

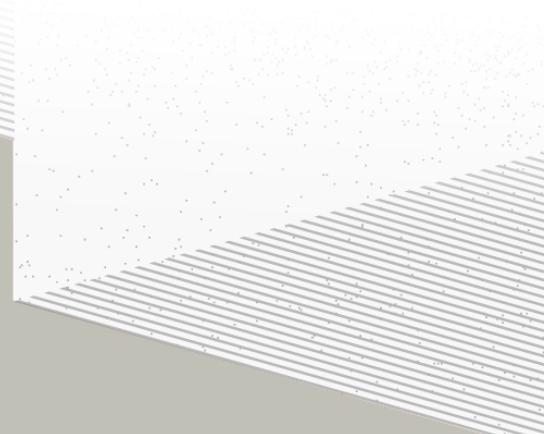
1

2

THÈSE DE DOCTORAT DE

NANTES UNIVERSITÉ

ÉCOLE DOCTORALE N° 596
Matière, Molécules, Matériaux
Spécialité : *Physique Subatomique et Instrumentation Nucléaire*



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

³ Contents

⁴	Contents	1
⁵	Remerciements	5
⁶	Introduction	7
⁷	1 Neutrino physics	9
⁸	1.1 Standard model	9
⁹	1.1.1 Limits of the standard model	9
¹⁰	1.2 Historic of the neutrino	9
¹¹	1.3 Oscillation	9
¹²	1.3.1 Phenomologies	9
¹³	1.4 Open questions	9
¹⁴	2 The JUNO experiment	11
¹⁵	2.1 Reactor Neutrinos physics in JUNO	12
¹⁶	2.1.1 Antineutrino spectrum measured in JUNO	12
¹⁷	2.1.2 Background spectra	15
¹⁸	2.2 Other physics	15
¹⁹	2.3 The JUNO detector	16
²⁰	2.3.1 Detection principle	17
²¹	2.3.2 Central Detector (CD)	18
²²	2.3.3 Veto detector	22
²³	2.4 Calibration strategy	23
²⁴	2.4.1 Energy scale calibration	23
²⁵	2.4.2 Calibration system	24
²⁶	2.4.3 Instrumental non-linearity calibration	25
²⁷	2.5 Satellite detectors	26
²⁸	2.5.1 TAO	26
²⁹	2.5.2 OSIRIS	27
³⁰	2.6 Software	27
³¹	2.7 Reactor anti-neutrino oscillation analysis	28
³²	2.7.1 IBD samples selection	28
³³	2.7.2 Synthetic overview of fit procedures developed at JUNO	29
³⁴	2.7.3 The spectrum model and sources of systematic uncertainties	31

35	2.7.4 Versions of the fit used in this thesis	33
36	2.8 State of the art of the Offline IBD reconstruction in JUNO	34
37	2.8.1 Interaction vertex reconstruction	34
38	2.8.2 Energy reconstruction	38
39	2.8.3 Machine learning for reconstruction	41
40	2.8.4 Physics results	44
41	2.9 Summary	44
42	3 Machine learning: Introduction to the methods and algorithms used in this thesis	45
43	3.1 Core concepts in machine learning and neural networks	46
44	3.1.1 Boosted Decision Tree (BDT)	46
45	3.1.2 Artificial Neural Network (NN)	46
46	3.1.3 Training procedure	48
47	3.1.4 Potential pitfalls	51
48	3.2 Neural networks architectures	54
49	3.2.1 Fully Connected Deep Neural Network (FCDNN)	54
50	3.2.2 Convolutional Neural Network (CNN)	54
51	3.2.3 Graph Neural Network (GNN)	56
52	3.2.4 Adversarial Neural Network (ANN)	58
53	4 Image recognition for IBD reconstruction with the SPMT system	59
54	4.1 Method and model	60
55	4.1.1 Model	61
56	4.1.2 Data representation	62
57	4.1.3 Dataset	64
58	4.1.4 Data characteristics	65
59	4.2 Training	67
60	4.3 Results	67
61	4.3.1 J21 results	68
62	4.3.2 J21 Combination of classic and ML estimator	70
63	4.3.3 J23 results	72
64	4.4 Conclusion and prospect	74
65	5 Graph representation of JUNO for IBD reconstruction	77
66	5.1 Data representation	78
67	5.2 Message passing algorithm	80
68	5.3 Data	82
69	5.4 Model	84
70	5.5 Training	84
71	5.6 Optimization	86
72	5.6.1 Software optimization	86
73	5.6.2 Hyperparameters optimization	87
74	5.7 performance of the final version	87
75	5.8 Conclusion	91

76	6 Reliability of machine learning methods	95
77	6.1 Method	96
78	6.2 Architecture	96
79	6.2.1 Back-propagation problematic	98
80	6.2.2 Reconstruction Network	99
81	6.2.3 Adversarial Neural Network	99
82	6.2.4 Training	99
83	6.3 Results	99
84	6.3.1 Back to identity	99
85	6.3.2 Breaking of the reconstruction	99
86	6.4 Conclusion and prospect	99
87	7 Dual calorimetric analysis for Precision Measurement	101
88	7.1 Motivations	103
89	7.1.1 Discrepancies between the SPMT and LPMT results	103
90	7.1.2 Charge Non-Linearity (QNL)	104
91	7.2 Approach	105
92	7.2.1 Data production	106
93	7.2.2 Individual fits	107
94	7.2.3 Joint fit	108
95	7.2.4 Data and theoretical spectrum generation	109
96	7.2.5 Limitations	110
97	7.3 Fit software	110
98	7.3.1 IBD generator	111
99	7.3.2 Fit	112
100	7.4 Technical challenges and development	113
101	7.5 Results	113
102	7.5.1 Validation	113
103	7.5.2 Covariance matrix	117
104	7.5.3 Statistical tests	122
105	7.6 Conclusion and perspectives	124
106	8 Conclusion	129
107	A Calculation of optimal α for estimator combination	131
108	A.1 Unbiased estimator	131
109	A.2 Optimal variance estimator	131
110	B Charge spherical harmonics analysis	133
111	C Additional spectrum smearing	141
112	D Correction of E_{vis} bias	143
113	List of Tables	145

114	List of Figures	154
115	List of Abbreviations	155
116	Bibliography	157

¹¹⁷ **Remerciements**

¹¹⁸ **Introduction**

¹¹⁹ **Chapter 1**

¹²⁰ **Neutrino physics**

¹²¹ *The neutrino, or ν 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.*

¹²² **Contents**

¹²³	1.1 Standard model	9
¹²⁴	1.1.1 Limits of the standard model	9
¹²⁵	1.2 Historic of the neutrino	9
¹²⁶	1.3 Oscillation	9
¹²⁷	1.3.1 Phenomologies	9
¹²⁸	1.4 Open questions	9
¹²⁹		
¹³⁰		
¹³¹		
¹³²		

¹³³ **1.1 Standard model**

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

¹³⁴ **1.1.1 Limits of the standard model**

¹³⁵ **1.2 Historic of the neutrino**

¹³⁶ **First theories**

¹³⁷ **Discovery**

¹³⁸ **Milestones and anomalies**

¹³⁹ **1.3 Oscillation**

¹⁴⁰ **1.3.1 Phenomologies**

¹⁴¹ **1.4 Open questions**

¹⁴² **Chapter 2**

¹⁴³ **The JUNO experiment**

¹⁴⁴

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

¹⁴⁵

Contents

¹⁴⁶ 2.1 Reactor Neutrinos physics in JUNO	¹⁴⁷	¹²
¹⁴⁸ 2.1.1 Antineutrino spectrum measured in JUNO	¹⁴⁹	¹²
¹⁴⁹ 2.1.2 Background spectra	¹⁵⁰	¹⁵
¹⁵⁰ 2.2 Other physics	¹⁵¹	¹⁵
¹⁵¹ 2.3 The JUNO detector	¹⁵²	¹⁶
¹⁵² 2.3.1 Detection principle	¹⁵³	¹⁷
¹⁵³ 2.3.2 Central Detector (CD)	¹⁵⁴	¹⁸
¹⁵⁴ 2.3.3 Veto detector	¹⁵⁵	²²
¹⁵⁵ 2.4 Calibration strategy	¹⁵⁶	²³
¹⁵⁶ 2.4.1 Energy scale calibration	¹⁵⁷	²³
¹⁵⁷ 2.4.2 Calibration system	¹⁵⁸	²⁴
¹⁵⁸ 2.4.3 Instrumental non-linearity calibration	¹⁵⁹	²⁵
¹⁵⁹ 2.5 Satellite detectors	¹⁶⁰	²⁶
¹⁶⁰ 2.5.1 TAO	¹⁶¹	²⁶
¹⁶¹ 2.5.2 OSIRIS	¹⁶²	²⁷
¹⁶² 2.6 Software	¹⁶³	²⁷
¹⁶³ 2.7 Reactor anti-neutrino oscillation analysis	¹⁶⁴	²⁸
¹⁶⁴ 2.7.1 IBD samples selection	¹⁶⁵	²⁸
¹⁶⁵ 2.7.2 Synthetic overview of fit procedures developed at JUNO	¹⁶⁶	²⁹
¹⁶⁶ 2.7.3 The spectrum model and sources of systematic uncertainties	¹⁶⁷	³¹
¹⁶⁷ 2.7.4 Versions of the fit used in this thesis	¹⁶⁸	³³
¹⁶⁸ 2.8 State of the art of the Offline IBD reconstruction in JUNO	¹⁶⁹	³⁴
¹⁶⁹ 2.8.1 Interaction vertex reconstruction	¹⁷⁰	³⁴
¹⁷⁰ 2.8.2 Energy reconstruction	¹⁷¹	³⁸
¹⁷¹ 2.8.3 Machine learning for reconstruction	¹⁷²	⁴¹
¹⁷² 2.8.4 Physics results	¹⁷³	⁴⁴
¹⁷³ 2.9 Summary	¹⁷⁴	⁴⁴

¹⁷⁴

¹⁷⁵ 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 Δm_{32}^2 . From this idea is born the Jiangmen Underground Neutrino Observatory (JUNO) experiment.

¹⁷⁶

JUNO is a neutrino detection experiment under construction located in China, in Guangdong proving, near the city of Kaiping. Its main objectives are the determination of the mass ordering at the $3\text{-}4\sigma$ level in 6 years of data taking and the measurement at the sub-percent precision of the oscillation parameters Δm_{21}^2 , $\sin^2 \theta_{12}$, Δm_{32}^2 and with less precision $\sin^2 \theta_{13}$ [2].



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

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

2.1 Reactor Neutrinos physics in JUNO

JUNO will try to determine the NMO and to bring at the few per mille level our knowledge of Δm_{31}^2 , Δm_{21}^2 and $\sin^2(2\theta_{12})$ via the precision analysis of the spectrum of the visible energy left by reactor antineutrinos in its detector.

2.1.1 Antineutrino spectrum measured in JUNO

To some extent, this analysis is equivalent to extracting from this spectrum the oscillation probability [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 Δm_{32}^2 and Δm_{31}^2 causing a phase shift of the spectrum as we can see in the figure 2.2.

In practice, a fit to the grey distribution of figure 2.3 will be performed. It is the sum of two components :signal (black) and backgrounds (colored). Reactor antineutrinos are detected by JUNO via Inverse Beta Decays (IBD) : $\bar{n}\bar{\nu}_e + p^- \rightarrow e^+ + n$. The energy spectrum under investigation is therefore that of the reconstructed e^+ visible energy. The black signal spectrum is therefore the sum of the antineutrino differential fluxes from all reactors and reaching the detecteur, weighted by the oscillation probability of Eq 2.1.1 and the IBD differential cross section and convoluted with detection

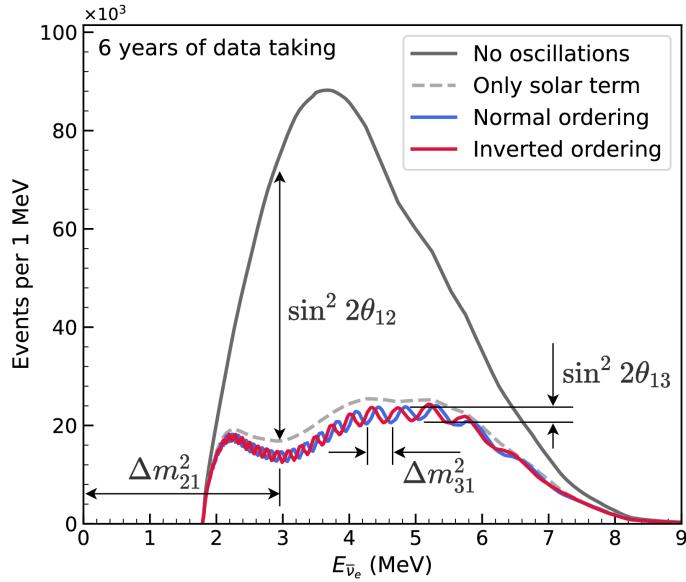


FIGURE 2.2 – Expected number of neutrinos event per MeV in JUNO after 6 years of data taking. The black curve shows the flux if there was no oscillation. The light gray curve shows the oscillation if only the solar terms are taken in account (θ_{12} , Δm_{21}^2). The blue and red curve shows the spectrum 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.

206 effects. These various ingredients are theoretically modelled in order to provide the probability
207 density function (PDF) to be used in the fit.

208 To reach JUNO's goals, it takes that this experimental spectrum still bears sizeable traces of the very
209 small phase shift mentioned above. Most notably, the following requirements must be fulfilled :

- 210 1. An energy resolution of $3\%/\sqrt{E(\text{MeV})}$ to be able to distinguish the fine structure of the fast
211 oscillation.
- 212 2. An energy scale known at the better than the 1% level.
- 213 3. A baseline between 40 and 65 km to maximise the $\bar{\nu}_e$ oscillation probability. The optimal
214 baseline would be 58 km and JUNO baseline is 53 km.
- 215 4. At least $\approx 100,000$ events. This is the necessary statistics to reach JUNO's canonical sensitivity
216 after 6 years of data taking.

217 $\bar{\nu}_e$ flux coming from nuclear power plants

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

222 The $\bar{\nu}_e$ coming from reactors are emitted from β -decay of unstable fission fragments. The Taishan
223 and Yangjiang reactors are Pressurised Water Reactor (PWR), the same type as Daya Bay. In those
224 type of reactor more the 99.7 % and $\bar{\nu}_e$ are produced by the fissions of four fuel isotopes ^{235}U , ^{238}U ,
225 ^{239}Pu and ^{241}Pu . The neutrino flux per fission of each isotope is determined by the inversion of the
226 measured β spectra of fission product [4–8] or by calculation using the nuclear databases [9, 10].

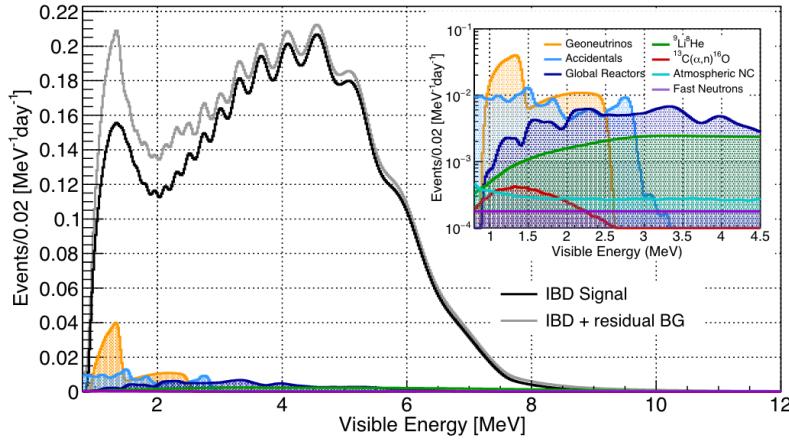


FIGURE 2.3 – Expected visible energy spectrum measured with the LPMT system with (grey) and without (black) backgrounds. The background amount for about 7% of the IBD candidate and are mostly localized below 3 MeV [3]

Reactor	Power (GW _{th})	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.

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

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

228 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
229 thermal energy released in each fission and $S_i(e_\nu)$ the neutrino flux per fission for this isotope.

230 The latter flux is difficult to predict. To evaluate JUNO's sensitivity and to serve as a starting point
231 in the spectrum PDF, the Huber-Mueller model is used [5], corrected using Daya Bays data [11] to
232 account for a $\sim 5\%$ deficit with respect to models, referred to as the reactor antineutrino anomaly [12],
233 and for a discrepancy between models and data in the spectral shape (the so call 5 MeV bump).

234 In addition to those prediction, a satellite experiment named TAO[13] will be setup near the reactor
235 core Taishan-1 to measure with an energy resolution of 2% at 1 MeV the neutrino flux coming from
236 the core, more details can be found in section 2.5.1. It will help identifying unknown fine structure
237 and give more insight on the $\bar{\nu}_e$ flux coming from this reactor.

238 2.1.2 Background spectra

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

- 241 — The geoneutrinos producing background in the $0.511 \sim 2.7$ MeV region.
- 242 — The neutrinos coming from the other nuclear reactors around Earth.

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

- 245 — The signal coming from radioactive decay (α , γ , β) from natural radioactive isotopes in the
 246 material of the detector.
- 247 — Cosmogenic event such as fast neutrons and activated isotopes induced by muons passing
 248 through the detector, most notably the spallation on ^{12}C .

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

250 2.2 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 2.2 and in the following sub-sections.

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 ν
Atmospheric neutrino	hundreds per year for ^{8}B	0–16 MeV	Radioactivity
Geoneutrino	hundreds per year	0.1–100 GeV	Negligible
	≈ 400 per year	0–3 MeV	Reactor ν

TABLE 2.2 – Detectable neutrino signal in JUNO and the expected signal rates and
 major background sources

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}U and ^{232}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.

260 Atmospheric neutrinos

261 Atmospheric neutrinos are neutrinos originating from the decay of π 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 from 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 [14].

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 from 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 num-
 296 ber. 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 $p \rightarrow \bar{\nu}K^+$ [15] thanks to the timing resolution of the SPMT system. Studies show
 300 that JUNO should be competitive with the current best limit at 5.9×10^{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 9.6×10^{33} years at 90 % C.L.

