . . .   | 66 |
| 2770 | 4.10 | . . . . .   | 67 |
| 2771 | a    | Distribution of “gen_30” reconstructed energy and true energy of the analysis dataset (J21) . . . . .   | 67 |
| 2772 | b    | Distribution of “gen_42” reconstructed energy and true energy of the analysis dataset (J23) . . . . .   | 67 |
| 2773 | 4.11 | Radius bias ( <b>on the left</b> ) and resolution( <b>on the right</b> ) of the classical algorithm in a $E$ , $R^3$ grid . . . . .   | 68 |
| 2774 | 4.12 | Reconstruction performance of the “gen_30” model on J21, the classic algorithm “Classical algorithm” from [64] and the combination of both using weighted mean. The top part of each plot is the resolution and the bottom part is the bias. . . . .                                      | 69 |
| 2775 | a    | Resolution and bias of energy reconstruction vs energy . . . . .  | 69 |
| 2776 | b    | Resolution and bias of energy reconstruction vs radius . . . . .  | 69 |
| 2777 | c    | Resolution and bias of radius reconstruction vs energy . . . . .  | 69 |
| 2778 | d    | Resolution and bias of radius reconstruction vs radius . . . . .  | 69 |
| 2779 | e    | Resolution and bias of radius reconstruction vs $\theta$ . . . . .  | 69 |
| 2780 | f    | Resolution and bias of radius reconstruction vs $\phi$ . . . . .  | 69 |
| 2781 | 4.13 | Correlation between CNN and classical method reconstruction ( <b>on the left</b> ) for energy and ( <b>on the right</b> ) for radius in a $E$ , $R^3$ grid . . . . .  | 70 |
| 2782 | 4.14 | Reconstruction performance of the “gen_42” model on J23 data and it’s comparison to the performances of the classic algorithm “Classical algorithm” from [64]. The top part of each plot is the resolution and the bottom part is the bias. . . . .                                       | 71 |
| 2783 | a    | Resolution and bias of energy reconstruction vs energy . . . . .  | 71 |
| 2784 | b    | Resolution and bias of energy reconstruction vs radius . . . . .  | 71 |
| 2785 | c    | Resolution and bias of radius reconstruction vs energy . . . . .  | 71 |
| 2786 | d    | Resolution and bias of radius reconstruction vs radius . . . . .  | 71 |
| 2787 | e    | Resolution and bias of radius reconstruction vs $\theta$ . . . . .  | 71 |
| 2788 | f    | Resolution and bias of radius reconstruction vs $\phi$ . . . . .  | 71 |
| 2789 | 5.1  | . . . . .   | 75 |
| 2790 | a    | Illustration of the different nodes in our graphs and their relations . . . . .   | 75 |
| 2791 | b    | Illustration of what a dense adjacency matrix would looks like and the part we are really interested in. Because Fired $\rightarrow$ Mesh and Mesh $\rightarrow$ I/O relations are undirected, we only consider in practice the top right part of the matrix for those relations. . . . . | 75 |
| 2792 | 5.2  | Illustration of the healpix segmentation. <b>On the left:</b> A segmentation of order 0. <b>On the right:</b> A segmentation of order 1 . . . . .   | 75 |
| 2793 | 5.3  | Features held by the nodes and edges in the graph. $D_{m_1 \rightarrow m_2}^{-1}$ is the inverse of the distance between two mesh center. The features $P_l^h$ , $\mathbb{A}$ and $\mathbb{B}$ are detailed in section 5.2  | 76 |
| 2794 | 5.4  | Illustration of the different update function needed by our GNN . . . . .   | 77 |
| 2795 | 5.5  | Distribution of the number of hits depending on the energy. <b>On the right:</b> for the LPMT system. <b>In the middle :</b> for the SPMT system. <b>On the left:</b> For both system . . . . .   | 78 |
| 2796 | a    | . . . . .   | 78 |
| 2797 | b    | . . . . .   | 78 |
| 2798 | c    | . . . . .   | 78 |

|              |  |     |
|--------------|--|-----|
| 2815    5.6  | Distribution of the number of hits depending on the radius. <b>On the right:</b> for the LPMT system. <b>On the right :</b> for the SPMT system. To prevent the superposition of structure of different scales we limit ourselves to the energy range $E_{true} \in [0, 9]$ . . . . .  | 79  |
| 2816    a    | .....  | 79  |
| 2817    b    | .....  | 79  |
| 2818    5.7  | Schema of the JWGv8.4.0 architecture, the colored triplet is the graph configuration after each JWG layers . . . . .   | 80  |
| 2819    5.8  | Comparison between Omilrec $E_{rec}$ and the true energy $E_{true}$ . The profile of the distribution $E_{true}/E_{rec}$ vs $E_{rec}$ is fitted with a 5th degree polynomial. . . . .  | 82  |
| 2820    5.9  | Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.6, JWGv8.4 presented in this chapter and the combination between the two as presented in section 4.4.2. The top part of each plot is the resolution and the bottom part is the bias. . . . .  | 83  |
| 2821    a    | Resolution and bias of energy reconstruction vs energy . . . . .   | 83  |
| 2822    b    | Resolution and bias of energy reconstruction vs radius . . . . .   | 83  |
| 2823    c    | Resolution and bias of radius reconstruction vs energy . . . . .   | 83  |
| 2824    d    | Resolution and bias of radius reconstruction vs radius . . . . .   | 83  |
| 2825    e    | Resolution and bias of radius reconstruction vs $\theta$ . . . . .   | 83  |
| 2826    f    | Resolution and bias of radius reconstruction vs $\phi$ . . . . .   | 83  |
| 2827    6.1  | Schema of the method to discover vulnerabilities in the reconstruction methods . . . . .   | 87  |
| 2828    7.1  | Two oscillated spectra of $1e7$ event expected in JUNO. In red the spectrum without supplementary QNL. In blue the same spectrum but where an event-wise QNL $\alpha_{qnl} = 10\%$ is introduced. . . . .  | 91  |
| 2829    7.2  | .....  | 92  |
| 2830    a    | Distribution of ratio of collected nPE after the additional QNL over the number of nPE that would be collected for different $\gamma_{qnl}$ . We select event with an interaction radius $R < 4m$ to not be affected by the non-uniformity. . . . .  | 92  |
| 2831    b    | Ratio of collected nPE after the additional QNL over the number of nPE that would be collected at different energies. We select event with an interaction radius $R < 4m$ to not be affected by the non-uniformity. The dots represent the mean of the distributions in figure 7.2a and the dashed line are the equivalent event-wise non-linearity from eq 7.2. The hatched zone is the residual non-linearity expected after calibration [29]. . . . . | 92  |
| 2832    7.3  | Theoretical LPMT spectrum at nominal oscillation values binned using 410 bins from 0.8 to 9 MeV. It is rescaled to 6 years statistic. The black line represent the 335 bin cut . . . . .   | 96  |
| 2833    7.4  | Schematic description of the fit framework . . . . .   | 97  |
| 2834    7.5  | Relative ( <b>On the left</b> ) and absolute ( <b>On the right</b> ) resolutions of the LPMT and SPMT systems used in this study. The number in parenthesis are the parameter $A$ , $B$ and $C$ respectively for each systems. . . . .   | 98  |
| 2835    7.6  | Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, Pearson $\chi^2$ , $\theta_{13}$ fixed. . . . .   | 101 |
| 2836    7.7  | Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, PearsonV $\chi^2$ , $\theta_{13}$ fixed. . . . .  | 102 |
| 2837    7.8  | Distribution of BFP - nominal value for 5000 toy Delta joint fit. 6 years exposure, all background, PearsonV $\chi^2$ , $\theta_{13}$ fixed. . . . .   | 102 |
| 2838    7.9  | <b>Top:</b> Theoretical spectrum without QNL (in red) and with $\alpha_{qnl} = 1\%$ (in blue). <b>Bottom:</b> Ratio between the theoretical spectrum with and without QNL. . . . .   | 103 |
| 2839    7.10 | Theoretical correlation matrix between the LPMT spectrum (bins 0-409) and the SPMT spectrum (410-819). The diagonal has been set to 0 (it was 1) for readability purpose. . . . .  | 105 |
| 2840    7.11 | Upper left corner of the estimated correlation matrix between the LPMT and SPMT spectrum for different configuration of $N$ toy with different number of $M$ events per toy  | 106 |
| 2841    a    | .....  | 106 |