303 **2.3 The JUNO detector**

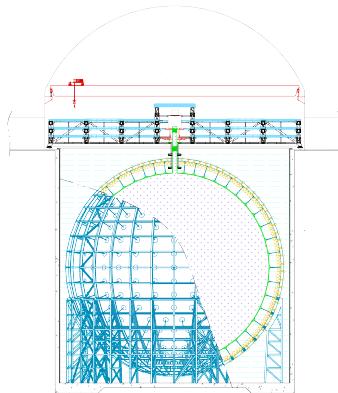
304 The JUNO detector is a scintillator detector buried 693.35 meters under the ground (1800 meters
 305 water equivalent). It consists of Central Detector (CD), a water pool and a Top Tracker (TT) as showed
 306 in figure 2.4a. The CD is an acrylic vessel containing the 20 ktons of Liquid Scintillator (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 ven-
 312 tilation system to get rid of the potential radon in the air. The CD is observed by two system of
 313 Photo-Multipliers Tubes (PMT). They are attached to the steel structure and their electronic readout
 314 is submersed 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 allow 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 allow two nuclear power plant to be used as sources.



(A) Schematics view of the JUNO detector.



(B) Top down view of the JUNO detector under construction

FIGURE 2.4

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

328 2.3.1 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}C and H nucleus in the LS:

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

329 Kinematics calculation shows that this interaction has an energy threshold for the $\bar{\nu}_e$ of $(m_n + m_e -$
 330 $m_p) \approx 1.806 \text{ MeV}$ [16]. This threshold make the experiment blind to very low energy neutrinos.
 331 The residual energy $E_\nu - 1.806 \text{ MeV}$ is be distributed as kinetic energy between the positron and the

332 neutron. The energy of the emitted positron E_e is given by [16]

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

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 2.3.2).
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.

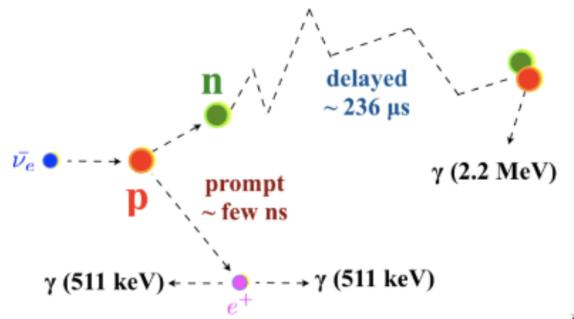


FIGURE 2.5 – Schematics of an IBD interaction in the central detector of JUNO

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 ($\sim 236 \mu s$).

347 2.3.2 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

processes is optimized to increase the assemblage transparency in water to >96%. The acrylic vessel is supported by a stainless steel structure via supporting node (fig 2.6). The structure and the nodes are designed to be resilient to natural catastrophic events such as earthquake and can support many times the effective load of the acrylic vessel.

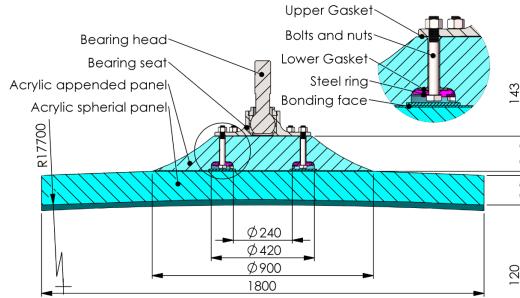


FIGURE 2.6 – Schematics of the supporting node for the acrylic vessel

364 Liquid scintillator

365 The Liquid Scintillator (LS) has a similar recipe as the one used in Daya Bay [17] 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 ~
371 98% of the LS, it is the main component with which ionizing particles and gamma interact.
372 Charged particles will collide with its electronic cloud transferring energy to the molecules,
373 gamma will interact via compton effect with the electronic cloud before finally be absorbed
374 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 [18]. 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 ~ 430 nm (full spectrum in figure 2.7) that
381 can thus be detected by our photo-multipliers systems. It amount for ~ 0.5% of the LS.

382 This formula has been optimized using dedicated studies with a Daya Bay detector [17, 20] to reach
383 the requirements for the JUNO experiment:

- 384 — A light yield / MeV of the amount of 10^4 photons to maximize the statistic in the energy
385 measurement.
- 386 — An attenuation length comparable to the size of the detector to prevent losing photons during
387 their propagation in the LS. The final attenuation length is 25.8m [21] to compare with the CD
388 diameter of 35.4m.
- 389 — Uranium/Thorium radiopurity to prevent background signal. The reactor neutrino program
390 require a contamination fraction $F < 10^{-15}$ while the solar neutrino program require $F <$
391 10^{-17} .

392 The LS will frequently be purified and tested in the Online Scintillator Internal Radioactivity In-
393 vestigation System (OSIRIS) [22] to ensure that the requirements are kept during the lifetime of the
394 experiment, more details to be found in section 2.5.2.

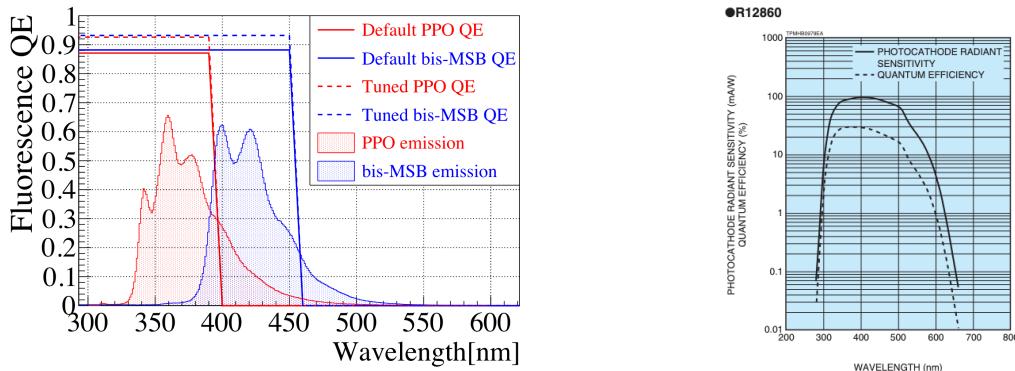


FIGURE 2.7 – On the left: Quantum efficiency (QE) and emission spectrum of the LAB and the bis-MSB [17]. On the right: Sensitivity of the Hamamatsu LPMT depending on the wavelength of the incident photons [19].

395 Large Photo-Multipliers Tubes (PMTs)

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

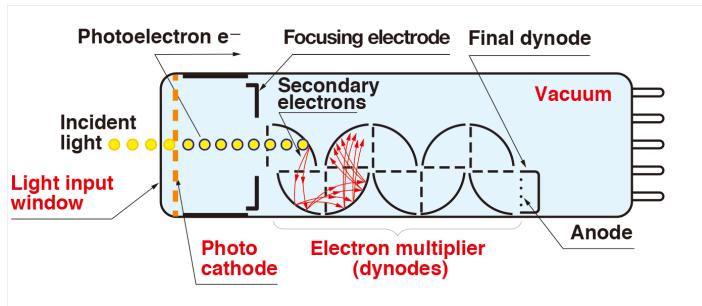


FIGURE 2.8 – Schematic of a PMT

403 The Large Photo-Multipliers Tubes (LPMT), used in the central detector and in the water pool, are
 404 20-inch (50.8 cm) radius PMTs. ~ 5000 dynode-PMTs [19] were produced by the Hamamatsu[©]
 405 company and ~ 15000 Micro-Channel Plate (MCP) [23] by the NNVT[©] company. This system is
 406 the one responsible for the energy measurement with a energy resolution of $3\%/\sqrt{E}$, resolution
 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 ~ 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 2.4 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)
- 425 — that complete the waveform with metadatas such as the local timestamp trigger, etc... Ths
- 426 — FPGA also act as a data buffer when needed by the DAQ and trigger system.
- 427 — Additional memory in order to temporally store the data in case of sudden burst of the input
- 428 — rate (such as in the case of nearby supernovae).

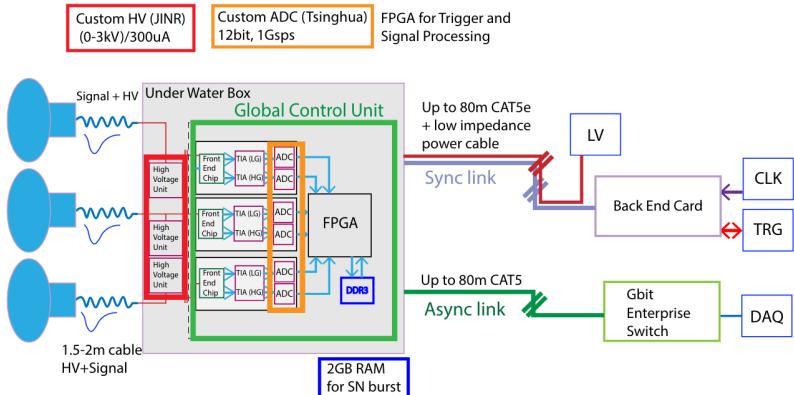


FIGURE 2.9 – The LPMT electronics scheme. It is composed of two part, the *wet* electronics on the left, located underwater and the *dry* electronics on the right. They are connected by Ethernet cable for data transmission and a dedicated low impedance cable for power distribution

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 ~ 43 PE/MeV for a final energy resolution of $\sim 17\%$.
 438 This resolution is not enough to measure the NMO, θ_{13} , Δm^2_{31} but will be sufficient to independently
 439 measure θ_{12} and Δm^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
 442 will explore in this thesis, for physics analysis. Due to the low PE rate, SPMTs will be running in
 443 photo-counting mode in the reactor range and thus will be insensitive to LPMT intrinsic effect (see
 444 section 2.4). Using this property, the intrinsic charge non linearity of the LPMTs can be measured by

445 comparing the PE count in the SPMTs and LPMTs [24]. Also, due to their smaller size and electronics,
 446 SPMTs have a better timing resolutions than the LPMTs. At higher energy range, like supernovae
 447 events, LPMTs will saturate where SPMTs due to their lower PE collection will to produce a reliable
 448 measure of the energy spectrum.

449 The SPMTs will be grouped by pack of 128 to an UWB hosting their electronics as illustrated in figure
 450 2.10. This underwater box host two high voltage splitter boards, each one supplying 64 SPMTs, an
 451 ASIC Battery Card (ABC) and a global control unit.

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

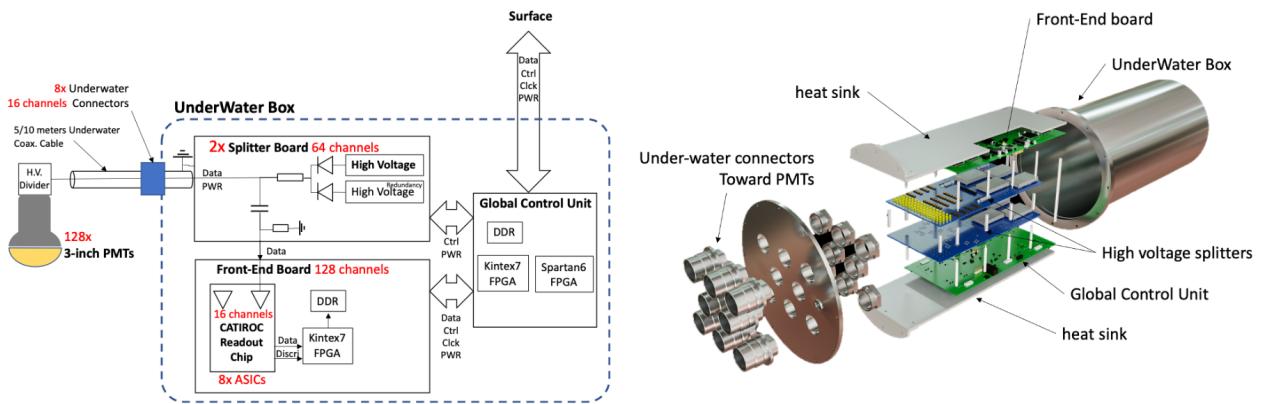


FIGURE 2.10 – Schematic of the JUNO SPMT electronic system (left), and exploded view of the main component of the UWB (right)

455 2.3.3 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 [25] 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.

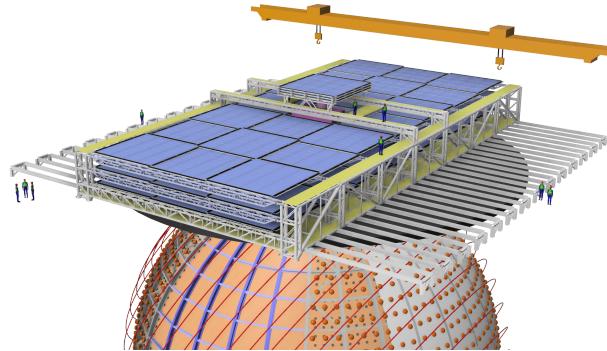


FIGURE 2.11 – The JUNO top tracker

476

477 **2.4 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,
 482 recognisable sources in the energy range of interest : [0-12] MeV for the positron signal and 2.2 MeV
 483 for the neutron capture. It is expected that the CD response will be different depending on the type
 484 of particle, due to the interaction with LS, the position on the event and the optical response of the
 485 acrylic sphere (see section 2.8). We also expect a non-linear energy response of the CD due to the LS
 486 properties [17] but also due to the reponse of the LPMTs system when collecting a large amount of
 487 PE [24].

488 **2.4.1 Energy scale calibration**

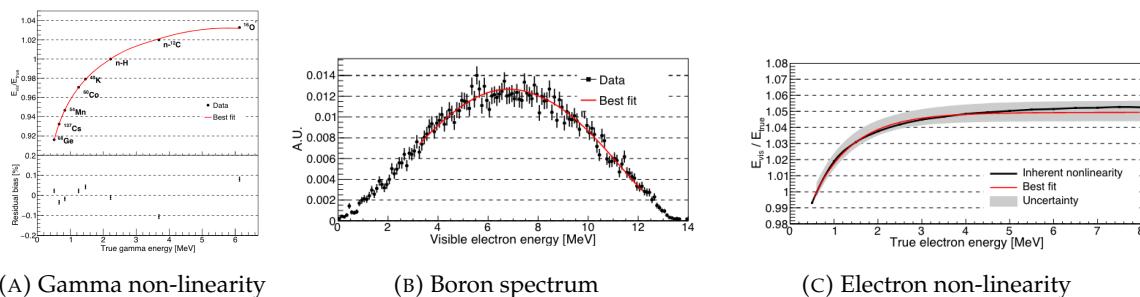
489 While electrons and positrons sources would be ideal, for a large LS detector thin-walled electrons
 490 or positrons sources could lead to leakage of radionucleides causing radioactive contamination.
 491 Instead, we consider gamma sources in the range of the prompt energy of IBDs. The sources are
 492 reported in table 2.3.

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

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

TABLE 2.3 – List of sources and their process considered for the energy scale calibration

497 From this calibration we call E_{vis} the "visible energy" that is reconstructed by our current algorithms
 498 and we compare it to the true energy deposited by the calibration source. The results shown in figure
 499 2.12 show the expected response of the detector from calibration sources. The non-linearity is clearly
 500 visible from the $E_{\text{vis}}/E_{\text{true}}$ shape. See [26] for more details.

FIGURE 2.12 – Fitted and simulated non linearity of gamma, electron sources and from the ^{12}B spectrum. Black points are simulated data. Red curves are the best fits. Figures taken from [26].

501 2.4.2 Calibration system

502 The non-uniformity due to the event position in the detector (more details in section 2.8) 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.

516 The preliminary calibration program is depicted in table 2.4.

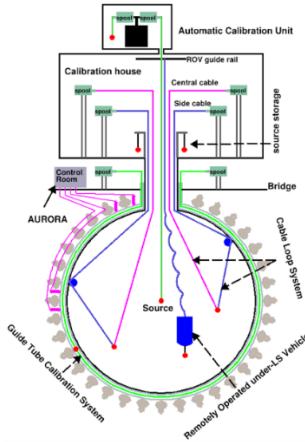


FIGURE 2.13 – Overview of the calibration system

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}Ge	ACU	75
	^{137}Cs	ACU	75
	^{54}Mn	ACU	75
	^{60}Co	ACU	75
	^{40}K	ACU	158

TABLE 2.4 – Calibration program of the JUNO experiment

517 2.4.3 Instrumental non-linearity calibration

518 One of the main interests of Dual Calorimetry is to calibrate away an instrumental effect called charge
 519 non linearity (QNL), which will be described in more detail in Chapter 7.

520 In short, during a typical IBD event, between 0 and 100 PEs can be produced in a given LPMT
 521 (depending on the position of the interaction and the positron energy). This is a large dynamic range.
 522 When the number of PEs is high, the reconstruction of the LPMT charge can become inaccurate,
 523 underestimating the actual number of PEs as illustrated in figure 2.14. This QNL is difficult to
 524 separate from other non linearities (like the non linearity in the LS photon yield as a function of
 525 the deposit energy). In chapter 5 and 6 of this thesis [24], a calibration method that constitutes the
 526 core of dual calorimetry are described. They are based on the comparisons between signals seen in
 527 LPMTs and signals seen in SPMTs. In the latter system, due to its small angular coverage, individual
 528 SPMT rarely see more than 1 PE per event, and therefore are essentially immune against QNL. The
 529 method described in [24] uses a tunable light source covering the range of 0 to 100 PE perLPMT
 530 channel

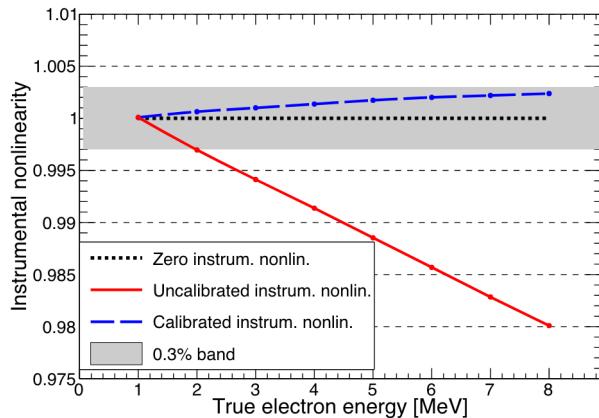


FIGURE 2.14 – Event-level instrumental non-linearity, defined as the ratio of the total measured LPMT charge to the true charge for events at the center of the detector. The solid red line represents event-level non-linearity without the channel-level correction in an extreme hypothetical scenario of 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 [26].

531 2.5 Satellite detectors

532 As introduced in section 2.1.1 and section 2.3.2, the precise knowledge and understanding of the
 533 detector condition is crucial for the measurements of the NMO and oscillation parameters. Thus two
 534 satellite detectors will be setup to monitor the experiment condition. TAO to monitor and understand
 535 the $\bar{\nu}_e$ flux and spectrum coming from the nuclear reactor and OSIRIS to monitor the LS response.

536 2.5.1 TAO

537 The Taishan Antineutrino Observatory (TAO) [13, 27] is a ton-level gadolinium doped liquid scin-
 538 tillator detector that will be located near the Taishan-1 reactor. It aim to measure the $\bar{\nu}_e$ spectrum at
 539 very low distance (44m) from the reactor to measure a quasi-unoscillated spectrum. TAO also aim to
 540 provide a major contribution to the so-called reactor anomaly [12]. Its requirement are to the level of
 541 2 % energy resolution at 1 MeV.

542 Detector

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

551 2.5.2 OSIRIS

552 The Online Scintillator Internal Radioactivity Investigation System (OSIRIS) [22] is an ultralow back-
 553 ground, 20 m^3 LS detector that will be located in JUNO cavern. It aim to monitor the radioactive
 554 contamination, purity and overall response of the LS before it is injected in JUNO. OSIRIS will
 555 be located at the end of the purification chain of JUNO, monitoring that the purified LS meet the
 556 JUNO requirements. The setup is optimized to detect the fast coincidences decay of $^{214}\text{Bi} - ^{214}\text{Po}$
 557 and $^{212}\text{Bi} - ^{212}\text{Po}$, indicators of the decay chains of U and Th respectively.

558 **Detector**

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

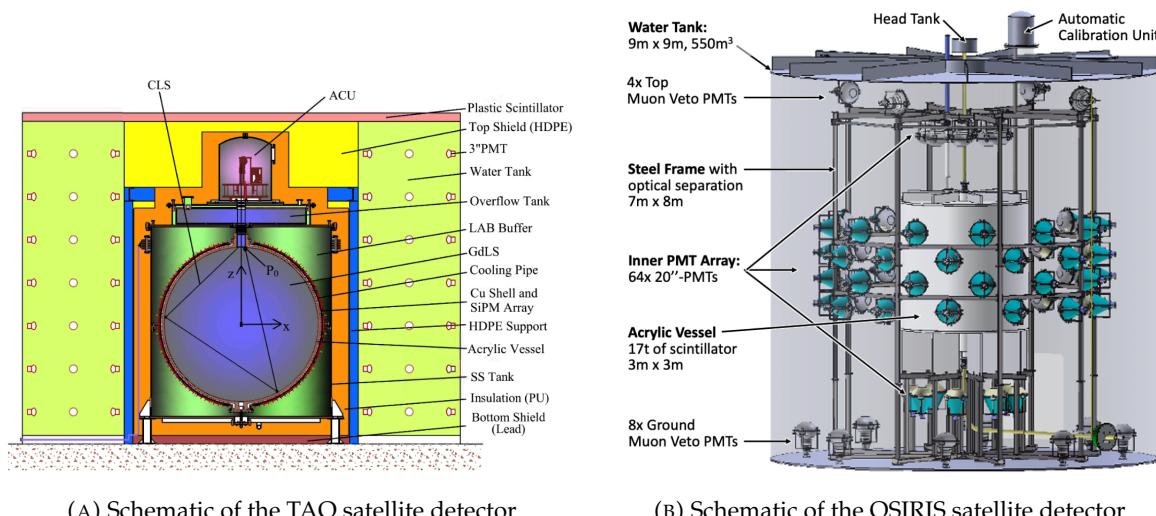


FIGURE 2.15

565 2.6 Software

566 The simulation, reconstruction and analysis algorithms are all packaged in the JUNO software,
 567 subsequently called the software. It is composed of multiple components integrated in the SNiPER
 568 [28] framework:

- 569 — Various primary particles simulators for the different kind of events, background and calibra-
 570 tion sources.
- 571 — A Geant4 [29–31] Monte Carlo (MC) simulation containing the detectors geometries, a custom
 572 optical model for the LS and the supporting structures of the detectors. The Geant4 simulation
 573 integrate all relevant physics process for JUNO, validated by the collaboration. This step of the
 574 simulation is commonly called *Detsim* and compute up to the production of photo-electrons

575 in the PMTs. The optics properties of the different materials and detector components have
 576 been measured beforehand to be used to define the material and surfaces in the simulation.

- 577 — An electronic simulation, simulating the response waveform of the PMTs, tracking it through
 578 the digitization process, accounting for effects such as non-linearity, dark noise, Time Trans-
 579 it Spread (TTS), pre-pulsing, after-pulsing and ringing of the waveform. It's also the step
 580 handling the event triggers and mixing. This step is commonly referenced as *Elecsim*.
- 581 — A waveform reconstruction where the digitized waveform are filtered to remove high-frequency
 582 white noise and then deconvoluted to yield time and charge informations of the photons hits
 583 on the PMTs. This step is commonly referenced as *Calib*.
- 584 — The charge and time informations are used by reconstruction algorithms to reconstruct the
 585 interaction vertex and the deposited energy. This step is commonly reported as *Reco*. See
 586 section 2.8 for more details on the reconstruction.
- 587 — Once the singular events are reconstructed, they go through event pairing and classification
 588 to select IBD events. This step is named Event Classification.
- 589 — The purified signal is then analysed by the analysis framework which depend of the physics
 590 topic of interest. An introduction to the reactor $\bar{n}u_e$ is presented in section 2.7.

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

596 2.7 Reactor anti-neutrino oscillation analysis

597 2.7.1 IBD samples selection

598 The $\bar{\nu}_e$ coming from nuclear reactor will, for the most part, interact with proton, hydrogen nucleus,
 599 via Inverse Beta Decay (IBD). The first step of the oscillation analysis is to constitute a sample of IBD
 600 candidates, dominated by actual IBDs. The IBD interaction, schematised in figure 2.5, will produce
 601 two particle, with differentiable signals.

602 The first signal comes from the positron slowdown and its annihilation with an electron of the LS.
 603 This is the *prompt* signal, happening a few ns after the IBD. The positron takes most of the $\bar{\nu}_e$ kinetic
 604 energy, as detailed in section 2.3.1.

605 The leftover kinetic energy is taken by the neutron that, after thermalisation in the LS, will be
 606 captured by an hydrogen and produce a 2.2 MeV gamma, or by a carbon emitting a 4.9 MeV gamma.
 607 This is the *delayed* signal, happening $\sim 236 \mu\text{s}$ after the IBD. This second mono-energetic event serve
 608 as a marker for the IBD.

609 The IBD selection is thus based on the selection of a prompt event, with an energy between 0.8 and
 610 12 MeV, and a delayed event with an energy in the ranges [1.9, 2.5] MeV or [4.4, 5.5] MeV. Those two
 611 signal needs to be in a 1 ms time window and within 1.5 m from each other. Additionally the two
 612 signal needs to be in a radius of 17.2m from the detector center (0.5 m from the edge) to protect from
 613 accidental background formed by two uncorrelated signals [32]. Those values will be further refined
 614 after once JUNO data-taking starts.

615 In addition, specials veto are setup to protect from cosmic muons and their aftermath. The details of
 616 those veto and selection can be found in [32].

617 The expected rate and selection efficiency on IBD can be found in table 2.5. After these selection, the
 618 residual background, including $\bar{\nu}_e$ coming from other sources than the reactor can be found in table
 619 2.6.

Selection Criterion	Efficiency [%]	IBD Rate [day ⁻¹]
All IBDs	100.0	57.4
Fiducial Volume	91.5	52.5
IBD Selection	98.1	51.5
Energy Range	99.8	-
Time Correlation (ΔT_{p-d})	99.0	-
Spatial Correlation (ΔR_{p-d})	99.2	-
Muon Veto (Temporal + Spatial)	91.6	47.1
Combined Selection	82.2	47.1

TABLE 2.5 – Summary of cumulative reactor antineutrino selection efficiencies. The reported IBD rates (with baselines <300 km) refer to the expected events per day after the selection criteria are progressively applied. Table taken from [32]

Backgrounds	Rate [day ⁻¹]	B/S [%]
Geoneutrinos	1.2	2.5
World reactors	1.0	2.1
Accidentals	0.8	1.7
⁹ Li/ ⁸ He	0.8	1.7
Atmospheric neutrinos	0.16	0.34
Fast neutrons	0.1	0.21
¹³ C(α, n) ¹⁶ O	0.05	0.01
Total backgrounds	4.11	8.7

TABLE 2.6 – Expected background rates, background to signal ratio (B/S), and rate and shape uncertainties. The B/S ratio is calculated by using the IBD signal rate of 47.1/day. Table taken from [32]

Once a sample is obtained, the oscillation analysis will consist essentially on the fit of a spectrum model to the spectrum observed in the selected sample. More specifically, the spectrum under analysis is the spectrum of the reconstructed visible energy of the positron : E_{vis}^{vis} . The reconstruction is presented in detail in section 2.8. For 6 years of data taking, it will resemble that on figure 2.3. In the next sections, I describe the fit procedures developed in JUNO. This will be the occasion to introduce notions useful for Chapter 7. Besides, I'll also describe the versions of the fit used in this Chapter 7.

2.7.2 Synthetic overview of fit procedures developed at JUNO

Several fit procedures are being developed by JUNO collaborators (half a dozen of groups work in parallel within the collaboration). We do not have the ambition of a thorough description here. Instead, we try to introduce the main elements useful to the reader to understand JUNO's future results, and the fit procedures used Chapter 7.

In most cases, the fit is a binned fit to the histogrammed spectrum of E_{vis}^{vis} , like the one in figure 2.3. It is based on the minimization of a χ^2 test statistics. Generically, it can be written this way :

$$\chi^2 = (\mathbf{T}(\boldsymbol{\theta}, \boldsymbol{\eta}) - \mathbf{D})^T \mathbf{V}^{-1} (\mathbf{T}(\boldsymbol{\theta}, \boldsymbol{\eta}) - \mathbf{D}) + \chi^2_{nuis}(\boldsymbol{\eta}) \quad (2.3)$$

where the components of data vector \mathbf{D} are the number of events found in individual bins of the fitted histogram, $\mathbf{T}(\boldsymbol{\theta}, \boldsymbol{\eta})$ is the vector of the predicted number of entries in each bins. This prediction is the integration over the width of the bins of the spectrum model for a given NMO (described latter in this section).

638 This model depends on the oscillation parameters $\theta = (\Delta m_{21}^2, \sin^2(2\theta_{12}), \Delta m_{31}^2, \sin^2(2\theta_{13}))$, and on
 639 nuisance parameters η involved in the fit model and associated with systematic uncertainties. Uncer-
 640 tainties are treated in two ways : statistical and some of the systematic uncertainties are accounted
 641 for via the covariance matrix $V = V_{stat} + V_{syst}$; remaining systematic uncertainties are treated via the
 642 penalty term χ^2_{nuis} , which is written this way :

$$\chi^2_{nuis}(\eta) = (\eta - \bar{\eta})^T \cdot V_\eta^{-1}(\eta) \cdot (\eta - \bar{\eta}) \quad (2.4)$$

643 where $\bar{\eta}$ is the vector containing the most probable values of the nuisance parameters according to
 644 our knowledge prior to the fit, and where V_η is the covariance matrix accounting of the uncertainty
 645 on these values, and the potential correlations between them. In principles, a likelihood could be
 646 used instead of a χ^2 . However, some of the systematic uncertainties are not trivial to parameterize,
 647 therefore treating them as nuisance parameters in not trivial.

648 An example of nuisance parameters are the A , B and C parameters of equation 7.19, which can be
 649 used to describe the resolution on the reconstructed energy. The fit model leading to $T(\theta, \eta)$ indeed
 650 incorporates this resolution.

651 Treatment of uncertainties

652 Differences between various fit procedures developed within JUNO often lies in the choice of the sys-
 653 tematic uncertainties that are treated via V or $\chi^2_{nuis}(\eta)$. Among the reasons behind these differences
 654 is the necessity to compare several approaches to ensure the robustness JUNO's oscillation analysis
 655 results. This approach was already adopted in the recent evaluations of JUNO's potential [3, 32].
 656 Studies carried out so far at Subatech assumes a treatment entirely via V .

657 Other differences lies in the choice of the way to evaluate V_{stat} . Two common approaches used in
 658 χ^2 fit are the Neyman and the Pearson approaches. If the size of the fitted sample is high enough,
 659 the variation of D_i , the number of entries in bin i , around its true expectation value \bar{D}_i is $\sqrt{\bar{D}_i}$.
 660 To evaluate this number, the Neyman approach uses simply the number of entries observed in the
 661 sample under analysis : $\sqrt{D_i}$. The Pearson approach uses the prediction by the fit model : $\sqrt{T(\theta, \eta)_i}$.

662 Both cases are approximations which lead to biases that are not tolerable given the precision JUNO
 663 must aim at for a successful oscillation analysis. To reduce this bias, most of JUNO groups employ the
 664 "Combined Neyman Pearson" approach introduced in [33]. Schematically, it consists on combining
 665 both approaches : $(V_{stat})_{ii} = 3 / \left(\frac{1}{T(\theta, \eta)_i} + \frac{2}{D_i} \right)$. Weights in this relation are chosen in order to cancel
 666 typical biases. The validity of this method is not guaranteed universally. In particular, limitations
 667 appear when a complex systematic matrix V_{syst} is added to V_{stat} .

668 This is the case in the approach followed at Subatech, were all sources of systematic uncertainties
 669 are treated via this matrix. Dedicated studies run at Subatech observed biases in the fitted oscillation
 670 parameters using CNP in this case. Subatech's group therefore adopted another approach (verified
 671 to be unbiased).

672 Originally, fitting the E_{vis}^{e+} spectrum should mean maximising a likelihood, equal to the product over
 673 all bins of the probabilities to find D_i in bin i . With a large enough samples, this product tends to a
 674 multidimensional gaussian (one dimension per bin) :

$$\mathcal{L} = 2\pi^{-\frac{N}{2}} |V|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{D} - \mathbf{T}(\theta, \eta))^T V^{-1}(\mathbf{D} - \mathbf{T}(\theta, \eta))} \quad (2.5)$$

675 Replacing \mathcal{L} by $-2 \ln \mathcal{L}$ one obtains :

$$\chi^2_{PV} = (\mathbf{T}(\theta, \eta) - \mathbf{D})^T V^{-1} (\mathbf{T}(\theta, \eta) - \mathbf{D}) + \ln(|V|) \quad (2.6)$$

676 where V is the total covariance matrix with its statistical component evaluated according to the
 677 Pearson approach. The $\ln |V|$ term, often neglected in χ^2 fits, ensures that biases, essentially related
 678 to the normalisation of the fitted distribution, are avoided. This "PearsonV" χ^2 is the one that we
 679 minimize in the fits used in Chapter 7.

680 Another difference between the various procedures developed at JUNO is the choice of the spectrum
 681 range and binning. So far, at Subatech, we use an histogram defined between 0.8 and 9 MeV, and a
 682 regular binning involving 20 keV wide bins.

683 Joint fit of JUNO and TAO spectra

684 Another difference between the various fit procedures developed in the collaboration is the inclusion
 685 of the data collected by TAO (see section 2.5.1). The spectrum prediction $T(\theta, \eta)$ involves predictions
 686 on the differential flux of $\bar{\nu}_e$ as a function of $E_{\bar{\nu}_e}$ produced in reactors. This is one of the main
 687 systematic uncertainties affecting the oscillation analysis. This can be constrained using the data
 688 of TAO. An efficient way to use them is via a simultaneous fit, which will constrain the part of the η
 689 parameters related to the reactor predictions. In this case, equation 2.3 becomes :

$$\chi^2 = \sum_d \left(T^d(\theta^d, \eta) - D^d \right)^T V^{-1} \left(T^d(\theta^d, \eta) - D^d \right) + \chi^2_{nuis}(\eta) \quad (2.7)$$

690 where the d superscript stands for the spectrum measured in JUNO or TAO.

691 Finally, it must be noted that JUNO's sensitivity to $\sin^2(2\theta_{13})$ is too weak for a competitive measure-
 692 ment. In most versions of the oscillation analyses carried out within JUNO, it will be considered as a
 693 nuisance parameter. In practice, the various χ^2 's presented earlier will receive an additional term :

$$\chi^2_{\sin^2(2\theta_{13})} = \frac{(\overline{\sin^2(2\theta_{13})} - \overline{\sin^2(2\theta_{13})})^2}{\sigma^2_{\sin^2(2\theta_{13})}} \quad (2.8)$$

694 where $\overline{\sin^2(2\theta_{13})}$ and the denominators can be provided, for instance, by the world average on this
 695 parameter.

696 2.7.3 The spectrum model and sources of systematic uncertainties

697 The E_{vis}^{e+} spectrum observed in data (Fig 2.3) is the sum of the IBD spectrum and of the various
 698 backgrounds spectra (see table 2.6). The spectrum prediction $T(\theta, \eta)$ is therefore the sum of IBD and
 699 backgrounds predictions. The latter are provided by MC simulations. The former results from the
 700 theoretical description of the series of phenomena that lead to the observed IBD spectrum. In a given
 701 bin i , it can be expressed this way :

$$T^i(\theta, \eta) = \sum_j C_{ij}^{E_{rec}} \int_{E_j^{vis}}^{E_{j+1}^{vis}} dE^{vis} \int_{-1}^1 d\cos\theta \Phi(E^\nu) \frac{d\sigma}{d\cos\theta}(E^\nu, \cos\theta) \frac{dE^\nu}{dE^{dep}} \frac{dE^{dep}}{dE^{vis}} \quad (2.9)$$

702 In the above equation, 4 kinds of energies appears: following the IBD, the antineutrino energy E^ν is
 703 quasi entirely transferred to the positron, of energy E_e . It eventually annihilates, so the actual energy
 704 released in the LS is E_{dep} , which includes the mass of the annihilated electron. The production optical
 705 photons is not linear in E_{dep} (see section 2.4), so that the visible energy (that will be reconstructed) is
 706 E_{vis} . This reconstruction comes with resolution effects, leading to E_{rec} .

707 Equation 2.9 describe the passage from the original differential flux (as a function of E^ν) of antineu-
 708 trinos reaching the detector to the reconstructed spectrum:

- $\Phi(E^\nu)$ is the differential antineutrino flux reaching JUNO.
- $\frac{d\sigma}{d\cos\theta}(E^\nu, \cos\theta)$ account for the IBD cross section, which depends on the antineutrino energy and on the incidence angle.
- The last two terms of the integrand are the differential relations linking E^ν , E^{dep} and E^{vis} .
- Reconstruction effects are described via C_{ij}^{rec} 's, that make the link between the true and reconstructed visible energy. In a simple case, it is equivalent to a convolution product. The matrix formalism here prepares the fact that a realistic analysis might employ a more empirical way, based on MC.

The differential flux is expressed this way:

$$\Phi(E^\nu) = \sum_r \left(\frac{\mathcal{P}_{\bar{\nu}_e \rightarrow \bar{\nu}_e}(E^\nu, L_r)}{4\pi L_r^2} \frac{W_r}{\sum_i f_{i,r} e_i} \sum_i f_{i,r} s_i(E^\nu) \right) \quad (2.10)$$

where:

- $\mathcal{P}_{\bar{\nu}_e \rightarrow \bar{\nu}_e}(E^\nu, L_r)$ is the antineutrino survival probability at distance L_r from the production point in reactor r , dictated by the oscillation probability.
- e_i stands for the mean energy released per fission for isotope i .
- W_r is the thermal power of reactor r .
- $f_{i,r}$ is the fission fraction in reactor r of isotope i among the four.
- $s_i(E^\nu)$ is the $\bar{\nu}_e$ energy spectrum - at emission point - per fission for each isotope, as emitted by the reactor.