|   |       |     |
|---|-------|-----|
| 2867            b   | ..... | 106 |
| 2868            c   | ..... | 106 |
| 2869 <b>7.12</b> Difference between the element of the theoretical and empiric correlation matrix                           | ..... | 107 |
| 2870            a   | ..... | 107 |
| 2871            b   | ..... | 107 |
| 2872 <b>7.13</b> Correlation on the reconstruction error between the LPMT and SPMT system as a                              |       |     |
| 2873            function of (On the left) the energy, (On the right) the radius. The SPMT recon-                            |       |     |
| 2874            struction comes from the NN presented in chapter 4 and the LPMT reconstruction                              |       |     |
| 2875            comes from OMILREC presented in section 2.6. To prevent effect due to the CNN bad                           |       |     |
| 2876            reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.   | ..... | 108 |
| 2877 <b>7.14</b> Correlation on the reconstruction error between the LPMT and SPMT system as a                              |       |     |
| 2878            function of the energy and the radius. The SPMT reconstruction comes from the NN                            |       |     |
| 2879            presented in chapter 4 and the LPMT reconstruction comes from OMILREC presented                             |       |     |
| 2880            in section 2.6. To prevent effect due to the CNN bad reconstruction, we select the event                    |       |     |
| 2881            with $1 < E_{dep} < 9$ MeV.   | ..... | 109 |
| 2882 <b>7.15</b> Distribution of the $\chi^2_{spe}$ for 1000 toys for different exposure. The dashed line represent         |       |     |
| 2883            the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$                 |       |     |
| 2884            distribution that are greater than those medians.   | ..... | 110 |
| 2885 <b>7.16</b> Distribution of the $\chi^2_{ind}$ for 1000 toys for different exposures. The dashed lines repre-          |       |     |
| 2886            sent the median of the distributions and the p-value are the percentage of the $\alpha_{qnl} = 0$           |       |     |
| 2887            distribution that are greater than those medians.   | ..... | 111 |
| 2888 <b>7.17</b> Distribution of the $\delta \sin^2(2\theta_{12})$ for 1000 toys for different exposure. The dashed line    |       |     |
| 2889            represent the median of the distribution and the p-value are the percentage of the                          |       |     |
| 2890 $\alpha_{qnl} = 0$ distribution that are greater than those medians.   | ..... | 112 |
| 2891            a      100 days exposure  | ..... | 112 |
| 2892            b      1 year exposure  | ..... | 112 |
| 2893            c      2 years exposure   | ..... | 112 |
| 2894            d      6 years exposure   | ..... | 112 |
| 2895 <b>7.18</b> Distribution of the $\delta \Delta m_{21}^2$ for 1000 toys for different exposure. The dashed line repre-  |       |     |
| 2896            sent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$            |       |     |
| 2897            distribution that are greater than those medians.   | ..... | 113 |
| 2898            a      100 days exposure  | ..... | 113 |
| 2899            b      1 year exposure  | ..... | 113 |
| 2900            c      2 years exposure   | ..... | 113 |
| 2901            d      6 years exposure   | ..... | 113 |
| 2902 <b>7.19</b> Distribution of $\chi^2_{H_0} - \chi^2_{H_1}$ for 1000 toys for different exposure. The dashed line repre- |       |     |
| 2903            sent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$            |       |     |
| 2904            distribution that are greater than those medians.   | ..... | 114 |
| 2905            a      100 days exposure  | ..... | 114 |
| 2906            b      1 year exposure  | ..... | 114 |
| 2907            c      2 years exposure   | ..... | 114 |
| 2908            d      6 years exposure   | ..... | 114 |
| 2909 <b>B.1</b> Illustration of the real part of the spherical harmonics  | ..... | 120 |
| 2910 <b>B.2</b> Scatter plot of the absolute and relative power, respectively on the left and right plot,                   |       |     |
| 2911            of each harmonic degree $l$ . The color indicate the radius of the event.                                   | ..... | 120 |
| 2912 <b>B.3</b> Error on the reconstructed radius vs the true radius by the harmonic method                                 | ..... | 121 |
| 2913 <b>B.4</b> Charge repartition in JUNO as seen by the Healpix segmentation. Those are Healpix                           |       |     |
| 2914            map of order 5 (i.e. 12288 pixels). The color represent the summed charge of the PMTs                       |       |     |
| 2915            in each pixels. The color scale is logarithmic. The view have been centered to prevent                      |       |     |
| 2916            event deformations.   | ..... | 122 |
| 2917            a   | ..... | 122 |
| 2918            b   | ..... | 122 |

---

|      |     |   |     |
|------|-----|---|-----|
| 2919 | c   | .....   | 122 |
| 2920 | d   | .....   | 122 |
| 2921 | e   | .....   | 122 |
| 2922 | f   | .....   | 122 |
| 2923 | g   | .....   | 122 |
| 2924 | h   | .....   | 122 |
| 2925 | B.5 | Scatter plot of the absolute and relative power, respectively on the left and right plot, of the $l = 0$ harmonic. The color indicate the radius of the event.  | 123 |
| 2926 | B.6 | Plot of the distribution of the relative power of each harmonic dependent on $R^3$ (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. <b>Part 1</b> | 124 |
| 2927 | B.7 | Plot of the distribution of the relative power of each harmonic dependent on $R^3$ (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. <b>Part 2</b> | 125 |
| 2928 |     |   |     |
| 2929 |     |   |     |
| 2930 |     |   |     |
| 2931 |     |   |     |
| 2932 |     |   |     |
| 2933 |     |   |     |
| 2934 |     |   |     |



# List of Abbreviations

|                |   |
|----------------|---|
| <b>ACU</b>     | Automatic Calibration Unit                                      |
| <b>BDT</b>     | Boosted Decision Tree   |
| <b>BFP</b>     | Best Fit Point  |
| <b>CD</b>      | Central Detector  |
| <b>CLS</b>     | Cable Loop System   |
| <b>CNN</b>     | Convolutional NN  |
| <b>DNN</b>     | Deep NN   |
| <b>DN</b>      | Dark Noise  |
| <b>EDM</b>     | Event Data Model  |
| <b>FCDNN</b>   | Fully Connected Deep NN   |
| <b>GNN</b>     | Graph NN  |
| <b>GT</b>      | Guiding Tube  |
| <b>IBD</b>     | Inverse Beta Decay  |
| <b>IO</b>      | Inverse Ordering  |
| <b>JUNO</b>    | Jiangmen Underground Neutrino Observatory                       |
| <b>LPMT</b>    | Large PMT   |
| <b>LR</b>      | Learning Rate   |
| <b>LS</b>      | Liquid Scintillator   |
| <b>MC</b>      | Monte Carlo simulation  |
| <b>ML</b>      | Machine Learning  |
| <b>MSE</b>     | Mean Squared Error  |
| <b>NMO</b>     | Neutrino Mass Ordering  |
| <b>NN</b>      | Neural Network  |
| <b>NO</b>      | Normal Ordering   |
| <b>NPE</b>     | Number of Photo Electron  |
| <b>OSIRIS</b>  | Online Scintillator Internal Radioactivity Investigation System |
| <b>PE</b>      | Photo Electron  |
| <b>PMT</b>     | Photo-Multipliers Tubes   |
| <b>PRelu</b>   | Parametrized Rectified Linear Unit                              |
| <b>QNL</b>     | Charge (Q) Non Linearity  |
| <b>ROV</b>     | Remotely Operated under-LS Vehicle                              |
| <b>ReLU</b>    | Rectified Linear Unit   |
| <b>ResNet</b>  | Residual Network  |
| <b>SGD</b>     | Stochastic Gradient Descent                                     |
| <b>SPMT</b>    | Small PMT   |
| <b>TAO</b>     | Taishan Antineutrino Oservatory                                 |
| <b>TR Area</b> | Total Reflexion Area  |
| <b>TTS</b>     | Time Transit Spread   |
| <b>TT</b>      | Top Tracker   |
| <b>UWB</b>     | Under Water Boxes   |
| <b>WCD</b>     | Water Cherenkov Detector  |