Sources of systematic uncertainties

The numerous quantities appearing in the spectrum model embody a good part of the systematic uncertainties. Among the leading contributions are those related to the knowledge of the reactor related quantities. Of importance are also the uncertainties related to the modelling of the non linearity of the photon emission (passage from E^{dep} to E^{vis}) and of the reconstruction resolution. The shape and rate of the backgrounds are also a leading source of systematic uncertainties. The uncertainty on IBD selection efficiency also has a notable role.

Sensitivities to NMO and oscillation parameters

JUNO will start taking data in 2025. During the months and years to come, oscillation analyses will naturally be optimized regularly. What we described here represent the state of the art mid 2024, and was used for the sensitivity studies published in [3, 32] and are presented in table 2.7

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

TABLE 2.7 – A summary of precision levels for the oscillation parameters. The reference value (PDG 2020 [34]) is compared with 100 days, 6 years and 20 years of JUNO data taking.

739 **Asimov studies**

740 To study the behavior and performance of fit procedures with enough realism, one should perform
 741 fits to a large number of toy spectra, generated with a number events equal to what one expects in
 742 real data, for the given exposure under consideration. This allows to study the impact of realistic
 743 statistical fluctuations. This is, however, time consuming, since thousands of spectra have to be
 744 generated and fitted.

745 When subtle details are not crucial, another approach is possible to estimate sensitivities to the NMO
 746 and oscillation parameters, as well as (for instance) to verify the technical implementation of fitter
 747 (as we will do in Chapter 7 for the implementation of the joint fit). It consists on generating only 1
 748 pseudo-data sample, where the content of each bin D^i is set to the predicted value T^i , computed with
 749 a reasonable choice for the values of the model parameters (for instance, with the recent PDG values
 750 for the oscillation parameters). This is equivalent to a spectrum with fluctuations. It provides valid
 751 sensitivities if the expected statistics in the real data sample is high enough in each bin to assume a
 752 gaussian behavior.

753 **2.7.4 Versions of the fit used in this thesis**

754 In Chapter 7, we'll study the potential of a particular application of Dual Calorimetry, call "Dual
 755 Calorimetry with neutrino oscillation." This approach require to perform fits to the E^{vis} spectrum
 756 reconstructed with the LPMT system, with the SPMT system, and a joint fit to both spectra.

757 In the two former cases, the PearsonV χ^2 introduced above will be used. In the latter case, it will
 758 be extended in the following way : The D data vector now possess 820 elements. Indeed, the fit is
 759 performed to a joint spectrum, where the LPMT spectrum is juxtaposed with the SPMT spectrum
 760 (see figure 2.16).

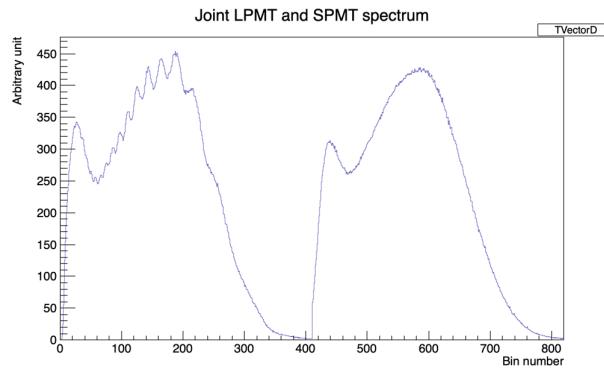


FIGURE 2.16 – Illustration of the spectrum considered when joint fitting

761 The prediction vector $T(\theta^d, \eta)$ is naturally extended in the same way. Its components 1 to 410 pre-
 762 dict the number of entries in the LPMT part of the LPMT+SPMT joint spectrum, while its components
 763 from 411 to 820 predict the contents of the SPMT part. Note that the list of oscillation parameters
 764 in $T_{411}(\theta^d, \eta)$ to $T_{820}(\theta^d, \eta)$ is the same as usual. However, $T_1(\theta^d, \eta)$ to $T_{410}(\theta^d, \eta)$ 2 additional
 765 parameters, $\delta(\sin^2(2\theta_{12}))$ and $\delta(\Delta m_{21}^2)$, are added to the corresponding oscillation parameters to
 766 account for a potential unexpected problem in the LPMT reconstruction or calibration.

767 In the case of this joint fit, the covariance matrix V is extended to a (820×820) matrix. It is a central
 768 element of this study, as will be explained in Chapter 7, since the LPMT and SPMT data spectrum
 769 are correlated, even at the statistical level. The determination of this matrix will be an important and
 770 original point.

771 Fits will be performed to an histogram spectrum defined over the 0.8-9 MeV range, with a flat binning
 772 (20 keV wide bins), often restricted to the 335 lowest E^{vis} bins.

773 In this section 2.7, we have provided a theoretical description of the fit procedures developed at
 774 JUNO. Software frameworks are necessary to use them in practice. The framework developed at
 775 Subatech will be described in Chapter 7.

776 2.8 State of the art of the Offline IBD reconstruction in JUNO

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

781 2.8.1 Interaction vertex reconstruction

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

784 Charge based algorithm

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

$$\vec{r}_{cb} = a \cdot \frac{\sum_i q_i \cdot \vec{r}_i}{\sum_i q_i} \quad (2.11)$$

786 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
 787 reconstructed interaction position. a is a scale factor introduced because a weighted average over
 788 a 3D sphere is inherently biased. Using calibration we can estimate $a \approx 1.3$ [37]. The results in
 789 figure 2.17b shows that the reconstruction is biased from around 15m and further. This is due to the
 790 phenomena called “total reflection area” or TR Area.

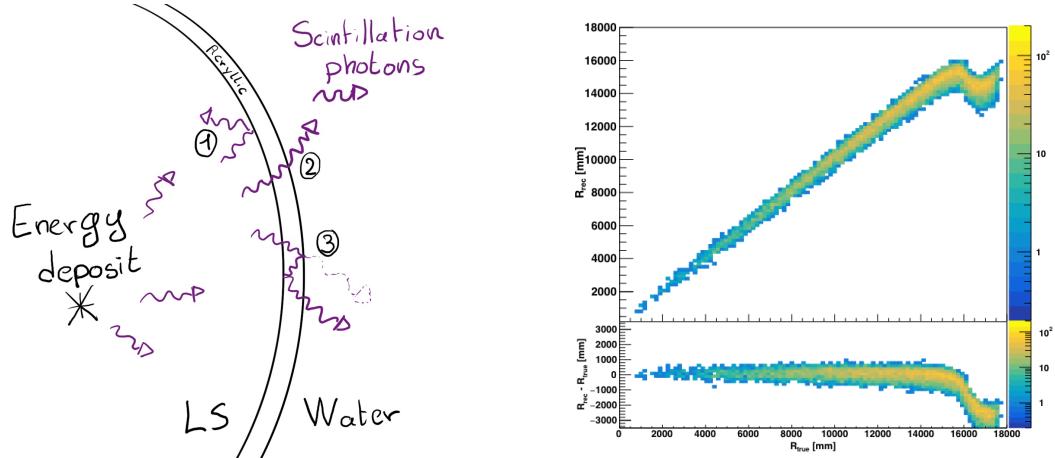
791 As depicted in the figure 2.17a the optical photons, given that they have a sufficiently large incidence
 792 angle, can be deviated of their trajectories when passing through the interfaces LS-acrylic and water-
 793 acrylic due to the optical index difference. This cause photons to be lost or to be detected by PMT
 794 further than anticipated if we consider their rectilinear trajectories. This cause the charge barycenter
 795 the be located closer to the center than the event really is.

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

799 Time based algorithm

800 The time based algorithm use the distribution of the time of flight corrections Δt (Eq 2.12) of an event
 801 to reconstruct its vertex and t_0 . It follow the following iterations:

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



(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 2.17

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

$$\Delta t_i(j) = t_i - \text{tof}_i(j) \quad (2.12)$$

804 where j is the iteration step, t_i is the timing of the i th PMT, and tof_i is the time-of-flight of the
805 photon considering an rectilinear trajectory and an effective velocity in the LS and water (see
806 [37] for detailed description of this effective velocity). Plot the Δt distribution and label the
807 peak position as Δt^{peak} (see fig 2.18a).

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

$$\vec{\delta}[\vec{r}(j)] = \frac{\sum_i \left(\frac{\Delta t(j) - \Delta t^{\text{peak}}(j)}{\text{tof}_i(j)} \right) \cdot (\vec{r}_0(j) - \vec{r}_i)}{N^{\text{peak}}(j)} \quad (2.13)$$

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

813 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
814 step 2.

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

$$t'_i = t_i - p_0 / \sqrt{N_i^{\text{pe}}} - p_1 - p_2 / N_i^{\text{pe}} \quad (2.14)$$

818 The parameters (p_0, p_1, p_2) were optimized to (9.42, 0.74, -4.60) for Hamamatsu PMTs and (41.31,
819 -12.04, -20.02) for NNVT PMTs [37]. The results presented in figure 2.18b shows that the time based
820 algorithm provide a more accurate vertex and is unbiased even in the TR area. This results (\vec{r}_0, t_0) is
821 used as initial value for the likelihood algorithm.

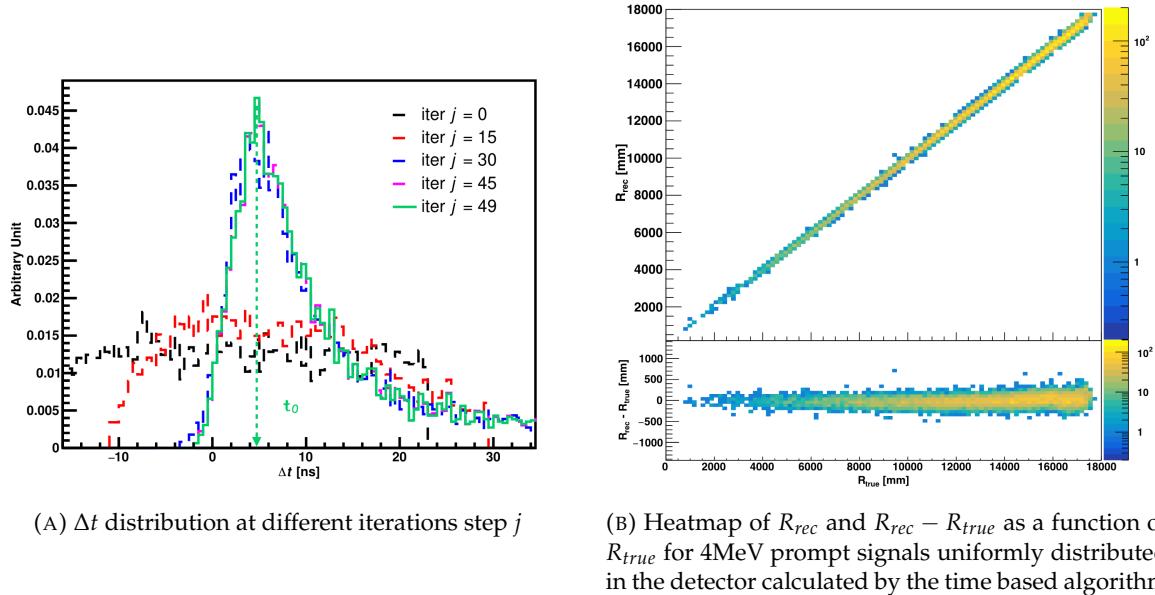


FIGURE 2.18

822 Time likelihood algorithm

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

$$t_{\text{res}}^i(\vec{r}_0, t_0) = t_i - \text{tof}_i - t_0 \quad (2.15)$$

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

$$f(t_{\text{res}}) = \sum_k \frac{\rho_k}{\tau_k} e^{-\frac{t_{\text{res}}}{\tau_k}}, \sum_k \rho_k = 1 \quad (2.16)$$

827 as the sum of the k component that emit light in the LS each one characterised by it's decay time τ_k
828 and intensity fraction ρ_k . The TTS component is expressed as a gaussian convolution

$$g(t_{\text{res}}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_{\text{res}}-\nu)^2}{2\sigma^2}} \cdot f(t_{\text{res}}) \quad (2.17)$$

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

$$p(t_{\text{res}}) = (1 - \epsilon_{dn}) \cdot g(t_{\text{res}}) + \epsilon(t_{\text{res}}) \quad (2.18)$$

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

$$\mathcal{L}(\vec{r}_0, t_0) = -\ln \left(\prod_i p(t_{\text{res}}^i) \right) \quad (2.19)$$

The parameter of Eq. 2.18 can be measured experimentally. The results shown in figure 2.19 used PDF from monte carlo simulation. The results shows that $R_{rec} - R_{true}$ is biased depending on the energy. While this could be corrected using calibration, another algorithm based on charge likelihood was developed to correct this problem.

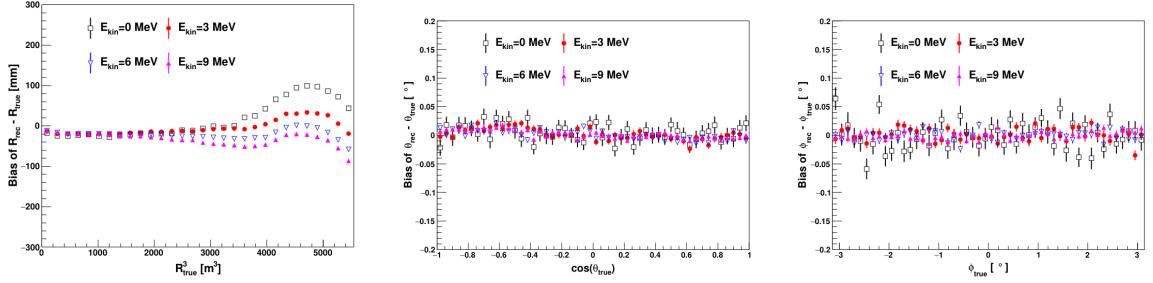


FIGURE 2.19 – Bias of the reconstructed radius R (left), θ (middle) and ϕ (right) for multiple energies by the time likelihood algorithm

Charge likelihood algorithm

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

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

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

$$P(\vec{r}_0, E) = \prod_j P_{nohit}^j(\vec{r}_0, E) \cdot \prod_i P_{hit}^i(\vec{r}_0, E) \quad (2.20)$$

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

$$\mathcal{L}(\vec{r}_0, E) = -\ln(P(\vec{r}_0, E)) \quad (2.21)$$

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

$$\mu_i(\vec{r}_0, E) = Y \cdot \frac{\Omega(\vec{r}_0, r_i)}{4\pi} \cdot \epsilon_i \cdot f(\theta_i) \cdot e^{-\sum_m \frac{d_m}{\zeta_m}} \cdot E + \delta_i \quad (2.22)$$

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

However Eq. 2.22 assume that the scintillation light yield is linear with energy and describe poorly the contribution of indirect light, shadow effect due to the supporting structure and the total reflection effects. The solution is to use data driven methods to produce the pdf by using the calibrations sources and position described in section 2.4. In the results presented in figures 2.20, the PDF was 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 resolution than the time likelihood. The figure 2.21 shows the radial resolution of the different

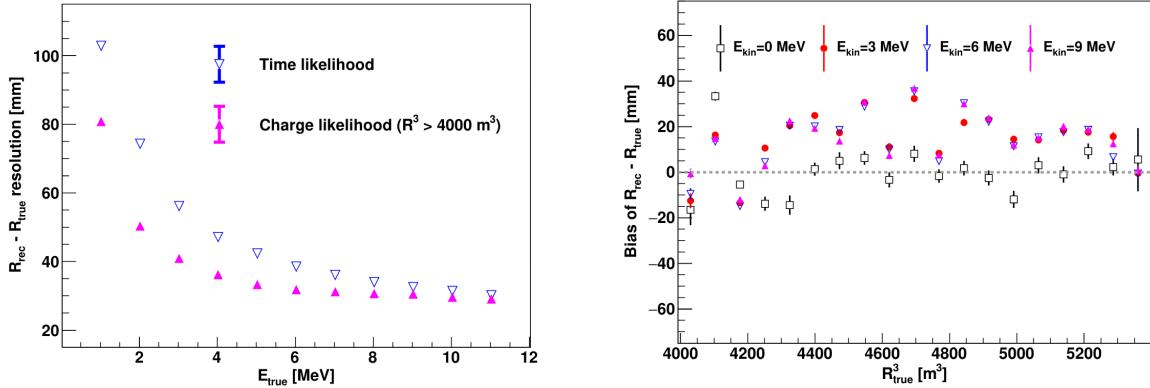


FIGURE 2.20 – On the left: Resolution of the reconstructed R as a function of the energy in the TR area ($R^3 > 4000 \text{ m}^3 \equiv R > 16 \text{ m}$) by the charge and time likelihood algorithms. On the right: Bias of the reconstructed R in the TR area for different energies by the charge likelihood algorithm

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

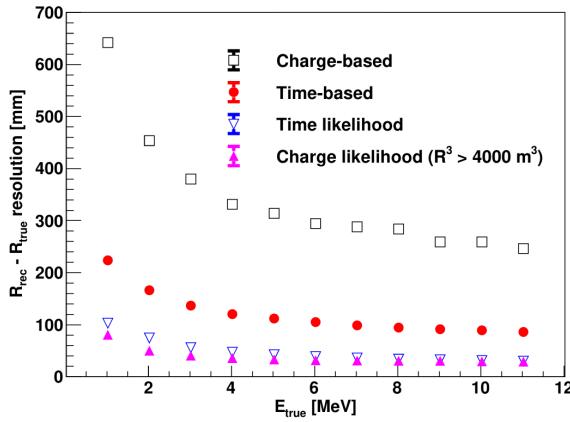


FIGURE 2.21 – Radial resolution of the different vertex reconstruction algorithms as a function of the energy

The charge based likelihood algorithms already give use some information on the energy as Eq. 2.21 is minimized but the energy can be further refined as shown in the next section.