# 2936 Bibliography

- 2937 [1] Liang Zhan, Yifang Wang, Jun Cao, and Liangjian Wen. "Determination of the Neutrino Mass  
2938 Hierarchy at an Intermediate Baseline". *Physical Review D* 78.11 (Dec. 10, 2008), 111103. ISSN:  
2939 1550-7998, 1550-2368. DOI: [10.1103/PhysRevD.78.111103](https://doi.org/10.1103/PhysRevD.78.111103). eprint: [0807.3203\[hep-ex, physics:hep-ph\]](https://arxiv.org/abs/0807.3203). URL: [http://arxiv.org/abs/0807.3203](https://arxiv.org/abs/0807.3203) (visited on 09/18/2023).
- 2940 [2] Fengpeng An et al. "Neutrino Physics with JUNO". *Journal of Physics G: Nuclear and Particle  
2941 Physics* 43.3 (Mar. 1, 2016), 030401. ISSN: 0954-3899, 1361-6471. DOI: [10.1088/0954-3899/43/3/030401](https://doi.org/10.1088/0954-3899/43/3/030401). eprint: [1507.05613\[hep-ex, physics:physics\]](https://arxiv.org/abs/1507.05613). URL: [http://arxiv.org/abs/1507.05613](https://arxiv.org/abs/1507.05613) (visited on 07/28/2023).
- 2942 [3] Liang Zhan, Yifang Wang, Jun Cao, and Liangjian Wen. "Experimental Requirements to Deter-  
2943 mine the Neutrino Mass Hierarchy Using Reactor Neutrinos". *Physical Review D* 79.7 (Apr. 14,  
2944 2009), 073007. ISSN: 1550-7998, 1550-2368. DOI: [10.1103/PhysRevD.79.073007](https://doi.org/10.1103/PhysRevD.79.073007). eprint: [0901.2976\[hep-ex\]](https://arxiv.org/abs/0901.2976). URL: [http://arxiv.org/abs/0901.2976](https://arxiv.org/abs/0901.2976) (visited on 09/18/2023).
- 2945 [4] A. A. Hahn, K. Schreckenbach, W. Gelletly, F. von Feilitzsch, G. Colvin, and B. Krusche. "Antineutrino spectra from 241Pu and 239Pu thermal neutron fission products". *Physics Letters B*  
2946 218.3 (Feb. 23, 1989), 365–368. ISSN: 0370-2693. DOI: [10.1016/0370-2693\(89\)91598-0](https://doi.org/10.1016/0370-2693(89)91598-0). URL:  
2947 <https://www.sciencedirect.com/science/article/pii/0370269389915980> (visited on  
2948 01/16/2024).
- 2949 [5] Th A. Mueller et al. "Improved Predictions of Reactor Antineutrino Spectra". *Physical Review C*  
2950 83.5 (May 23, 2011), 054615. ISSN: 0556-2813, 1089-490X. DOI: [10.1103/PhysRevC.83.054615](https://doi.org/10.1103/PhysRevC.83.054615).  
2951 eprint: [1101.2663\[hep-ex, physics:nucl-ex\]](https://arxiv.org/abs/1101.2663). URL: [http://arxiv.org/abs/1101.2663](https://arxiv.org/abs/1101.2663)  
2952 (visited on 01/16/2024).
- 2953 [6] F. von Feilitzsch, A. A. Hahn, and K. Schreckenbach. "Experimental beta-spectra from 239Pu  
2954 and 235U thermal neutron fission products and their correlated antineutrino spectra". *Physics  
2955 Letters B* 118.1 (Dec. 2, 1982), 162–166. ISSN: 0370-2693. DOI: [10.1016/0370-2693\(82\)90622-0](https://doi.org/10.1016/0370-2693(82)90622-0).  
2956 URL: <https://www.sciencedirect.com/science/article/pii/0370269382906220> (visited  
2957 on 01/16/2024).
- 2958 [7] K. Schreckenbach, G. Colvin, W. Gelletly, and F. Von Feilitzsch. "Determination of the antineu-  
2959 trino spectrum from 235U thermal neutron fission products up to 9.5 MeV". *Physics Letters B*  
2960 160.4 (Oct. 10, 1985), 325–330. ISSN: 0370-2693. DOI: [10.1016/0370-2693\(85\)91337-1](https://doi.org/10.1016/0370-2693(85)91337-1). URL:  
2961 <https://www.sciencedirect.com/science/article/pii/0370269385913371> (visited on  
2962 01/16/2024).
- 2963 [8] Patrick Huber. "On the determination of anti-neutrino spectra from nuclear reactors". *Physical  
2964 Review C* 84.2 (Aug. 29, 2011), 024617. ISSN: 0556-2813, 1089-490X. DOI: [10.1103/PhysRevC.84.024617](https://doi.org/10.1103/PhysRevC.84.024617). eprint: [1106.0687\[hep-ex, physics:hep-ph, physics:nucl-ex, physics:nucl-th\]](https://arxiv.org/abs/1106.0687).  
2965 URL: [http://arxiv.org/abs/1106.0687](https://arxiv.org/abs/1106.0687) (visited on 01/16/2024).
- 2966 [9] P. Vogel, G. K. Schenter, F. M. Mann, and R. E. Schenter. "Reactor antineutrino spectra and  
2967 their application to antineutrino-induced reactions. II". *Physical Review C* 24.4 (Oct. 1, 1981).  
2968 Publisher: American Physical Society, 1543–1553. DOI: [10.1103/PhysRevC.24.1543](https://doi.org/10.1103/PhysRevC.24.1543). URL:  
2969 <https://link.aps.org/doi/10.1103/PhysRevC.24.1543> (visited on 01/16/2024).
- 2970 [10] D. A. Dwyer and T. J. Langford. "Spectral Structure of Electron Antineutrinos from Nuclear  
2971 Reactors". *Physical Review Letters* 114.1 (Jan. 7, 2015), 012502. ISSN: 0031-9007, 1079-7114. DOI:  
2972 [10.1103/PhysRevLett.114.012502](https://doi.org/10.1103/PhysRevLett.114.012502). eprint: [1407.1281\[hep-ex, physics:nucl-ex\]](https://arxiv.org/abs/1407.1281). URL:  
2973 [http://arxiv.org/abs/1407.1281](https://arxiv.org/abs/1407.1281) (visited on 01/16/2024).

- [11] JUNO Collaboration et al. "Sub-percent Precision Measurement of Neutrino Oscillation Parameters with JUNO". *Chinese Physics C* 46.12 (Dec. 1, 2022), 123001. ISSN: 1674-1137, 2058-6132. DOI: [10.1088/1674-1137/ac8bc9](https://doi.org/10.1088/1674-1137/ac8bc9). eprint: [2204.13249 \[hep-ex\]](https://arxiv.org/abs/2204.13249). URL: <http://arxiv.org/abs/2204.13249> (visited on 08/11/2023).
- [12] JUNO Collaboration et al. *TAO Conceptual Design Report: A Precision Measurement of the Reactor Antineutrino Spectrum with Sub-percent Energy Resolution*. May 18, 2020. DOI: [10.48550/arXiv.2005.08745](https://doi.org/10.48550/arXiv.2005.08745). eprint: [2005.08745 \[hep-ex, physics:nucl-ex, physics:physics\]](https://arxiv.org/abs/2005.08745). URL: <http://arxiv.org/abs/2005.08745> (visited on 01/18/2024).
- [13] G. Mention, M. Fechner, Th. Lasserre, Th. A. Mueller, D. Lhuillier, M. Cribier, and A. Letourneau. "Reactor antineutrino anomaly". *Physical Review D* 83.7 (Apr. 29, 2011). Publisher: American Physical Society, 073006. DOI: [10.1103/PhysRevD.83.073006](https://doi.org/10.1103/PhysRevD.83.073006). URL: <https://link.aps.org/doi/10.1103/PhysRevD.83.073006> (visited on 03/05/2024).
- [14] V. Kopeikin, M. Skorokhvatov, and O. Titov. "Reevaluating reactor antineutrino spectra with new measurements of the ratio between  $^{235}\text{U}$  and  $^{239}\text{Pu}$   $\beta^-$  spectra". *Physical Review D* 104.7 (Oct. 25, 2021), L071301. ISSN: 2470-0010, 2470-0029. DOI: [10.1103/PhysRevD.104.L071301](https://doi.org/10.1103/PhysRevD.104.L071301). eprint: [2103.01684 \[hep-ph, physics:nucl-ex, physics:nucl-th\]](https://arxiv.org/abs/2103.01684). URL: <http://arxiv.org/abs/2103.01684> (visited on 01/18/2024).
- [15] A. Letourneau et al. "On the origin of the reactor antineutrino anomalies in light of a new summation model with parameterized  $\beta^-$  transitions". *Physical Review Letters* 130.2 (Jan. 10, 2023), 021801. ISSN: 0031-9007, 1079-7114. DOI: [10.1103/PhysRevLett.130.021801](https://doi.org/10.1103/PhysRevLett.130.021801). eprint: [2205.14954 \[hep-ex, physics:hep-ph\]](https://arxiv.org/abs/2205.14954). URL: <http://arxiv.org/abs/2205.14954> (visited on 01/16/2024).
- [16] Particle Data Group et al. "Review of Particle Physics". *Progress of Theoretical and Experimental Physics* 2020.8 (Aug. 14, 2020), 083C01. ISSN: 2050-3911. DOI: [10.1093/ptep/ptaa104](https://doi.org/10.1093/ptep/ptaa104). URL: <https://doi.org/10.1093/ptep/ptaa104> (visited on 12/04/2023).
- [17] Super-Kamiokande Collaboration et al. "Diffuse Supernova Neutrino Background Search at Super-Kamiokande". *Physical Review D* 104.12 (Dec. 10, 2021), 122002. ISSN: 2470-0010, 2470-0029. DOI: [10.1103/PhysRevD.104.122002](https://doi.org/10.1103/PhysRevD.104.122002). eprint: [2109.11174 \[astro-ph, physics:hep-ex\]](https://arxiv.org/abs/2109.11174). URL: <http://arxiv.org/abs/2109.11174> (visited on 02/28/2024).
- [18] JUNO Collaboration et al. "JUNO Sensitivity on Proton Decay  $p \rightarrow \bar{\nu}K^+$  Searches". *Chinese Physics C* 47.11 (Nov. 1, 2023), 113002. ISSN: 1674-1137, 2058-6132. DOI: [10.1088/1674-1137/ace9c6](https://doi.org/10.1088/1674-1137/ace9c6). eprint: [2212.08502 \[hep-ex, physics:hep-ph\]](https://arxiv.org/abs/2212.08502). URL: <http://arxiv.org/abs/2212.08502> (visited on 08/09/2024).
- [19] Alessandro Strumia and Francesco Vissani. "Precise quasielastic neutrino/nucleon cross section". *Physics Letters B* 564.1 (July 2003), 42–54. ISSN: 03702693. DOI: [10.1016/S0370-2693\(03\)00616-6](https://doi.org/10.1016/S0370-2693(03)00616-6). eprint: [astro-ph/0302055](https://arxiv.org/abs/astro-ph/0302055). URL: <http://arxiv.org/abs/astro-ph/0302055> (visited on 01/16/2024).
- [20] Daya Bay et al. *Optimization of the JUNO liquid scintillator composition using a Daya Bay antineutrino detector*. July 1, 2020. DOI: [10.48550/arXiv.2007.00314](https://doi.org/10.48550/arXiv.2007.00314). eprint: [2007.00314 \[hep-ex, physics:physics\]](https://arxiv.org/abs/2007.00314). URL: <http://arxiv.org/abs/2007.00314> (visited on 07/26/2023).
- [21] J. B. Birks. "CHAPTER 3 - THE SCINTILLATION PROCESS IN ORGANIC MATERIALS—I". *The Theory and Practice of Scintillation Counting*. Ed. by J. B. Birks. International Series of Monographs in Electronics and Instrumentation. Jan. 1, 1964, 39–67. ISBN: 978-0-08-010472-0. DOI: [10.1016/B978-0-08-010472-0.50008-2](https://doi.org/10.1016/B978-0-08-010472-0.50008-2). URL: <https://www.sciencedirect.com/science/article/pii/B9780080104720500082> (visited on 02/07/2024).
- [22] Photomultiplier tube R12860 | Hamamatsu Photonics. URL: [https://www.hamamatsu.com/eu/en/product/optical-sensors/pmt/pmt\\_tube-alone/head-on-type/R12860.html](https://www.hamamatsu.com/eu/en/product/optical-sensors/pmt/pmt_tube-alone/head-on-type/R12860.html) (visited on 02/08/2024).
- [23] Yan Zhang, Ze-Yuan Yu, Xin-Ying Li, Zi-Yan Deng, and Liang-Jian Wen. "A complete optical model for liquid-scintillator detectors". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 967 (July 2020), 163860. ISSN: 01689002. DOI: [10.1016/j.nima.2020.163860](https://doi.org/10.1016/j.nima.2020.163860). eprint: [2003.12212 \[physics\]](https://arxiv.org/abs/2003.12212). URL: <http://arxiv.org/abs/2003.12212> (visited on 02/07/2024).