2.8.2 Energy reconstruction

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

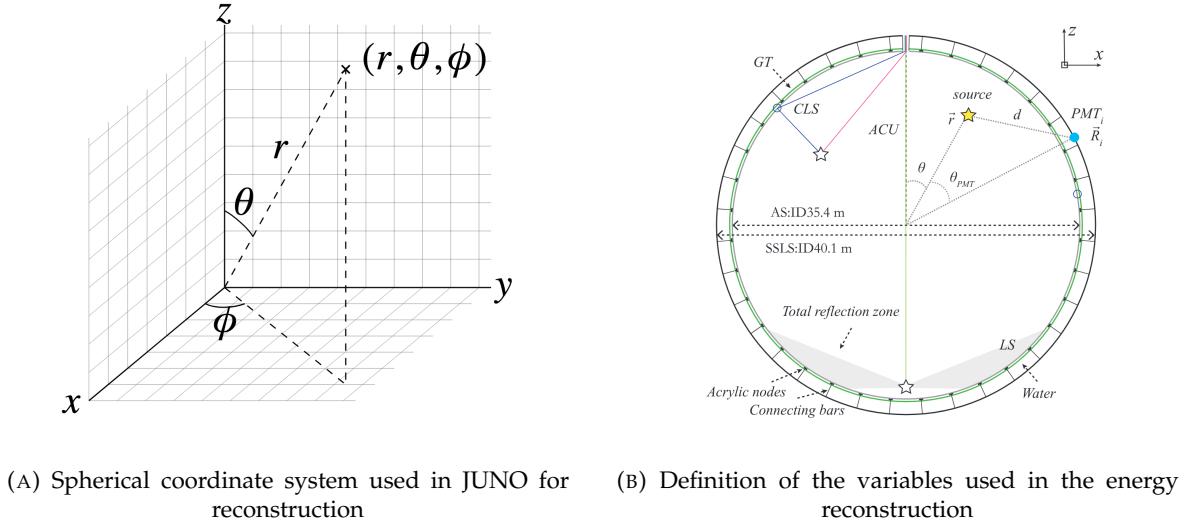


FIGURE 2.22

867 Charge estimation

868 The most important element in the energy reconstruction is $\mu_i(\vec{r}_0, E)$ described in Eq. 2.22. For
 869 realistic cases, we also need to take into account the electronics effect that were omitted in the
 870 previous section. Those effect will cause a charge smearing due to the uncertainties in the N_{pe}
 871 reconstruction. Thus we define $\hat{\mu}^L(\vec{r}_0, E)$ which is the expected N_{pe}/E in the whole detector for an
 872 event with visible energy E_{vis} and position \vec{r}_0 . The position of the event and PMTs are now defined
 873 using $(r, \theta, \theta_{pmt})$ as defined in figure 2.22b.

$$\hat{\mu}(r, \theta, \theta_{pmt}, E_{vis}) = \frac{1}{E_{vis}} \frac{1}{M} \sum_i^M \frac{\bar{q}_i - \mu_i^D}{DE_i}, \quad \mu_i^D = DNR_i \cdot L \quad (2.23)$$

874 where i runs over the PMTs with the same θ_{pmt} , DE_i is the detection efficiency of the i th PMT. μ_i^D
 875 is the expected number of dark noise photoelectrons in the time window L . The time window have
 876 been optimized to $L = 280$ ns [41]. \bar{q}_i is the average recorded photoelectrons in the time window
 877 and \hat{Q}_i is the expected average charge for 1 photoelectron. The N_{pe} map is constructed following the
 878 procedure described in [36].

879 Time estimation

880 The second important observable is the hit time of photons that was previously defined in Eq. 2.15.
 881 It is here refined as

$$t_r = t_h - \text{tof} - t_0 = t_{LS} + t_{TT} \quad (2.24)$$

882 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
 883 by a gaussian

$$t_{TT} = \mathcal{N}(\mu_{TT} + t_d, \sigma_{TT}) \quad (2.25)$$

884 where μ_{TT} is the mean transit time in PMTs, σ_{TT} is the Transit Time Spread (TTS) of the PMTs and t_d
 885 is the delay time in the electronics. The effective refraction index of the LS is also corrected to take
 886 into account the propagation distance in the detector.

887 The timing PDF $P_T(t_r | r, d, \mu_l, \mu_d, k)$ can now be generated using calibration sources [41]. This PDF

888 describe the probability that the residual time of the first photon hit is in $[t_r, t_r + \delta]$ with r the radius
 889 of the event vertex, $d = |\vec{r} - \vec{r}_{PMT}|$ the propagation distance, μ_l and μ_d the expected number of PE
 890 and dark noise in the electronic reading window and k is the detected number of PE.

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

$$F(t, r, d) = \int_t^L f(t', r, d) dt' \quad (2.26)$$

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

$$P_T^l(t|k = n) = I_n^l [f_l(t) F_l^{n-1}(t)] \quad (2.27)$$

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

$$f_d(t) = 1/L, F_d(t) = 1 - \frac{t}{L} \quad (2.28)$$

897 and so for the case where only one photon is detected by the PMT ($k = 1$)

$$P_T(t|\mu_l, \mu_d, k = 1) = I_1[P(1, \mu_l)P(0, \mu_d)f_l(t) + P(0, \mu_l)P(1, \mu_d)f_d(t)] \quad (2.29)$$

898 where $P(k_\alpha, \mu_\alpha)$ is the Poisson probability to detect k_α PE from $\alpha \in \{l, d\}$ with the condition $k_l + k_d =$
 899 k .

900 Now that we have the individual timing and charge probability we can construct the charge likeli-
 901 hood referred as QMLE:

$$\mathcal{L}(q_1, q_2, \dots, q_N | \vec{r}, E_{vis}) = \prod_{j \in \text{unfired}} e^{-\mu_j} \prod_{i \in \text{fired}} \left(\sum_{k=1}^K P_Q(q_i|k) \cdot P(k, \mu_i) \right) \quad (2.30)$$

902 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
 903 charge pdf for k PE. And we can also construct the time likelihood referred as TMLE:

$$\mathcal{L}(t_{1,r}, t_{2,r}, \dots, t_{N,r} | \vec{r}, t_0) = \prod_{i \in \text{hit}} \frac{\sum_{k=1}^K P_T(t_{i,r}|r, d, \mu_i^l, \mu_i^d, k) \cdot P(k, \mu_i^l + \mu_i^d)}{\sum_{k=1}^K P(k, \mu_i^l + \mu_i^d)} \quad (2.31)$$

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

905 Merging those two likelihood give the charge-time likelihood QTMLE

$$\mathcal{L}(q_1, q_2, \dots, q_N; t_{1,r}, t_{2,r}, \dots, t_{N,r} | \vec{r}, t_0, E_{vis}) = \mathcal{L}(q_1, q_2, \dots, q_N | \vec{r}, E_{vis}) \cdot \mathcal{L}(t_{1,r}, t_{2,r}, \dots, t_{N,r} | \vec{r}, t_0) \quad (2.32)$$

906 The radial and energy resolutions of the different likelihood are presented in figure 2.23 (from [41]).
 907 We can see the improvement of adding the time information to the vertex reconstruction and that
 908 an increase in vertex precision can bring improvement in the energy resolution, especially at low
 909 energies.

910 Data driven methods prove to be performant in the energy and vertex reconstruction given that we
 911 have enough calibrations sources to produce the PDF. In the next section, we'll see another type of
 912 data-driven method based on machine learning.

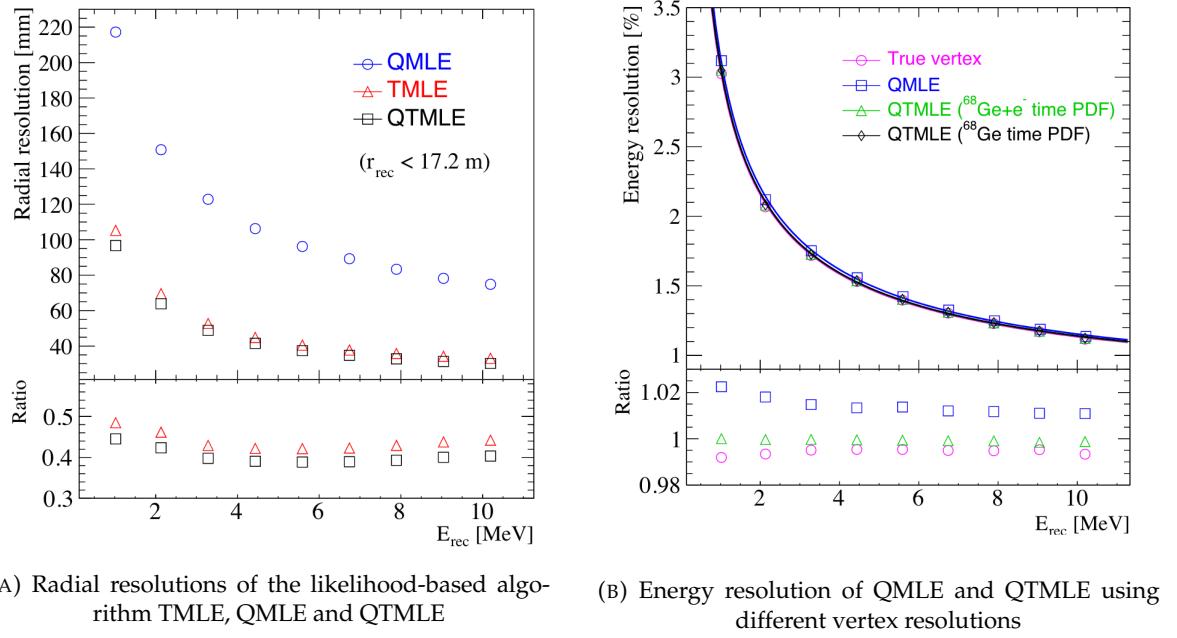


FIGURE 2.23

2.8.3 Machine learning for reconstruction

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

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

$$\mathbb{R}^{5 \times N_{\text{pmt}}} \mapsto \mathbb{R}^4 \quad (2.33)$$

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

For machine learning to converge to performant algorithm, a large dataset exploring all the phase space of interest is needed. For the following studies, data from the monte carlo simulation presented in section 2.6 are used for training. When the detector will be finished calibrations sources will be complementarily be used.

932 **Boosted Decision Tree (BDT)**

933 On of the most classic ML method used in physics in last years is the Boosted Decision Tree (see
 934 Chapter 3.1.1). They have been explored for vertex reconstruction [42] et for energy reconstruction
 935 [42, 43].

936 For vertex and energy reconstruction a BDT was developed using the aggregated informations pre-
 937 sented in 2.8.

Parameter	description
$nHits$	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 2.8 – Features used by the BDT for vertex reconstruction

938 Its reconstruction performances are presented in figure 2.25.

939 A second and more advanced BDT, subsequently named BDTE, that only reconstruct energy use a
 940 different set of features [43]. They are presented in the table 2.9

941 **Neural Network (NN)**

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

946 The CNN are using 2D projection of the detector representing it as an image with two channel, one
 947 for the charge Q and one for the time t . The position of the PMTs is structurally encoded in the pixel
 948 containing the information of this PMT. In [42], the pixel is chosen based on a transformation of θ
 949 and ϕ coordinates to the 2D plane and rounded to the nearest pixel. A sufficiently large image has
 950 been chosen to prevent two PMT to be located in the same pixel. An example of this projection can
 951 be found in figure 2.24. The performances of the CNN can be found in figure 2.25.

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

956 One of the way to present structurally the sphericity of JUNO to a NN is to use a graph: A collection
 957 of objects V called nodes and relations E called edges, each relation associated to a couple v_1, v_2

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

TABLE 2.9 – 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

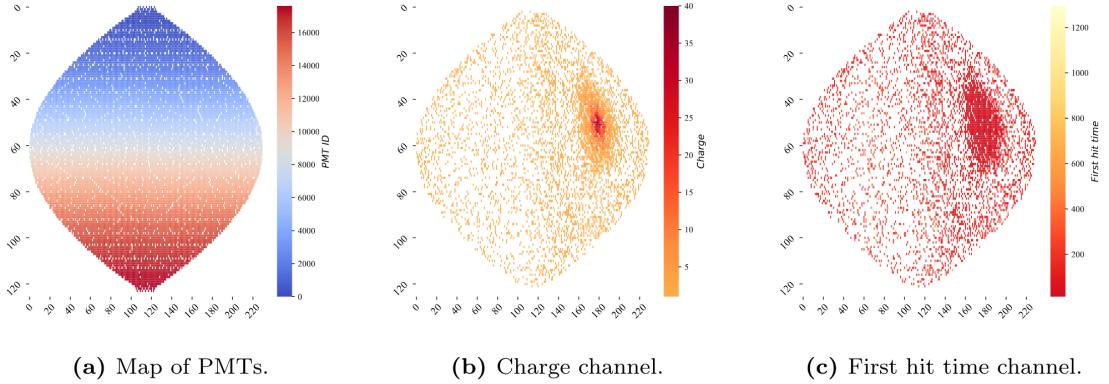


FIGURE 2.24 – Projection of the LPMTs in JUNO on a 2D plane. (a) Show the distribution of all PMTs and (b) and (c) are example of what the charge and time channel looks like respectively

958 forming the graph $G(E, V)$. Nodes and edges can hold informations or features. In [42] the nodes,
959 are geometrical region of the detector as defined by the HealPix [48]. The features of the nodes are
960 aggregated informations from the PMTs it contains. The edges contains geographic informations of
961 the nodes relative positions.

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

965 The neural network then process the graph using Chebyshev Convolutions [49]. The performances
966 of the GNN are presented in figure 2.25.

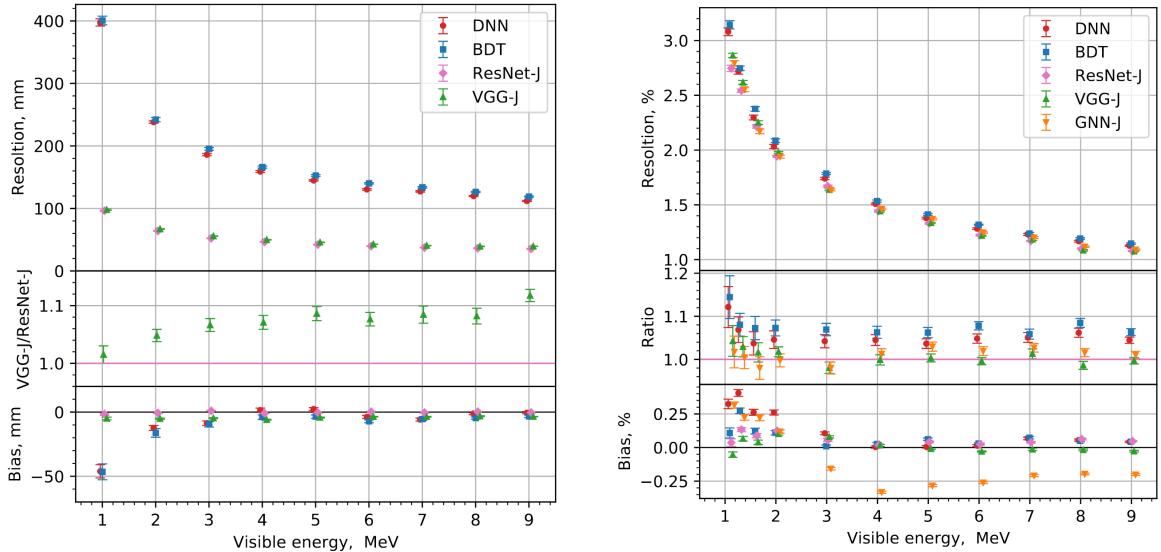


FIGURE 2.25 – Radial (left) and energy (right) resolutions of different ML algorithms.
The results presented here are from [42]. DNN is a deep neural network, BDT is a BDT,
ResNet-J and VGG-J are CNN and GNN-J is a GNN.

967 Overall ML algorithms show similar performances as classical algorithms in term of energy recon-
968 structions with the more complex structure CNN and GNN showing better performances than BDT

and DNN. For vertex reconstruction, the BDT and DNN show poor performance while CNN are on the level of the classical algorithms.

2.8.4 Physics results

The oscillation parameters are directly extracted from the minimization procedure and the error can be estimated directly from the procedure. For the NMO, the data are fitted under the two assumption of NO and IO. The difference in χ^2 give us the preferred ordering and the significance of our test. 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 Δm_{31}^2 , Δm_{21}^2 , $\sin^2 \theta_{12}$ and $\sin^2 \theta_{13}$ respectively [3]. The expected sensitivity to mass ordering is 3σ after 6.5 years [50].

2.9 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.

988 **Chapter 3**

989 **Machine learning: Introduction to the
990 methods and algorithms used in this
991 thesis**

992 “I have the shape of a human being and organs equivalent to those of a
993 human being. My organs, in fact, are identical to some of those in a
994 prostheticized human being. I have contributed artistically, literally, and
995 scientifically to human culture as much as any human being now
996 alive. What more can one ask?”

997 Isaac Asimov, *The Complete Robot*

998 **Contents**

<small>999</small> 3.1 Core concepts in machine learning and neural networks	<small>1000</small> 46
3.1.1 Boosted Decision Tree (BDT)	46
3.1.2 Artificial Neural Network (NN)	46
3.1.3 Training procedure	48
3.1.4 Potential pitfalls	51
<small>1001</small> 3.2 Neural networks architectures	<small>1002</small> 54
3.2.1 Fully Connected Deep Neural Network (FCDNN)	54
3.2.2 Convolutional Neural Network (CNN)	54
3.2.3 Graph Neural Network (GNN)	56
3.2.4 Adversarial Neural Network (ANN)	58

1003 Machine Learning (ML) and more specifically Neural Network (NN) are families of data-driven
1004 algorithms. They are used in a wide variety of domains including natural language processing,
1005 computer vision, speech recognition and, the subject of this thesis, scientific studies.

1006 They are used to model complex distributions from a finite dataset to extract a generalist behavior.
1007 For example, in our case, it could be an algorithm that would differentiate the nature of a particle
1008 interacting in the liquid scintillator, between a positron and an electron, based on the readout charge
1009 and time (Q, t) of the 17612 LPMT of the JUNO experiment. During a first training phase, it would
1010 learn the discriminative features between the two in the 35224-dimensional charge and time distri-
1011 bution, built from samples of e^+ and e^- events.

1012 It would learn to derive from a complex, highly dimensional set of data the essential few informations
1013 characterizing the interactions: a three body energy deposition (the positron and two annihilation
1014 gammas) and the single deposit from an electron.

1015 Ideally, the algorithm would learn to recognize those informations on its own, regardless of the input
1016 size and complexity. In practice, however, these algorithms are guided by human design through

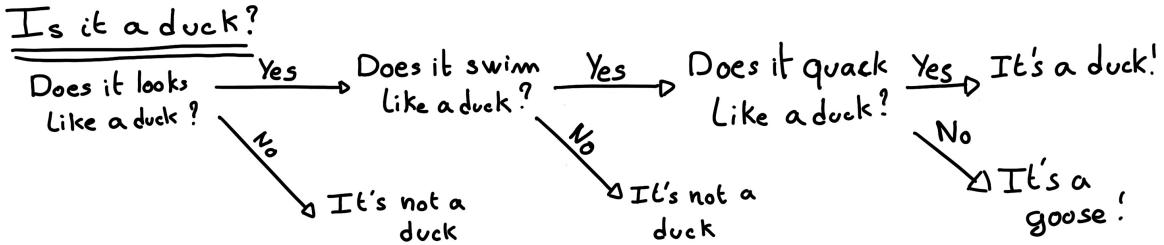


FIGURE 3.1 – Example of a BDT that determine if the given object is a duck

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 2.8).