- [24] Hai-Bo Yang et al. "Light Attenuation Length of High Quality Linear Alkyl Benzene as Liquid Scintillator Solvent for the JUNO Experiment". *Journal of Instrumentation* 12.11 (Nov. 27, 2017), T11004–T11004. ISSN: 1748-0221. DOI: [10.1088/1748-0221/12/11/T11004](https://doi.org/10.1088/1748-0221/12/11/T11004). eprint: [1703.01867](https://arxiv.org/abs/1703.01867) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/1703.01867> (visited on 07/28/2023).
- [25] JUNO Collaboration et al. *The Design and Sensitivity of JUNO's scintillator radiopurity pre-detector OSIRIS*. Mar. 31, 2021. DOI: [10.48550/arXiv.2103.16900](https://doi.org/10.48550/arXiv.2103.16900). eprint: [2103.16900](https://arxiv.org/abs/2103.16900) [physics]. URL: <http://arxiv.org/abs/2103.16900> (visited on 02/07/2024).
- [26] Angel Abusleme et al. "Mass Testing and Characterization of 20-inch PMTs for JUNO". *The European Physical Journal C* 82.12 (Dec. 24, 2022), 1168. ISSN: 1434-6052. DOI: [10.1140/epjc/s10052-022-11002-8](https://doi.org/10.1140/epjc/s10052-022-11002-8). eprint: [2205.08629](https://arxiv.org/abs/2205.08629) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/2205.08629> (visited on 02/08/2024).
- [27] Yang Han. "Dual Calorimetry for High Precision Neutrino Oscillation Measurement at JUNO Experiment". AstroParticule et Cosmologie, France, Paris U. VII, APC, June 2021.
- [28] R. Acquaferredda et al. "The OPERA experiment in the CERN to Gran Sasso neutrino beam". *Journal of Instrumentation* 4.4 (Apr. 2009), P04018. ISSN: 1748-0221. DOI: [10.1088/1748-0221/4/04/P04018](https://doi.org/10.1088/1748-0221/4/04/P04018). URL: <https://dx.doi.org/10.1088/1748-0221/4/04/P04018> (visited on 02/29/2024).
- [29] JUNO collaboration et al. "Calibration Strategy of the JUNO Experiment". *Journal of High Energy Physics* 2021.3 (Mar. 2021), 4. ISSN: 1029-8479. DOI: [10.1007/JHEP03\(2021\)004](https://doi.org/10.1007/JHEP03(2021)004). eprint: [2011.06405](https://arxiv.org/abs/2011.06405) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/2011.06405> (visited on 08/10/2023).
- [30] Hans Th J. Steiger. *TAO – The Taishan Antineutrino Observatory*. Sept. 21, 2022. DOI: [10.48550/arXiv.2209.10387](https://doi.org/10.48550/arXiv.2209.10387). eprint: [2209.10387](https://arxiv.org/abs/2209.10387) [physics]. URL: <http://arxiv.org/abs/2209.10387> (visited on 01/16/2024).
- [31] Tao Lin et al. "The Application of SNiPER to the JUNO Simulation". *Journal of Physics: Conference Series* 898.4 (Oct. 2017). Publisher: IOP Publishing, 042029. ISSN: 1742-6596. DOI: [10.1088/1742-6596/898/4/042029](https://doi.org/10.1088/1742-6596/898/4/042029). URL: <https://dx.doi.org/10.1088/1742-6596/898/4/042029> (visited on 02/27/2024).
- [32] S. Agostinelli et al. "Geant4—a simulation toolkit". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 506.3 (July 1, 2003), 250–303. ISSN: 0168-9002. DOI: [10.1016/S0168-9002\(03\)01368-8](https://doi.org/10.1016/S0168-9002(03)01368-8). URL: <https://www.sciencedirect.com/science/article/pii/S0168900203013688> (visited on 02/27/2024).
- [33] J. Allison et al. "Geant4 developments and applications". *IEEE Transactions on Nuclear Science* 53.1 (Feb. 2006). Conference Name: IEEE Transactions on Nuclear Science, 270–278. ISSN: 1558-1578. DOI: [10.1109/TNS.2006.869826](https://doi.org/10.1109/TNS.2006.869826). URL: <https://ieeexplore.ieee.org/document/1610988?isnumber=33833&arnumber=1610988&count=33&index=7> (visited on 02/27/2024).
- [34] J. Allison et al. "Recent developments in Geant4". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 835 (Nov. 1, 2016), 186–225. ISSN: 0168-9002. DOI: [10.1016/j.nima.2016.06.125](https://doi.org/10.1016/j.nima.2016.06.125). URL: <https://www.sciencedirect.com/science/article/pii/S0168900216306957> (visited on 02/27/2024).
- [35] Wenjie Wu, Miao He, Xiang Zhou, and Haoxue Qiao. "A new method of energy reconstruction for large spherical liquid scintillator detectors". *Journal of Instrumentation* 14.3 (Mar. 8, 2019), P03009–P03009. ISSN: 1748-0221. DOI: [10.1088/1748-0221/14/03/P03009](https://doi.org/10.1088/1748-0221/14/03/P03009). eprint: [1812.01799](https://arxiv.org/abs/1812.01799) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/1812.01799> (visited on 07/28/2023).
- [36] Guihong Huang et al. "Improving the energy uniformity for large liquid scintillator detectors". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1001 (June 11, 2021), 165287. ISSN: 0168-9002. DOI: [10.1016/j.nima.2021.165287](https://doi.org/10.1016/j.nima.2021.165287). URL: <https://www.sciencedirect.com/science/article/pii/S0168900221002710> (visited on 03/01/2024).
- [37] Ziyuan Li et al. "Event vertex and time reconstruction in large volume liquid scintillator detector". *Nuclear Science and Techniques* 32.5 (May 2021), 49. ISSN: 1001-8042, 2210-3147. DOI:

- 3086        [10.1007/s41365-021-00885-z](https://doi.org/10.1007/s41365-021-00885-z). eprint: [2101.08901 \[hep-ex, physics:physics\]](https://arxiv.org/abs/2101.08901). URL: <http://arxiv.org/abs/2101.08901> (visited on 07/28/2023).
- 3087
- 3088 [38] Gioacchino Ranucci. "An analytical approach to the evaluation of the pulse shape discrimination properties of scintillators". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 354.2 (Jan. 30, 1995), 389–399. ISSN: 0168-9002. DOI: [10.1016/0168-9002\(94\)00886-8](https://doi.org/10.1016/0168-9002(94)00886-8). URL: <https://www.sciencedirect.com/science/article/pii/0168900294008868> (visited on 03/07/2024).
- 3089
- 3090
- 3091
- 3092
- 3093 [39] C. Galbiati and K. McCarty. "Time and space reconstruction in optical, non-imaging, scintillator-based particle detectors". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 568.2 (Dec. 1, 2006), 700–709. ISSN: 0168-9002. DOI: [10.1016/j.nima.2006.07.058](https://doi.org/10.1016/j.nima.2006.07.058). URL: <https://www.sciencedirect.com/science/article/pii/S0168900206013519> (visited on 03/07/2024).
- 3094
- 3095
- 3096
- 3097
- 3098 [40] M. Moszyński and B. Bengtson. "Status of timing with plastic scintillation detectors". *Nuclear Instruments and Methods* 158 (Jan. 1, 1979), 1–31. ISSN: 0029-554X. DOI: [10.1016/S0029-554X\(79\)90170-8](https://doi.org/10.1016/S0029-554X(79)90170-8). URL: <https://www.sciencedirect.com/science/article/pii/S0029554X79901708> (visited on 03/07/2024).
- 3099
- 3100
- 3101
- 3102 [41] Gui-Hong Huang, Wei Jiang, Liang-Jian Wen, Yi-Fang Wang, and Wu-Ming Luo. "Data-driven simultaneous vertex and energy reconstruction for large liquid scintillator detectors". *Nuclear Science and Techniques* 34.6 (June 17, 2023), 83. ISSN: 2210-3147. DOI: [10.1007/s41365-023-01240-0](https://doi.org/10.1007/s41365-023-01240-0). URL: <https://doi.org/10.1007/s41365-023-01240-0> (visited on 08/17/2023).
- 3103
- 3104
- 3105
- 3106 [42] Zhen Qian et al. "Vertex and Energy Reconstruction in JUNO with Machine Learning Methods". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1010 (Sept. 2021), 165527. ISSN: 01689002. DOI: [10.1016/j.nima.2021.165527](https://doi.org/10.1016/j.nima.2021.165527). eprint: [2101.04839 \[hep-ex, physics:physics\]](https://arxiv.org/abs/2101.04839). URL: <http://arxiv.org/abs/2101.04839> (visited on 07/24/2023).
- 3107
- 3108
- 3109
- 3110
- 3111 [43] Arsenii Gavrikov, Yury Malyshkin, and Fedor Ratnikov. "Energy reconstruction for large liquid scintillator detectors with machine learning techniques: aggregated features approach". *The European Physical Journal C* 82.11 (Nov. 14, 2022), 1021. ISSN: 1434-6052. DOI: [10.1140/epjc/s10052-022-11004-6](https://doi.org/10.1140/epjc/s10052-022-11004-6). eprint: [2206.09040 \[physics\]](https://arxiv.org/abs/2206.09040). URL: <http://arxiv.org/abs/2206.09040> (visited on 07/24/2023).
- 3112
- 3113
- 3114
- 3115
- 3116 [44] R. Abbasi et al. "Graph Neural Networks for low-energy event classification & reconstruction in IceCube". *Journal of Instrumentation* 17.11 (Nov. 2022). Publisher: IOP Publishing, P11003. ISSN: 1748-0221. DOI: [10.1088/1748-0221/17/11/P11003](https://doi.org/10.1088/1748-0221/17/11/P11003). URL: <https://dx.doi.org/10.1088/1748-0221/17/11/P11003> (visited on 04/04/2024).
- 3117
- 3118
- 3119
- 3120 [45] S. Reck, D. Guderian, G. Vermarien, A. Domi, and on behalf of the KM3NeT collaboration on behalf of the. "Graph neural networks for reconstruction and classification in KM3NeT". *Journal of Instrumentation* 16.10 (Oct. 2021). Publisher: IOP Publishing, C10011. ISSN: 1748-0221. DOI: [10.1088/1748-0221/16/10/C10011](https://doi.org/10.1088/1748-0221/16/10/C10011). URL: <https://dx.doi.org/10.1088/1748-0221/16/10/C10011> (visited on 04/04/2024).
- 3121
- 3122
- 3123
- 3124
- 3125 [46] The IceCube collaboration et al. "A convolutional neural network based cascade reconstruction for the IceCube Neutrino Observatory". *Journal of Instrumentation* 16.7 (July 2021). Publisher: IOP Publishing, P07041. ISSN: 1748-0221. DOI: [10.1088/1748-0221/16/07/P07041](https://doi.org/10.1088/1748-0221/16/07/P07041). URL: <https://dx.doi.org/10.1088/1748-0221/16/07/P07041> (visited on 04/04/2024).
- 3126
- 3127
- 3128
- 3129 [47] DUNE Collaboration et al. "Neutrino interaction classification with a convolutional neural network in the DUNE far detector". *Physical Review D* 102.9 (Nov. 9, 2020). Publisher: American Physical Society, 092003. DOI: [10.1103/PhysRevD.102.092003](https://doi.org/10.1103/PhysRevD.102.092003). URL: <https://link.aps.org/doi/10.1103/PhysRevD.102.092003> (visited on 04/04/2024).
- 3130
- 3131
- 3132
- 3133 [48] K. M. Górski, E. Hivon, A. J. Banday, B. D. Wandelt, F. K. Hansen, M. Reinecke, and M. Bartelmann. "HEALPix: A Framework for High-Resolution Discretization and Fast Analysis of Data Distributed on the Sphere". *The Astrophysical Journal* 622 (Apr. 1, 2005). ADS Bibcode: 2005ApJ...622..759G, 759–771. ISSN: 0004-637X. DOI: [10.1086/427976](https://doi.org/10.1086/427976). URL: <https://ui.adsabs.harvard.edu/abs/2005ApJ...622..759G> (visited on 04/04/2024).
- 3134
- 3135
- 3136
- 3137

- 3138 [49] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. *Convolutional Neural Networks*  
3139 *on Graphs with Fast Localized Spectral Filtering*. Feb. 5, 2017. DOI: [10.48550/arXiv.1606.09375](https://doi.org/10.48550/arXiv.1606.09375).  
3140 eprint: [1606.09375\[cs, stat\]](https://arxiv.org/abs/1606.09375). URL: <http://arxiv.org/abs/1606.09375> (visited on  
3141 04/04/2024).
- 3142 [50] JUNO Collaboration et al. "JUNO Physics and Detector". *Progress in Particle and Nuclear Physics*  
3143 123 (Mar. 2022), 103927. ISSN: 01466410. DOI: [10.1016/j.ppnp.2021.103927](https://doi.org/10.1016/j.ppnp.2021.103927). eprint: [2104.02565\[hep-ex\]](https://arxiv.org/abs/2104.02565). URL: <http://arxiv.org/abs/2104.02565> (visited on 09/18/2023).
- 3144 [51] Leo Breiman, Jerome Friedman, R. A. Olshen, and Charles J. Stone. *Classification and Regression*  
3145 *Trees*. New York: Chapman and Hall/CRC, Oct. 25, 2017. 368 pp. ISBN: 978-1-315-13947-0. DOI:  
3146 [10.1201/9781315139470](https://doi.org/10.1201/9781315139470).
- 3147 [52] Jerome H. Friedman. "Greedy function approximation: A gradient boosting machine." *The Annals of Statistics* 29.5 (Oct. 2001). Publisher: Institute of Mathematical Statistics, 1189–1232. ISSN:  
3148 0090-5364, 2168-8966. DOI: [10.1214/aos/1013203451](https://doi.org/10.1214/aos/1013203451). URL: <https://projecteuclid.org/journals/annals-of-statistics/volume-29/issue-5/Greedy-function-approximation-A-gradient-boosting-machine/10.1214/aos/1013203451.full> (visited on 04/29/2024).
- 3149 [53] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. Jan. 29, 2017.  
3150 DOI: [10.48550/arXiv.1412.6980](https://doi.org/10.48550/arXiv.1412.6980). eprint: [1412.6980\[cs\]](https://arxiv.org/abs/1412.6980). URL: <http://arxiv.org/abs/1412.6980> (visited on 05/13/2024).
- 3151 [54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep Residual Learning for Image  
3152 Recognition". *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016  
3153 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). ISSN: 1063-6919. June  
3154 2016, 770–778. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90). URL: <https://ieeexplore.ieee.org/document/7780459> (visited on 07/17/2024).
- 3155 [55] Olga Russakovsky et al. *ImageNet Large Scale Visual Recognition Challenge*. Jan. 29, 2015. DOI:  
3156 [10.48550/arXiv.1409.0575](https://doi.org/10.48550/arXiv.1409.0575). eprint: [1409.0575\[cs\]](https://arxiv.org/abs/1409.0575). URL: <http://arxiv.org/abs/1409.0575>  
3157 (visited on 05/17/2024).
- 3158 [56] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image*  
3159 *Recognition*. Apr. 10, 2015. DOI: [10.48550/arXiv.1409.1556](https://doi.org/10.48550/arXiv.1409.1556). eprint: [1409.1556\[cs\]](https://arxiv.org/abs/1409.1556). URL:  
3160 <http://arxiv.org/abs/1409.1556> (visited on 05/17/2024).
- 3161 [57] Anna Allen. *generic-github-user/Image-Convolution-Playground*. original-date: 2018-09-28T22:42:55Z.  
3162 July 15, 2024. URL: <https://github.com/generic-github-user/Image-Convolution-Playground> (visited on 07/16/2024).
- 3163 [58] Jason Ansel et al. *PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Trans-*  
3164 *formation and Graph Compilation*. Publication Title: 29th ACM International Conference on Ar-  
3165 *chitectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS*  
3166 '24) original-date: 2016-08-13T05:26:41Z. Apr. 2024. DOI: [10.1145/3620665.3640366](https://doi.org/10.1145/3620665.3640366). URL:  
3167 <https://pytorch.org/assets/pytorch2-2.pdf> (visited on 07/16/2024).
- 3168 [59] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document  
3169 recognition". *Proceedings of the IEEE* 86.11 (Nov. 1998). Conference Name: Proceedings of the  
3170 IEEE, 2278–2324. ISSN: 1558-2256. DOI: [10.1109/5.726791](https://doi.org/10.1109/5.726791). URL: <https://ieeexplore.ieee.org/document/726791> (visited on 07/16/2024).
- 3171 [60] NVIDIA T4 Tensor Core GPUs for Accelerating Inference. NVIDIA. URL: <https://www.nvidia.com/en-gb/data-center/tesla-t4/> (visited on 07/16/2024).
- 3172 [61] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. *Neural*  
3173 *Message Passing for Quantum Chemistry*. June 12, 2017. DOI: [10.48550/arXiv.1704.01212](https://doi.org/10.48550/arXiv.1704.01212).  
3174 eprint: [1704.01212\[cs\]](https://arxiv.org/abs/1704.01212). URL: <http://arxiv.org/abs/1704.01212> (visited on 05/22/2024).
- 3175 [62] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. *Diffusion Convolutional Recurrent Neural*  
3176 *Network: Data-Driven Traffic Forecasting*. Feb. 22, 2018. DOI: [10.48550/arXiv.1707.01926](https://doi.org/10.48550/arXiv.1707.01926).  
3177 eprint: [1707.01926\[cs, stat\]](https://arxiv.org/abs/1707.01926). URL: <http://arxiv.org/abs/1707.01926> (visited on  
3178 05/22/2024).
- 3179 [63] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
3180 Ozair, Aaron Courville, and Yoshua Bengio. *Generative Adversarial Networks*. June 10, 2014. DOI:  
3181 [10.48550/arXiv.1406.2891](https://doi.org/10.48550/arXiv.1406.2891).