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.

3.1.1 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.

3.1.2 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

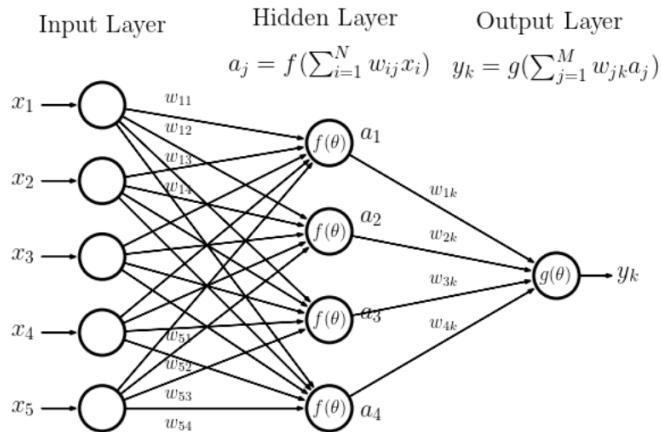


FIGURE 3.2 – Schema of a simple neural network

1052 set of values from the preceding layer, here the a_i takes every informations of the x_i input layer,
 1053 and aggregate those values following learnable *parameters* w_{ij} . The aggregation procedure is core of
 1054 defining the architecture of the NN. The different architectures used in this thesis will be discussed
 1055 in section 3.2. The process is repeated until reaching the output layer.

1056 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
 1057 output layer. We try to produce a vertex reconstruction algorithm that will approach the charge
 1058 barycentre. Let's limit the input x_i to the charge of the i th PMT, one of the solution is to aggregate
 1059 on a_1 the x coordinate of the barycenter. The network would thus adapt the w_{i1} parameters so
 1060 they correspond to the x coordinates of the i th PMT. Same for the y and z coordinate on a_2 and
 1061 a_3 respectively.

1062 The layers used in the example above are designated as *Fully connected* layers, where every neurons
 1063 of the layer is connected to the every neurons of the preceding layer. The layer can be expressed
 1064 using the Einstein summation and in bold the learnable parameters

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

1065 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}
 1066 is the parameters, or weights, matrix (composed of the w_{ij}). In practice, this fully connected layer is
 1067 often adjoined a bias \mathbf{B} and an *activation function* F .

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

1068 This is the fundamental component of the Fully Connected Deep NN (FCDNN) family presented in
 1069 section 3.2.1.

1070 This description of neural networks as layers introduce the principles of *depth* and *width*, the number
 1071 of layers in the NN and the number of neurons in each layer respectively. Those quantities that not
 1072 directly used for the computation of the results but describes the NN or its training are designated
 1073 as *hyperparameters*.

1074 Now we just need to adapt the parameters so that this network learn that w_{ij} are the PMT coordinate.
 1075 We describe the space produced by the parameters of the network as the *parameter phase space* or *latent*
 1076 *space*. The optimization of the network and exploration of this phase space is done through training
 1077 over a *training dataset* as described in next section.

1078

3.1.3 Training procedure

1079 To adapt the parameters we need an object that describe how well the network perform. This is
 1080 the *loss* of our neural networks \mathcal{L} . In our barycenter example, it could be the distance between the
 1081 reconstructed and real barycenter. Using this metric we can adjust the parameters of our network.

1082 Depending if we try to minimize or maximize it, it need to posses a minima or a maxima. For example
 1083 when doing *regression*, i.e. produce a scalar result like the coordinates of a barycenter, a common loss
 1084 is the Mean Square Error (MSE). Let i be our dataset, the N events considered for training, y_i be the
 1085 target scalar, the barycenter positions of each events, x_i the input data, the charge vector, and $f(x_i, \theta)$
 1086 the result of the network. The network here is modelled by f , and its parameter θ

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

1087 Another common loss function is the Mean Absolute Error (MAE)

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

1088 We see that those loss function possess a minima when $f(x_i, \theta) = y_i$.

1089 Most of the modern neural networks use gradient descent to optimize their parameters, i.e. the
 1090 gradient of the parameter w , designated in literature as θ , with respect of the loss function \mathcal{L} is
 1091 subtracted each optimisation step t

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

1092 This induce \mathcal{L} needs to be differentiable with respect to θ , thus the layers and their activation func-
 1093 tions also need to be differentiable. This simple gradient descent, designated as Stochastic Gradient
 1094 Descent (SGD), can be extended with first and second order momentums like in the Adam optimizer
 1095 [53]. More details about the optimizers can be found in section 3.1.3.

1096

Training lifecycle

1097 The training of NN does not follow strict rules, you could imagine totally different lifecycle but I will
 1098 describe here the one used in this thesis, the most common one.

1099 As illustrated in figure 3.3, the training is split into *epochs*. Each epochs is split into *step* where the
 1100 NN will optimize its parameters over a *batch*, a sub-sample of the training datasets. The ideal batch
 1101 size, number of event in a batch, would be the entire dataset, as the NN optimization would not be
 1102 biased by the specificity of a sub-sample, but due to memory limitations the batch size is driven by
 1103 technical limitations.

1104 At the end of each epochs, the neural network is evaluated over a validation dataset, a dataset from
 1105 which no optimisation is done. It is used as reference for the network performance as and monitor
 1106 overtraining (see section 3.1.4).

1107 Hyperparameters that can be optimized during the training can be optimized at each epoch, for
 1108 example the learning rate, or each step, the optimizer momentum for example.

1109 There is not really a typical number of epochs or steps for the training. The number steps can be
 1110 defined such as in one epoch, the NN see the entirety of the dataset but the number of steps and
 1111 epochs are hyperparameters that are optimized over the each subsequent training. We adjust them
 1112 by looking at the loss evolution profile over time.

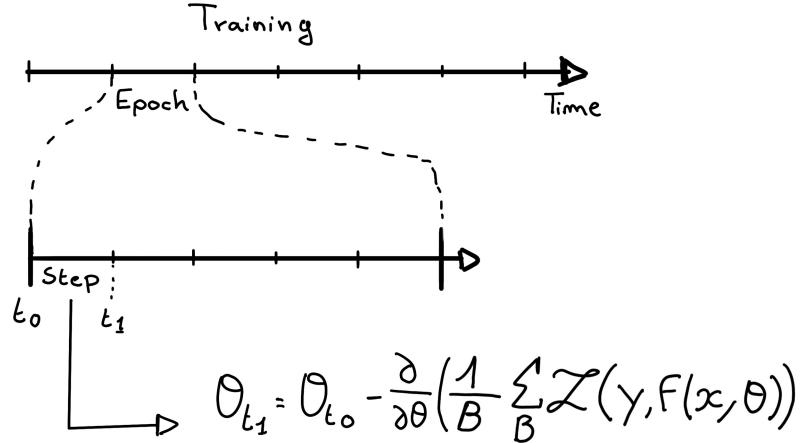


FIGURE 3.3 – Illustration of the training lifecycle

1113 Most training are started with a fixed number of epochs, i.e. from what we've seen from precedent
 1114 training, the network stop learning, the loss is constant, after N epoch so we run the training for
 1115 $N + \delta$ epochs to see if the modification brings improvements to the loss profile. We can setup what's
 1116 called *early stopping policies* that'll stop the training early in specific cases like loss explosion or loss
 1117 stability but this require fine tuning and don't bring much in our case as we are not really limited in
 1118 training time.

1119 The optimizer

1120 As briefly introduced at the beginning of this section, the parameters of the neural network are
 1121 optimized using the gradient descent method. We compute the gradient of the mean loss over the
 1122 batch with respect of each parameters and we update the parameters in accord to minimize the loss.
 1123 The gradient is computed backward from the loss up to the first layer parameters using the chain
 1124 rule, in this case with only one parameter at each step for simplicity:

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

1125 where θ is a parameter, i is the layer index. We see here that the gradient of the first layer is
 1126 dependent of the gradient of all the following layers. Because the only value known at the start
 1127 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
 1128 propagation*.

1129 This update of the parameters is done following an optimizer policy. Those optimizers depends on
 1130 hyperparameters. The ones used in this thesis are:

- 1131 1. SGD (Stochastic Gradient Descent). This is the simplest optimizer, it depend on only one
 1132 hyperparameter, the learning rate λ (LR) and update the parameters θ following

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

1133 where t is the step index. It is a powerful optimizer but is very sensible to local minima of the
 1134 loss in the parameters phase space as illustrated in figure 3.4a.

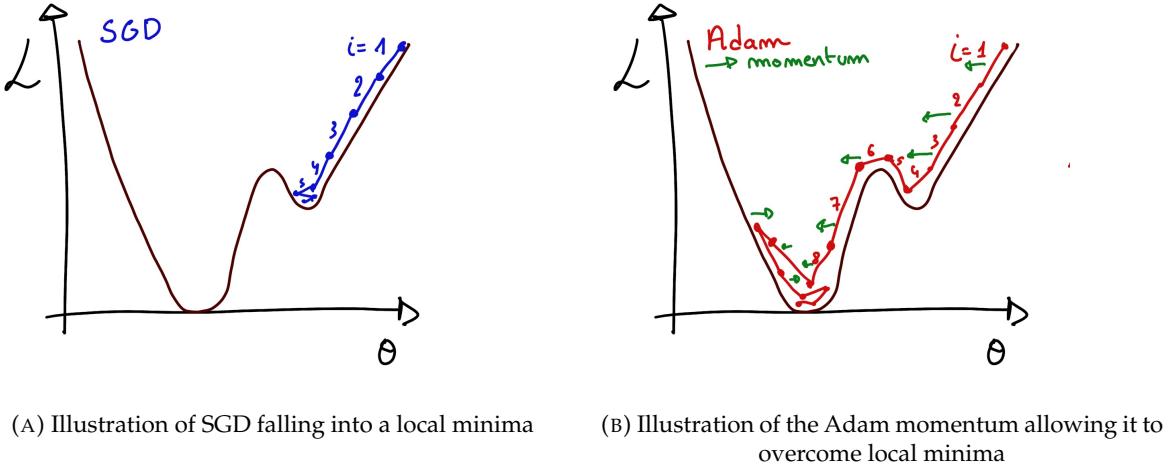


FIGURE 3.4

2. Adam [53]. The concept is, in short, to have and SGD but with momentum. Adam possess two momentum $m(\beta_1)$ and $v(\beta_2)$ which are respectively proportional to $\frac{\partial \mathcal{L}}{\partial \theta}$ and $(\frac{\partial \mathcal{L}}{\partial \theta})^2$. β_1 and β_2 are hyperparameters that dictate the moment update at each optimization step. The parameters are then upgraded following

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

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

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

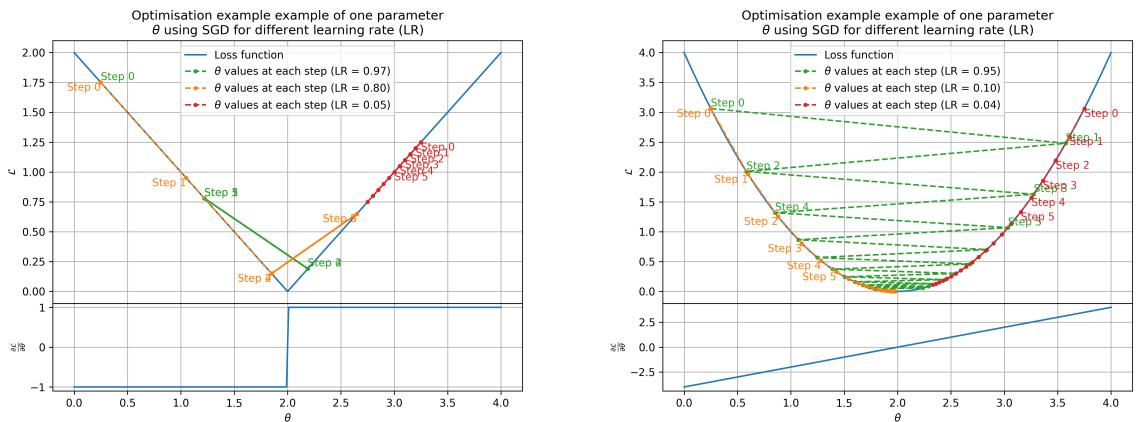
where ϵ is a small number to prevent divergence when v is close to 0. These momentums allow to overcome small local minima in the parameters phase. Imagine ball going down a slope as illustrated in 3.4a, if you ignore the stored momentum you get SGD and get stuck as on the left plot. Now if you consider 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 rat, 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 use is the save of the best model. In some situation, the loss value after each epoch will strongly oscillate or can even worsen. This policy allow us to keep the best version



(A) Illustration of the SGD optimizer on one parameter θ 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 θ 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.

of the model attained during the training phase.

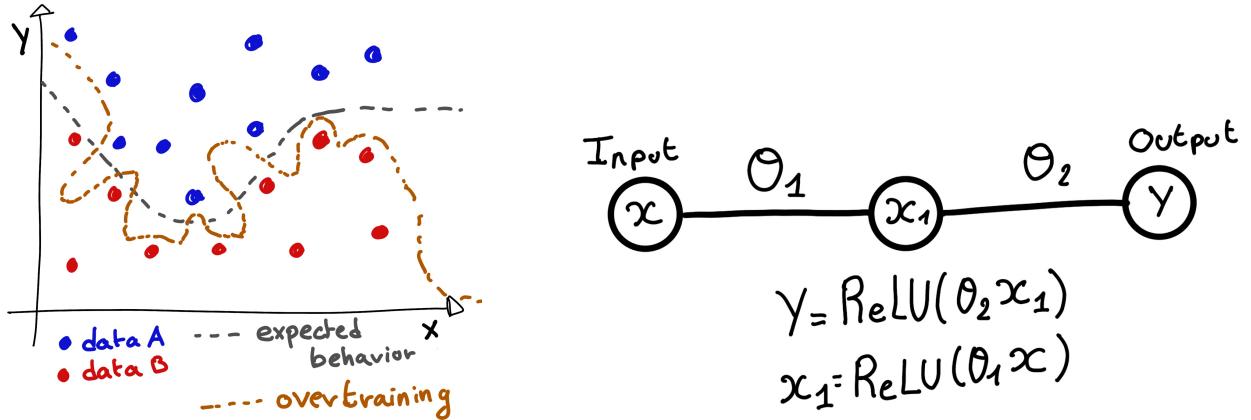
3.1.4 Potential pitfalls

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

Overtraining

This happen when the network learn 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 posses 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 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 parameter 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 diverge 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

1173 Gradient vanishing

1174 Gradient vanishing is the effect of the gradient being so small for the early layers that the parameters
 1175 are barely updated after each step. This cause the network to be unable to converge to the minima.

1176 This comes from the way the gradient descent is calculated. Imagine a simple network composed of
 1177 three fully connected layers: the input layer, a intermediate layer and the output layer. Let L be the
 1178 loss, θ_1 the parameter between the input and the intermediate layer and θ_2 the parameter between
 1179 the intermediate and output layer. This network is schematized in figure 3.6b.

1180 The gradient for θ_1 will be computed using the chain rule presented in equation 3.6. Because θ_1
 1181 depends on θ_2 , if the gradient of θ_2 is small, so will be the gradient of θ_1 . Now if we would have
 1182 much more layer, we can see how the subsequent multiplication of small gradients would lead to
 1183 very small update of the parameters thus "vanishing gradient".

1184 Multiple actions can be taken to prevent this effect such as:

- 1185 — **Batch normalization:** In this case we apply a normalization layer that will normalize the data.
 1186 It means that we transform the input variable X into a variable D which distribution follow
 1187 $\langle D \rangle = 0$ and $\sigma_D = 1$. This helps the parameters of the network to maintain an appropriate
 1188 scale.
- 1189 — **Residual Network (ResNet)** [54]: Residual network is a technique for neural network in
 1190 which, instead of just sequentially feeding the results of each layer to the next one, you
 1191 compute a residual over the input data. This technique is illustrated in figure 3.7. The
 1192 reference [54] show empirical evidence of its relevance.

1193 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 θ

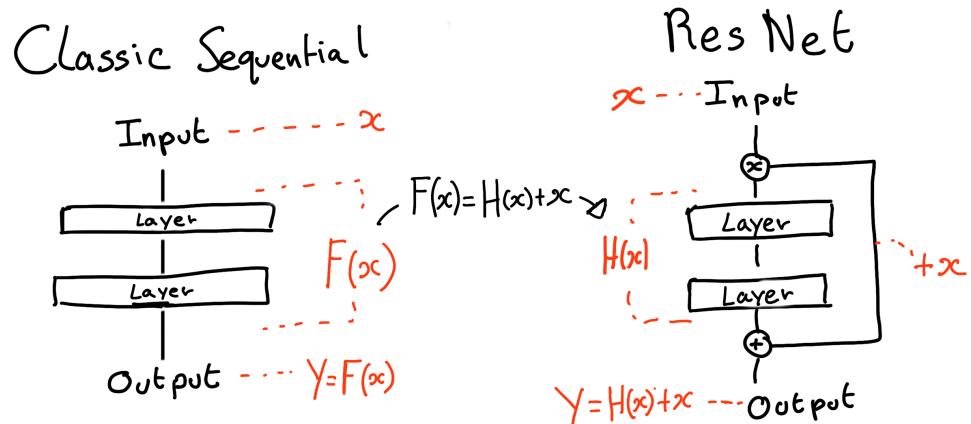


FIGURE 3.7 – Illustration of the ResNet framework

follow

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

1194 The explosion is illustrated in figure 3.8 where we can see that the loss degrades with each step of
 1195 optimization. In this illustration it is clear that reducing the learning rate suffice but this behaviour
 1196 can happens in the middle of the training where the learning rate schedule does not permit reactivity.

1197 There exist solutions to prevent this explosions:

- 1198 — **Gradient clipping:** Is this case we work on the gradient so that the norm of gradient vector
 1199 does not exceed a certain threshold. In our illustration in figure 3.8 the gradient for $\theta > 0$
 1200 could be clipped at 3 for example.
- 1201 — **Batch normalization:** For the same reasons as for gradient vanishing, normalizing the input
 1202 data help reduce erratic behaviour.

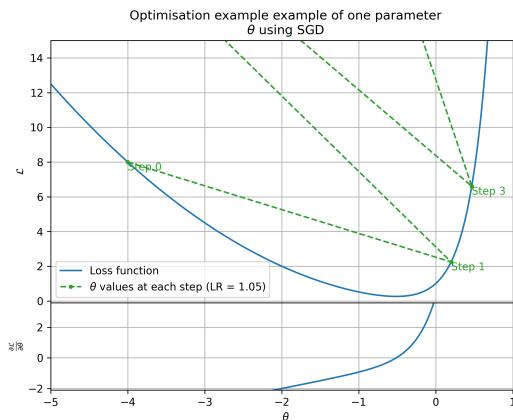


FIGURE 3.8 – Illustration of the gradient explosion. Here it can be solved with a lower learning rate but its not always the case.

1203 **3.2 Neural networks architectures**

1204 **3.2.1 Fully Connected Deep Neural Network (FCDNN)**

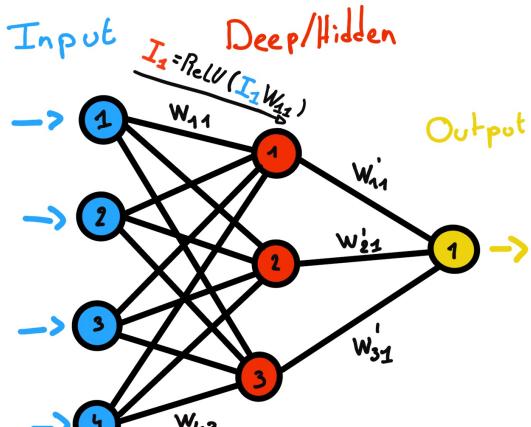
1205 The Fully Connected Deep Neural Network (FCDNN) architecture is the stack of multiple fully
 1206 connected layers as presented in the figure 3.9a. Most of the time, the classic ReLU function

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

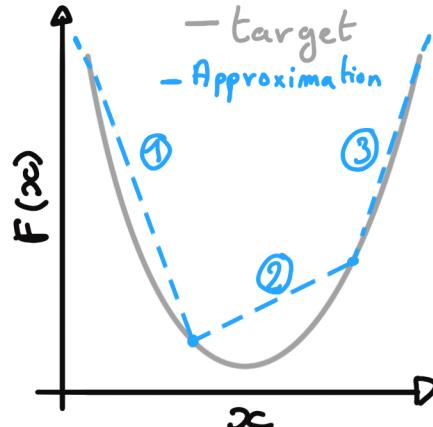
1207 is used as activation function. PReLU and Sigmoid are also popular choices:

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

1209 The reasoning behind ReLU and PReLU is that with enough of them, you can mimic any continuous
 1210 function as illustrated in figure 3.9b. Sigmoid is more used in case of classification, its behavior going
 1211 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

1212 Due to its simplicity, FCDNN are also used as basic pieces for more complex architectures such as
 1213 the CNN and GNN that will be presented in the next sections.

1214 **3.2.2 Convolutional Neural Network (CNN)**

1215 It's not trivial to describe in text the principles of Convolutional Neural Network (CNN) and how
 1216 they works. We try a general description below followed by a step by step description of a concrete
 1217 example.

1218 Convolutional Neural Networks are a family of neural networks that use discrete convolution filters,
 1219 as illustrated in an example in figure 3.10, to process the input data, often images. They are com-
 1220 monly used in image recognition [55] for classification or regression problematics. Concretely, you
 1221 multiply element-wise a portion of the input data, in the case of an image, a small part of the image,

1222 with a kernel of same dimension. In figure 3.10, we multiply the 3×3 pixels sub-image with the
 1223 3×3 kernel.

1224 Their filters scan the input data, highlighting patterns of interest, this scanning procedure making
 1225 them translation-invariant. In the concrete case of figure 3.10, for each pixel of the input image, we
 1226 group it with the 8 neighbours pixel and produce a new pixel that correspond to the output image.
 1227 For the pixel on the edges that do not have neighbours, we either create “imaginary” pixel with the
 1228 value 0 or we just ignore them. If we ignore them, the output image will posses fewer pixels than the
 1229 input image. We see that the operation do not care where is the pattern of interest in the images, the
 1230 filter output will be *invariant* whatever *translation* is applied to the image.

1231 This invariance mean that they are capable of detecting oriented features independently of their
 1232 location on the image. Again taking 3.10 as an example, with only the 9 parameters composing the
 1233 kernel, we can highlight the contour of the duck by looking at the “yellowness” of the pixels.

1234 The learning parameters of CNNs are the kernels components, the network thus learn the optimal
 1235 filters to extract the desired features.

1236 The convolution layers are commonly chained [56], reducing the input dimension while increasing
 1237 the number of filters. The idea behind is that the first layers will process local informations and
 1238 the latest layers will process more global informations, as the latest convolution filters will process
 1239 the results of the preceding that themself have processed local information. To try to preserve the
 1240 amount of information, we tend to grow the numbers of filters for each division of the input data.
 1241 The results of the convolution filters is commonly then flattened and feed to a smaller FCDNN which
 1242 will process the filters results to yield the desired output.

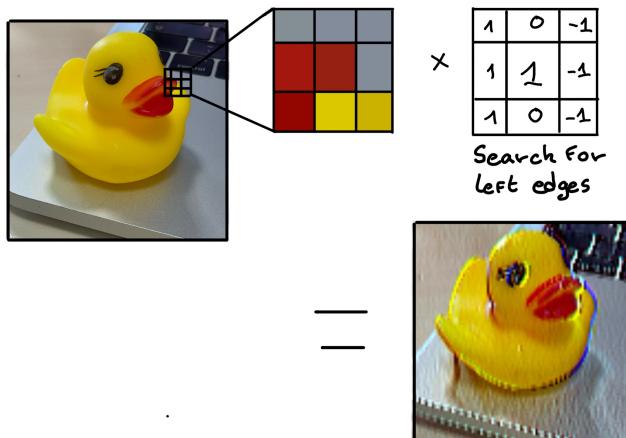


FIGURE 3.10 – Illustration of the effect of a convolution filter. Here we apply a filter with the aim do detect left edges. We see in the resulting image that the left edges of the duck are bright yellow where the right edges are dark blue indicating the contour of the object. The convolution was calculated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

1243 As an example, let’s take the Pytorch [58] example for the MNIST [59], a dataset of black and white
 1244 images of handwritten digits. Those images are 28×28 pixels with only one channel corresponding
 1245 to the grey level of the pixel. Example of images from this dataset are presented in figure 3.11a

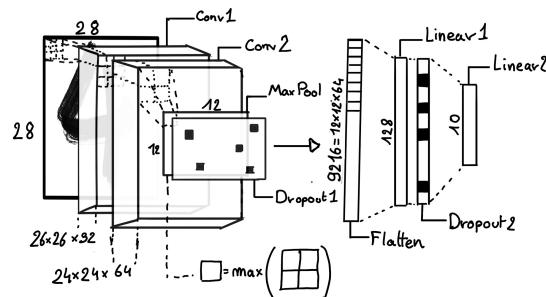
1246 A schema of the CNN used in the Pytorch example is presented in figure 3.11b. Using this schema
 1247 as a reference, the trained network is made of:

- 1248 1. A convolutional layer of (3×3) filters yielding 32 channels. A bias parameter is applied
 1249 to each channel for a total of $(32 \cdot (3 \times 3) + 32) = 320$ parameters. The resulting image is
 1250 $(26 \times 26 \times 32)$ (26 per 26 pixels with 32 channels). The ReLU activation function is applied to
 1251 each pixel.

- 1252 2. A second convolutional layer of (3×3) filters yielding 64 channels. This channel also posses
 1253 a bias parameter for a total of $(64 \cdot (3 \times 3) + 64) = 640$ parameters. Resulting image is $(24 \times$
 1254 $24 \times 64)$. This channel also apply a ReLU activation function.
- 1255 3. Then comes a (2×2) max pool layer with a stride of 1 meaning that for each channel the max
 1256 value of pixels in a (2×2) block is condensed in a single resulting pixel. The resulting image
 1257 is $(12 \times 12 \times 64)$.
- 1258 4. This image goes through a dropout layer which will set the pixel to 0 with a probability of
 1259 0.25. This help prevent overtraining the neural network (see section 3.1.4 for more details).
- 1260 5. The data is the flattened i.e. condensed into a vector of $(12 \times 12 \times 64) = 9216$ values.
- 1261 6. Then comes a fully connected linear layer (Eq. 3.2) with a ReLU activation that output 128
 1262 feature. It needs $(9216 \cdot 128) + 128 = 1'179'776$ parameters.
- 1263 7. This 128 item vector goes through another dropout layer with a probability of 0.5
- 1264 8. The vector is then transformed through a linear layer with ReLU activation. It output 10
 1265 values, one for each digit class $(0, 1, 2, \dots, 9)$. It need $(128 \cdot 10) + 128 = 1408$ parameters.
- 1266 9. Finally the 10 values are normalized using a log softmax function $\text{LogSoftmax}(x_i) = \log \left(\frac{\exp(x_i)}{\sum_j \exp(x_j)} \right)$.
- 1267 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

1268 The final network needs 1'182'144 parameters or, if we consider each parameters to be a double
 1269 precision floating point, 9.45 MB of data. To gives a order of magnitude, such neural network is
 1270 considered "simple", train in a matter of minutes on T4 GPU [60] (14 epochs) and reach an accuracy
 1271 in its prediction of 99%.

3.2.3 Graph Neural Network (GNN)

1272 As seen in the previous section, the CNNs are powerful for image processing, and more generally
 1273 any data that can be expressed as a regular, discrete space and from which the information reside
 1274 in the dispersion in this space. For an image, the edges of an object and how they assemble. A red
 1275 square, straight edges with a sharp angle between them, is much less representative of a duck than
 1276 an yellow sphere, round edges without sharp angles.

1277 This "image" projection is not fitted for every problematics. The signals produced by a detector does
 1278 not always have the properties of images. In the case of JUNO for example, we can create an image
 1279 of two channels, one for the charge Q and one for the timing t but this image should be spheric.
 1280 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

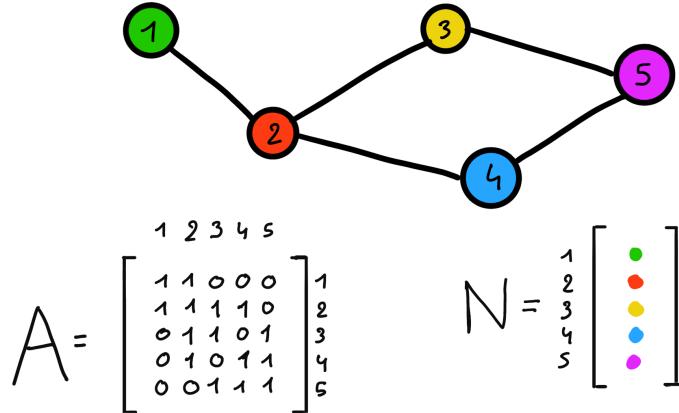


FIGURE 3.12 – Illustration of a graph and its tensor representation.

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.

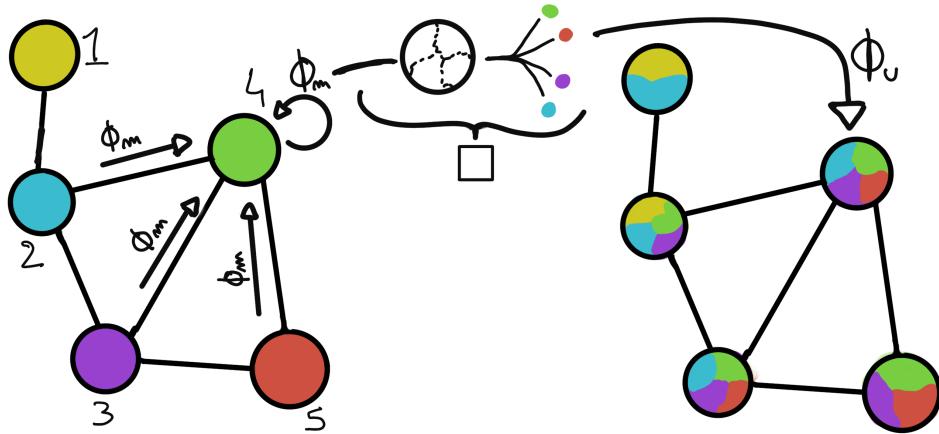


FIGURE 3.13 – Illustration of the message passing algorithm. The detailed explanation can be found in section 3.2.3

1298 To process such object we need specific machine learning algorithms we call Graph neural network.
 1299 To efficiently manipulate graph we need to structurally encode their property in the neural network
 1300 computing architecture: each node is equivalent (as opposite to ordered data in a vector), each node
 1301 has a set of neighbours, ... One of this method is the message passing algorithm presented historically
 1302 in "Neural Message Passing for Quantum Chemistry" [61]. In this algorithm, with each layer of
 1303 message passing a new set of features is computed for each node following

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

1304 where ϕ_u is a differentiable *update* function, \square_j is a differentiable *aggregation* function and ϕ_m is a
 1305 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
 1306 nodes n_j from which it exist an edge $e_{ij} \rightarrow (n_i, n_j)$. k is the layer on which the message passing
 1307 algorithm is applied. The update function need also a few other property if we want to keep the
 1308 graph property, most notably the permutational invariance of its parameters (example: mean, std,
 1309 sum, ...). The differents message, update and aggregation functions can really be any kind of function
 1310 if they follow the constraint presented before, even small Neural Network.

1311 The edges features can also be updated, either by directly taking the results of ϕ_m or by using another
 1312 message function ϕ_e .

1313 To explain this process, let's take the situation presented in figure 3.13. We start with an input graph
 1314 on left, in this case the message passing algorithm is mixing the color on each nodes and produce
 1315 nodes of mixed color. For simplicity, the ϕ_m and ϕ_u function are the identity, they take a color and
 1316 output the same color.

1317 Let's look at what's happening in the node 4. It has 3 neighbours and is a neighbour of itself. The four
 1318 resulting ϕ_m extract the color of each nodes and then feed them to the \square function. The \square function
 1319 just equally distribute the color in the node. Finally the ϕ_u function just update the node with the
 1320 output of \square .

1321 Interestingly we see that the new node 4 does not have any yellow, the color of node 1. But if we were
 1322 to run the message passing algorithm again, it would get some as node 2 is now partially yellow. If
 1323 color here represent information, we see that multiple step are needed so that each node is "aware"
 1324 of the informations the other nodes possess.

1325 Message passing is a very generic way of describing the process of GNN and it can be specialized
 1326 for convolutional filtering [49], diffusion [62] and many other specific operation. GNN are used in a
 1327 wide variety of application such as regression problematics, node classification, edge classification,
 1328 node and edge prediction, ...

1329 It is a very versatile but complex tool.

1330 3.2.4 Adversarial Neural Network (ANN)

1331 The adversarial machine learning, Adversarial Neural Networks (ANN) in the case of neural net-
 1332 work, is a family of unsupervised machine learning algorithms where the learning algorithm (gen-
 1333 erator) is competing against another algorithm (discriminator). Taking the example of Generative
 1334 Adversarial Networks, concept initially developed by Goodfellow et al. [63], the discriminator goal
 1335 is to discriminate between data coming from a reference dataset and data produced by the generator.
 1336 The generator goal, on the other hand, is to produce data that the discriminator would not be able to
 1337 differentiate from data from the reference dataset. The expression of duality between the two models
 1338 is represented in the loss where, at least a part of it, is driven by the results of the discriminator.

1339 **Chapter 4**

1340 **Image recognition for IBD
reconstruction with the SPMT system**

1341 *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 - $\sigma_x \sigma_p \geq \frac{\hbar}{2}$

1342 **Contents**

<small>1344</small>	4.1 Method and model	<small>60</small>
<small>1345</small>	4.1.1 Model	<small>61</small>
<small>1346</small>	4.1.2 Data representation	<small>62</small>
<small>1347</small>	4.1.3 Dataset	<small>64</small>
<small>1348</small>	4.1.4 Data characteristics	<small>65</small>
<small>1349</small>		
<small>1350</small>	4.2 Training	<small>67</small>
<small>1351</small>	4.3 Results	<small>67</small>
<small>1352</small>	4.3.1 J21 results	<small>68</small>
<small>1353</small>	4.3.2 J21 Combination of classic and ML estimator	<small>70</small>
<small>1354</small>	4.3.3 J23 results	<small>72</small>
<small>1355</small>	4.4 Conclusion and prospect	<small>74</small>

1356 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 2.3.2), a thousands versus a few dozen, and their front-end electronics (see section 2.3.2).

1357 The SPMT system is essential to the deployment of the Dual Calorimetry techniques, already men-
1358 tioned in Section 2.8 and described in [24, 26, 64]. It is indeed less subject than the LPMTs to
1359 charge non linearity effects (QNL). This topic will be studied in more detail in Chapter 7, where the
1360 potential of one of the Dual Calorimetry techniques is explored. It consists on combined oscillation
1361 analyses based on two antineutrino energy spectra : one reconstructed with the LPMT system, the
1362 other one with the SPMT system. For that purpose, it is therefore necessary to have reconstruction
1363 tools available. Well maintained tools using the LPMT are available in the collaboration's official
1364 software. This is not the case concerning the SPMT system, where algorithms were developed more
1365 sporadically. This is one of the reasons why we developed the CNN described in this chapter.

1373 Our efforts on it were limited to the early months of this thesis: it was above all a way to learn about
 1374 ML and about JUNO's detector and software. We benchmarked its performance against a classical
 1375 algorithm developed in [65] but not yet implemented in JUNO's software.

1376 As discussed in Chapter 3, Machine Learning (ML) algorithms shine when modeling highly dimen-
 1377 sional data from a given dataset. In our case, we have access to complete monte-carlo simulation of
 1378 our detector to produce large datasets that could represent multiple years of data taking. Ideally ML
 1379 algorithms would be able to consider the entirety of the information in the detector and converge on
 1380 the best parameters to yield optimal results.

1381 The difference between this ideal and what can be achieved in reality is an important subject. In
 1382 particular, we wonder if an exhaustive usage of the information present in the detector could lead to
 1383 use informations that are mismodelled in our simulated training samples (or present only in these
 1384 samples) and therefore lead to biases when the algorithm is applied to real data. A simple way
 1385 to start addressing this reliability issue is to try to evaluate to which extent various reconstruction
 1386 methods use the same information. An attempt at this is presented at the end of this chapter. This is
 1387 also the subject of Chapter 6.

1388 4.1 Method and model

1389 One of simplest way to look at JUNO data is to consider the detector as an array of geometrically
 1390 distributed sensors on a sphere. Their repartition is almost homogeneous, on this sphere surface
 1391 providing an almost equal amount of information per unit surface. It is then tempting to represent
 1392 the detector as a spherical image with the PMTs in place of pixels. Two events with two different
 1393 energy or position would produce two different images.