- 3190        10.48550/arXiv.1406.2661. eprint: 1406.2661[cs,stat]. URL: <http://arxiv.org/abs/1406.2661> (visited on 05/29/2024).
- 3191
- 3192 [64] Victor Lebrin. "Towards the Detection of Core-Collapse Supernovae Burst Neutrinos with the  
3193 3-inch PMT System of the JUNO Detector". These de doctorat. Nantes Université, Sept. 5, 2022.  
3194 URL: <https://theses.fr/2022NANU4080> (visited on 05/22/2024).
- 3195 [65] Dan Cireşan, Ueli Meier, and Juergen Schmidhuber. *Multi-column Deep Neural Networks for  
3196 Image Classification*. version: 1. Feb. 13, 2012. DOI: 10.48550/arXiv.1202.2745. eprint: 1202.  
3197 2745[cs]. URL: <http://arxiv.org/abs/1202.2745> (visited on 06/27/2024).
- 3198 [66] R. Abbasi et al. "A Convolutional Neural Network based Cascade Reconstruction for the Ice-  
3199 Cube Neutrino Observatory". *Journal of Instrumentation* 16.7 (July 1, 2021), P07041. ISSN: 1748-  
3200 0221. DOI: 10.1088/1748-0221/16/07/P07041. eprint: 2101.11589[hep-ex]. URL: <http://arxiv.org/abs/2101.11589> (visited on 06/27/2024).
- 3201 [67] D. Maksimović, M. Nieslony, and M. Wurm. "CNNs for enhanced background discrimination  
3202 in DSNB searches in large-scale water-Gd detectors". *Journal of Cosmology and Astroparticle  
3203 Physics* 2021.11 (Nov. 2021). Publisher: IOP Publishing, 051. ISSN: 1475-7516. DOI: 10.1088/  
3204 1475-7516/2021/11/051. URL: <https://dx.doi.org/10.1088/1475-7516/2021/11/051>  
3205 (visited on 06/27/2024).
- 3206 [68] Taco S. Cohen, Mario Geiger, Jonas Koehler, and Max Welling. *Spherical CNNs*. Feb. 25, 2018.  
3207 DOI: 10.48550/arXiv.1801.10130. eprint: 1801.10130[cs,stat]. URL: <http://arxiv.org/abs/1801.10130> (visited on 07/13/2024).
- 3208 [69] NVIDIA A100 GPUs Power the Modern Data Center. NVIDIA. URL: <https://www.nvidia.com/en-gb/data-center/a100/> (visited on 08/06/2024).
- 3209 [70] NVIDIA V100. NVIDIA. URL: <https://www.nvidia.com/en-gb/data-center/v100/> (visited on 08/06/2024).
- 3210 [71] Leonard Imbert. *leonard-IMBERT/datamo*. original-date: 2023-10-17T12:37:38Z. Aug. 9, 2024. URL:  
3211 <https://github.com/leonard-IMBERT/datamo> (visited on 08/09/2024).
- 3212 [72] "IEEE Standard for Floating-Point Arithmetic". *IEEE Std 754-2019 (Revision of IEEE 754-2008)*  
3213 (July 2019). Conference Name: IEEE Std 754-2019 (Revision of IEEE 754-2008), 1–84. DOI: 10.  
3214 1109/IEEESTD.2019.8766229. URL: <https://ieeexplore.ieee.org/document/8766229>  
3215 (visited on 07/03/2024).
- 3216 [73] Chuanya Cao et al. "Mass production and characterization of 3-inch PMTs for the JUNO experiment".  
3217 *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers,  
3218 Detectors and Associated Equipment* 1005 (July 2021), 165347. ISSN: 01689002. DOI: 10.1016/j.nima.  
3219 2021.165347. eprint: 2102.11538[hep-ex,physics:physics]. URL: <http://arxiv.org/abs/2102.11538> (visited on 02/08/2024).
- 3220 [74] K. M. Gorski, E. Hivon, A. J. Banday, B. D. Wandelt, F. K. Hansen, M. Reinecke, and M. Bartelman.  
3221 "HEALPix – a Framework for High Resolution Discretization, and Fast Analysis of Data  
3222 Distributed on the Sphere". *The Astrophysical Journal* 622.2 (Apr. 2005), 759–771. ISSN: 0004-  
3223 637X, 1538-4357. DOI: 10.1086/427976. eprint: astro-ph/0409513. URL: <http://arxiv.org/abs/astro-ph/0409513> (visited on 08/10/2023).
- 3224 [75] Teng Li, Xin Xia, Xing-Tao Huang, Jia-Heng Zou, Wei-Dong Li, Tao Lin, Kun Zhang, and Zi-Yan  
3225 Deng. "Design and development of JUNO event data model\*\*". *Chinese Physics C* 41.6 (June  
3226 2017). Publisher: IOP Publishing, 066201. ISSN: 1674-1137. DOI: 10.1088/1674-1137/41/  
3227 6/066201. URL: <https://dx.doi.org/10.1088/1674-1137/41/6/066201> (visited on  
3228 08/16/2024).
- 3229 [76] Martin Reinecke. *Ducc0*. original-date: 2021-04-12T15:35:50Z. Aug. 9, 2024. URL: <https://gitlab.mpcdf.mpg.de/mtr/ducc> (visited on 08/16/2024).
- 3230 [77] Mario Schwarz, Sabrina M. Franke, Lothar Oberauer, Miriam D. Plein, Hans Th J. Steiger,  
3231 and Marc Tippmann. *Measurements of the Lifetime of Orthopositronium in the LAB-Based Liquid  
3232 Scintillator of JUNO*. Apr. 25, 2018. DOI: 10.1016/j.nima.2018.12.068. eprint: 1804.  
3233 09456[physics]. URL: <http://arxiv.org/abs/1804.09456> (visited on 09/17/2024).
- 3234 [78] Narongkiat Rodphai, Zhimin Wang, Narumon Suwonjandee, and Burin Asavapibhop. "20-  
3235 inch photomultiplier tube timing study for JUNO". *Journal of Physics: Conference Series* 2145.1  
3236 (visited on 09/17/2024).

- 3243 (Dec. 2021). Publisher: IOP Publishing, 012017. ISSN: 1742-6596. DOI: [10.1088/1742-6596/2145/1/012017](https://doi.org/10.1088/1742-6596/2145/1/012017). URL: <https://dx.doi.org/10.1088/1742-6596/2145/1/012017> (visited on 09/17/2024).
- 3244 [79] Dong-Hao Liao et al. "Study of TTS for a 20-inch dynode PMT\*". *Chinese Physics C* 41.7 (July 3245 2017). Publisher: IOP Publishing, 076001. ISSN: 1674-1137. DOI: [10.1088/1674-1137/41/7/076001](https://doi.org/10.1088/1674-1137/41/7/076001). URL: <https://dx.doi.org/10.1088/1674-1137/41/7/076001> (visited on 09/17/2024).
- 3246 [80] Nan Li et al. "Characterization of 3-inch photomultiplier tubes for the JUNO central detector". 3247 *Radiation Detection Technology and Methods* 3.1 (Nov. 22, 2018), 6. ISSN: 2509-9949. DOI: [10.1007/s41605-018-0085-8](https://doi.org/10.1007/s41605-018-0085-8). URL: <https://doi.org/10.1007/s41605-018-0085-8> (visited on 3248 09/17/2024).
- 3249 [81] Anatael Cabrera et al. *Multi-Calorimetry in Light-based Neutrino Detectors*. Dec. 20, 2023. DOI: 3250 [10.48550/arXiv.2312.12991](https://arxiv.org/abs/2312.12991). eprint: [2312.12991\[hep-ex, physics:physics\]](https://arxiv.org/abs/2312.12991). URL: [http://arxiv.org/abs/2312.12991](https://arxiv.org/abs/2312.12991) (visited on 08/19/2024).
- 3251 [82] Angel Abusleme et al. "Potential to Identify the Neutrino Mass Ordering with Reactor 3252 Antineutrinos in JUNO" (May 2024). eprint: [2405.18008](https://arxiv.org/abs/2405.18008).
- 3253 [83] Rene Brun et al. *root-project/root: v6.26/06*. Version v6-26-06. Mar. 3, 2022. DOI: [10.5281/zenodo.3895860](https://doi.org/10.5281/zenodo.3895860). URL: <https://zenodo.org/records/3895860> (visited on 09/05/2024).
- 3254 [84] X. B. Ma, W. L. Zhong, L. Z. Wang, Y. X. Chen, and J. Cao. "Improved calculation of the energy 3255 release in neutron-induced fission". *Physical Review C* 88.1 (July 12, 2013). Publisher: American 3256 Physical Society, 014605. DOI: [10.1103/PhysRevC.88.014605](https://doi.org/10.1103/PhysRevC.88.014605). URL: <https://link.aps.org/doi/10.1103/PhysRevC.88.014605> (visited on 09/06/2024).
- 3257 [85] Daya Bay Collaboration et al. "Measurement of the Reactor Antineutrino Flux and Spectrum at 3258 Daya Bay". *Physical Review Letters* 116.6 (Feb. 12, 2016). Publisher: American Physical Society, 061801. DOI: [10.1103/PhysRevLett.116.061801](https://doi.org/10.1103/PhysRevLett.116.061801). URL: <https://link.aps.org/doi/10.1103/PhysRevLett.116.061801> (visited on 09/06/2024).
- 3259 [86] Timo Gnambs. "A Brief Note on the Standard Error of the Pearson Correlation". *Collabra: 3260 Psychology* 9.1 (Sept. 6, 2023). Ed. by Thomas Evans, 87615. ISSN: 2474-7394. DOI: [10.1525/collabra.87615](https://doi.org/10.1525/collabra.87615). URL: <https://doi.org/10.1525/collabra.87615> (visited on 09/10/2024).
- 3261 [87] "Note Sur Une Méthode de Résolution des équations Normales Provenant de L'Application 3262 de la MéThode des Moindres Carrés a un Système D'équations Linéaires en Nombre Inférieur 3263 a Celui des Inconnues. — Application de la Méthode a la Résolution D'un Système Defini 3264 D'éQuations LinéAires". *Bulletin géodésique* 2.1 (Apr. 1, 1924), 67–77. ISSN: 1432-1394. DOI: [10.1007/BF03031308](https://doi.org/10.1007/BF03031308). URL: <https://doi.org/10.1007/BF03031308> (visited on 09/10/2024).
- 3265 [88] Pauli Virtanen et al. "SciPy 1.0: fundamental algorithms for scientific computing in Python". 3266 *Nature Methods* 17.3 (Mar. 2020). Publisher: Nature Publishing Group, 261–272. ISSN: 1548-7105. 3267 DOI: [10.1038/s41592-019-0686-2](https://doi.org/10.1038/s41592-019-0686-2). URL: <https://www.nature.com/articles/s41592-019-0686-2> (visited on 08/14/2024).

3281

3282

**Titre :** Méthode Deep Learning and analyse Double Calorimétrique pour la mesure de haute précision des paramètres d'oscillation des neutrinos dans JUNO

**Mot clés :** Neutrinos; expérience JUNO; Deep Learning; reconstruction d'IBD; oscillations des neutrinos; double calorimetrie

**Résumé :** JUNO est un observatoire de neutrinos à scintillateur liquide, polyvalent et medium baseline (environ 52 km), situé en Chine. Ses principaux objectifs sont de mesurer les paramètres d'oscillation  $\theta_{12}$ ,  $\Delta m_{21}^2$  et  $\Delta m_{31}^2$  avec une précision au pour-mille et de déterminer l'ordre des masses des neutrinos avec un niveau de confiance de  $3\sigma$ . Atteindre ces objectifs nécessite une résolution énergétique sans précédent de  $3\%/\sqrt{E(\text{MeV})}$  avec cette technologie. Cela demande une compréhension approfondie des divers effets au sein du détecteur.

Le système de double calorimetrie, composé de deux systèmes de mesure distincts observant le même événement, permet non seulement une calibration mais aussi une détection des effets du détecteur avec une grande précision, comme démontré dans cette thèse. Le Deep Learning, un outil de plus en plus utilisé en physique expérimentale, joue un rôle crucial dans cet effort. Dans cette thèse, je présente le développement, l'application et l'analyse des techniques de Deep Learning pour la reconstruction d'évènements dans l'expérience JUNO.

3314

**Title:** Deep learning methods and Dual Calorimetric analysis for high precision neutrino oscillation measurements at JUNO

**Keywords:** Neutrinos; JUNO experiment; Deep learning; IBD reconstruction; neutrinos Oscillation; dual Calorimetry

**Abstract:** JUNO is a multipurpose, medium baseline ( $\sim 52$  km) liquid scintillator neutrino observatory located in China. Its primary objectives are to measure the oscillation parameters  $\theta_{12}$ ,  $\Delta m_{21}^2$ , and  $\Delta m_{31}^2$  with per mil precision and to determine the neutrino mass ordering at a  $3\sigma$  confidence level. Achieving these goals requires an unprecedented energy resolution of  $3\%/\sqrt{E(\text{MeV})}$  with this technology. This demands a comprehensive understanding of the various effects within the

detect. The Dual Calorimetry system—two distinct measurement systems observing the same event-enables not only high-precision calibration but also detection of detector effects, as demonstrated in this thesis. Deep learning, an increasingly powerful tool in physics, plays a critical role in this effort. In this thesis, I present the development, application, and analysis of Deep Learning techniques for reconstruction in the JUNO experiment.