1394 The most common approach in machine learning for image processing and image recognition is the
 1395 Convolutional Neural Network (CNN). It is widely used in research and industry [56, 66–68] due to
 1396 its strengths (see section 3.2.2) and has proven its relevance in image processing.

1397 Some CNN are developed to process spherical images [69] but for the sake of simplicity and as a
 1398 first approach we decided to go with a planar projection of the detector, approach that has proven its
 1399 efficiency using the LPMT system (see section 2.8.3). The details about this planar projection will be
 1400 discussed in section 4.1.2.

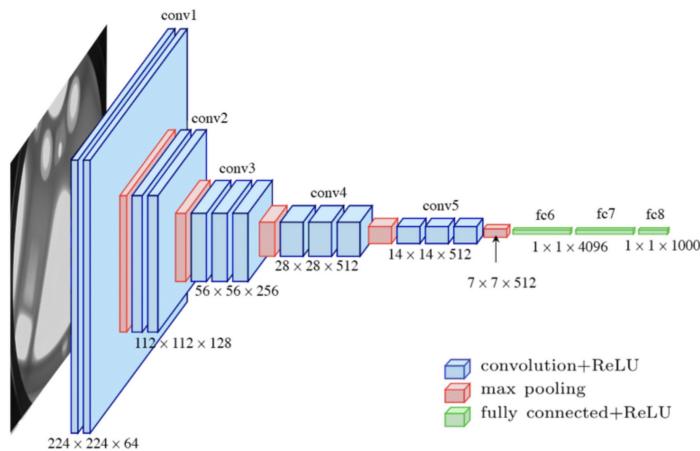


FIGURE 4.1 – Graphic representation of the VGG-16 architecture, presenting the different kind of layer composing the architecture.

4.1.1 Model

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

- N_{blocks} : the number of convolution blocks, a block being composed of two convolutional layers with 3×3 filters using ReLU activation function, a 3×3 kernel max-pooling layer (except for the last block).
- $N_{channels}$: The number of channels in the first block. The number of channels in the subsequent blocks is computed using $N_{channels}^i = i * N_{channels}$, $i \in [1..N_{blocks}]$.
- **FCDNN configuration:** The result of the last convolution layer is flattened then fed to a FCDNN. Its configuration is expressed as the ouputs of sequenced fully connected linear layer using the PReLU activation function. For example $2 * 1024 + 2 * 512$ is the sequence of 2 layers which output is 1024 followed by 2 other layers with an output of 512. Finally the last layer is a linear layer outputting 4 features without activation function. Each feature 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.

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

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

where E_{dep} is the deposited energy and R is the radius of JUNO's CD. With the energy in MeV and the distance in meters, we use the factor 0.85 and 10 to balance 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. The loss $(E_r + V_r)$ 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 2.23. 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.

The above losses and their parameters values results from fine-tuning after multiples runs and adjustments of the full random search.

Each combinations of those hyperparameters (for example ($N_{blocks} = 2, N_{channels} = 32$, FCDNN = $(2 * 1024)$, Loss = $(E + V)$)) produce models, hereinafter referred as configurations, are then tested and compared to each other over an analysis sample.

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

- Gen₀: $N_{blocks} = 4, N_{channels} = 64$, FCDNN configuration: $1024 * 2 + 512 * 2$, Loss $\equiv E + V$
- Gen₁: $N_{blocks} = 4, N_{channels} = 64$, FCDNN configuration: $1024 * 2 + 512 * 2$, Loss $\equiv E_r + V_r$
- Gen₂: $N_{blocks} = 5, N_{channels} = 64$, FCDNN configuration: $4096 * 2 + 1024 * 2$, Loss $\equiv E + V$
- Gen₃: $N_{blocks} = 5, N_{channels} = 64$, FCDNN configuration: $4096 * 2 + 1024 * 2$, Loss $\equiv E_r + V_r$

The resulting models possess between 2'041'034, for Gen₅₂ and Gen₅₃, and 5'759'839'242 parameters, for Gen₂₆ and Gen₂₇. The models of interest in this thesis, from which the results are discussed in section 4.3, possess 86'197'196 parameters for Gen₃₀ and 332'187'530 parameters for Gen₄₂. For comparison the model of CNN developed in JUNO before posses 38'352'403 parameters [42].

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

TABLE 4.1 – Sets of hyperparameters values considered in this study

To rank the various configuration we cannot used directly the mean loss over the validation dataset as ($E + V$) and ($E_r + V_r$) are not numerically comparable. We thus use the following quantities, directly related to the reconstruction performances:

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

The models were developped in Python using the Pytorch framework [58] using NVIDIA A100 [70] and NVIDIA V100 [71] gpus. The A100 was split in two, thus the accessible gpu memory was the same as V100, 20 Gb, making it impossible to train some of the architectures due to memory consumption.

The training was monitored in realtime by a custom tooling that was developed during this thesis, DataMo [72].

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

4.1.2 Data representation

This data is represented as 240×240 images with a charge Q channel and a time t channel. The SPMTs are then projected on the plane as illustrated in figure 4.2b using the coordinate system presented in 4.2a. The P_y coordinate, the row corresponding to the SPMT in the projection, is proportional to θ . The P_x coordinate, the column corresponding to the SPMT in the projection, is defined by $\phi \sin \theta$ in spherical coordinates. $\theta = 0$ is defined as being the top of the detector and $\phi = 0$ is defined as an arbitrary direction in the detector. In practice, $\phi = 0$ is given by the MC simulation.

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

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

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

This projection keep the SPMT position in the image proportional to their spherical coordinates while keeping the neighbouring information. This proportionality allow us to keep the specificities of the detector structure, the vertical bands visible in 4.2b.

When two SPMTs in the same pixel are hit in the event time window, the charges are summed and the lowest of the hit-time is chosen. The time window depends on the datasets and are detailed in section 4.1.2. 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 timing is thus more representative of the event location.

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 4.2b 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 4.1.1, the convolution filter we use are 3×3 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) [73] 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. It's important to keep in mind

1518 the fact that a part of the detector that has not been hit is also an information: There is no signal in
 1519 this part of the detector. This problematic will be explored in more details in Chapter 5.

1520 The second kind of pixels are the one that do not represent parts of the detector such as the corners
 1521 of the image. The question is basically the same, what to put in the charge and the time channel. The
 1522 decision is to set the charge and time to 0 following the above reasoning.

1523 Another problematic that happens with this representation, and this is not dependent of the chosen
 1524 projection, is the deformation in the edges of the image and the loss of the neighbouring information
 1525 in the for the SPMTs at the edge of the image $\phi \sim 180^\circ$. This deformation and neighbouring loss
 1526 could be partially circumvented as explained in section 4.4

1527 4.1.3 Dataset

1528 In this study we will discuss two datasets of one millions prompt signal of IBD events.

1529 J21

1530 The first one comes from the JUNO official MC simulation J21v1r0-Pre2 (released the 18th August
 1531 2021). This historical version is the one on which the classical SPMT reconstruction algorithm was
 1532 developed. This classical methods is based on the time likelihood presented section 2.8 for the vertex
 1533 reconstruction, and compute the energy by correcting the detector effect on the ration N_{pe}/E_{dep} . It is
 1534 detailed in Chapter 4 of [65]. This dataset is used as a reference for comparison to classical algorithm
 1535 performances. The data in this dataset is *detsim* level (see section 2.6) which includes no digitization,
 1536 no DAQ and therefore no reconstruction of PMT signals. Only the number of PEs that hit a PMT and
 1537 the hit times are provided. A fast simulation based on gaussian drawings produces charges, with
 1538 bias and variability, and the equivalent for times. The drawings parameters were adjusted based on
 1539 [23, 74]. Because there is no charge reconstruction, the timing on the event is based on the Geant4
 1540 simulation, and so $t = 0$ is the moment the positron is created in the CD. To prevent correlation
 1541 between the numerical value of the time of the first hit t_0 and the radius of the event, we offset all
 1542 time by this first hit time. Without simulation of the charge reconstruction, we cannot simulate the
 1543 event trigger, we thus add an arbitrary time cut at a $t_0 + 1000$ ns.

1544 J23

1545 The second comes from the JUNO official monte-carlo simulations J23.0.1-rc8.dc1 (released the 7th
 1546 January 2024). The data is *calib* level (see section 2.6). Here the charge comes from the waveform
 1547 integration, the time window resolution and trigger decision are all simulated inside the software.

1548 To put in perspective this amount of data, the expected IBD rate in JUNO is 47 / days. Taking into
 1549 account the calibration time, and the source reactor shutdown, it amount to $\sim 94'000$ IBD events
 1550 in 6 years. With this million of event, we are training the equivalent of ~ 10 years of data. With
 1551 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
 1552 event scenarios that could be happening in the detector.

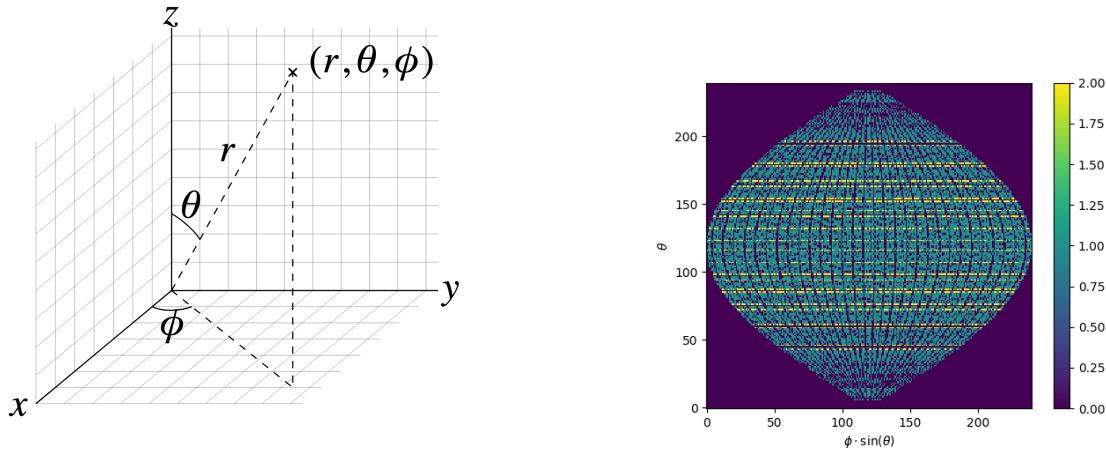
1553 While we expect and hope the MC simulation to give use a realistic representation of the detector,
 1554 there could be effect, even after the fine-tuning on calibration data, that the simulation cannot handle.
 1555 Thus, once the calibration will be available, we will need to evaluate, and if needed retrain, the
 1556 network on calibration data to establish definitive performances.

1557 The simulated data is composed of positron events, uniformly distributed in the CD volume and in
 1558 kinetic energy over $E_k \in [0; 9]$ MeV producing a deposited energy $E_{dep} \in [1.022; 10.022]$ MeV. This is
 1559 done to mimic the signal produced by the IBD prompt signal. Uniform distributions are used so that

1560 the CNN does not learn a potential energy distribution, favoring some part of the energy spectrum
 1561 instead of other.

1562 4.1.4 Data characteristics

1563 To delve a bit into the kind of data we will use, you can find in figure 4.2b the repartition of the
 1564 SPMTs in the image. The color represent the number of SPMTs per pixel.



(A) Spherical coordinate system used in JUNO for reconstruction

(B) Repartition of SPMTs in the image projection. The color scale is the number of SPMTs per pixel

FIGURE 4.2

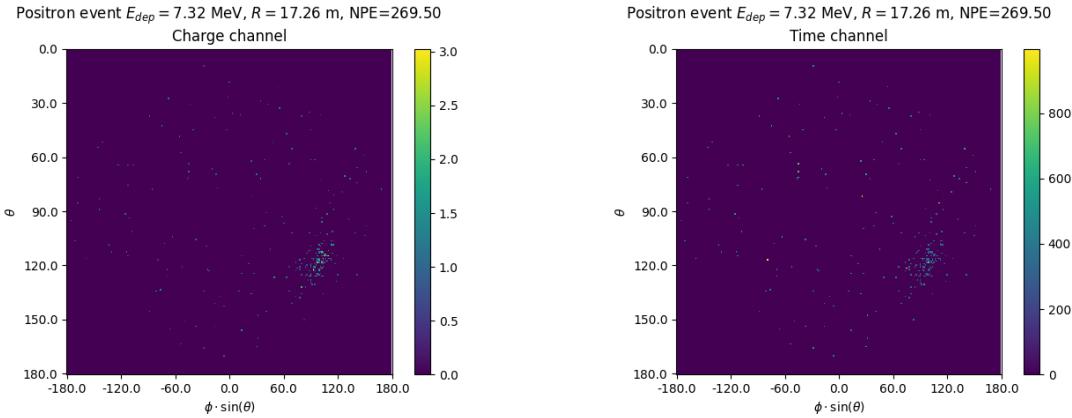


FIGURE 4.3 – Example of a high energy, radial event. We see a concentration of the charge on the bottom right of the image, clear indication of a high radius event. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

1565 See also figures 4.3 to 4.6 - and the explanation in their captions - which present events from J23 for
 1566 different positions and energies. We see some characteristics and we can instinctively understand
 1567 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

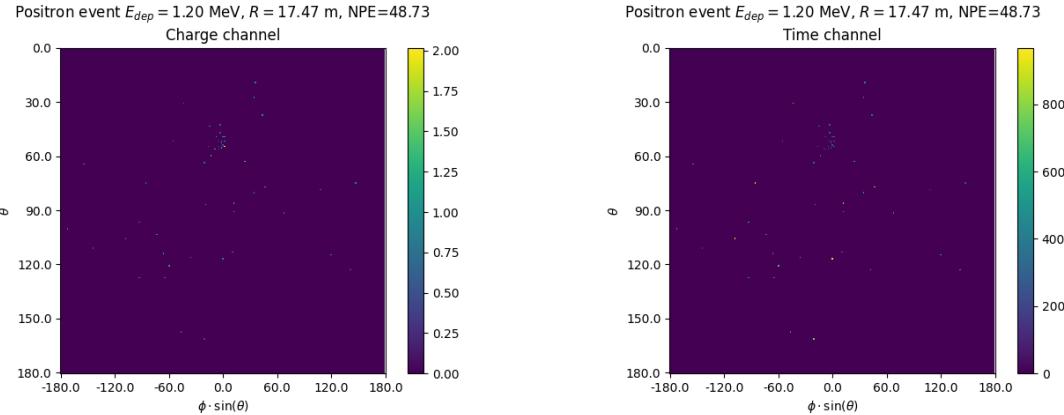


FIGURE 4.4 – Example of a low energy, radial event. The signal here is way less explicit, we can kind of guess that the event is located in the top middle of the image. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in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

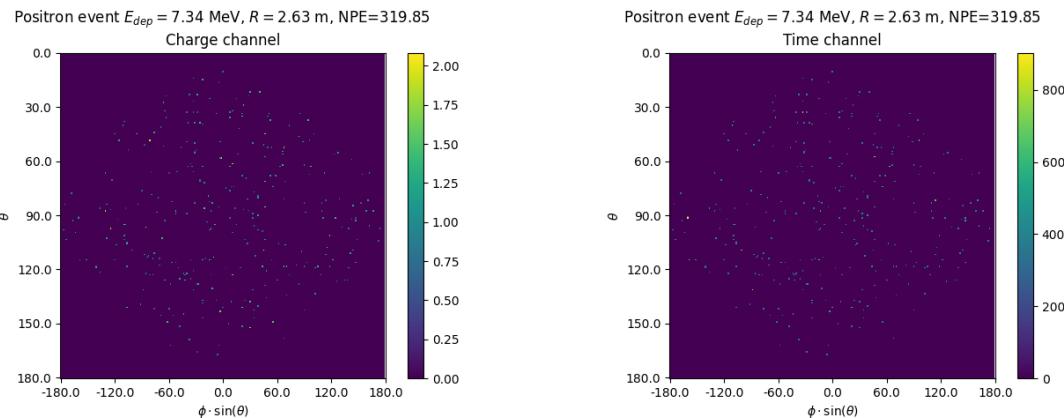


FIGURE 4.5 – Example of a high energy, central event. In this image we can see a lot of signal but uniformly spread, this is indicative of a central event. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

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 are heavily dominated by the stochastic behavior of light emission and detection.

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.

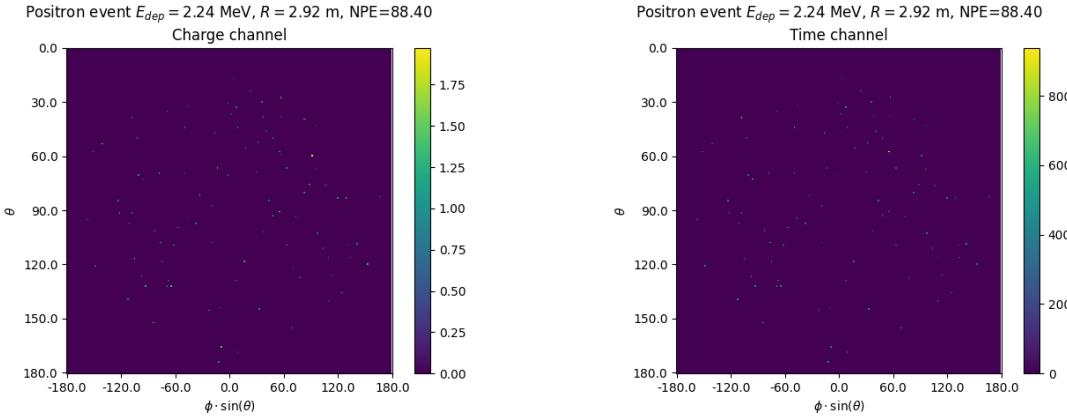


FIGURE 4.6 – Example of a low energy, central event. Here there is no clear signal, the uniformity of the distribution should make it central. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

4.2 Training

1575 The optimizer used for the training is the Adam [53] optimizer, with a learning rate λ of $1e-3$. The
 1577 other hyperparameters were left to their default value ($\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 1e^{-8}$). The
 1578 learning rate was reduced exponentially during the training at a rate of $\gamma = 0.95$, thus $\lambda_{i+1} = 0.95\lambda_i$
 1579 where i is the epoch.

1580 Following the lifecycle presented in section 3.1.3, the training used a batch size of 64 events meaning
 1581 that, each step, the loss is computed on 64 events before updating the NN parameters. An epoch is
 1582 composed of 10k steps, thus each epoch, the NN sees 640k events. The training last for 30 epochs, so
 1583 overall the NN goes through 19.2 millions events or 19.2 times the dataset.

1584 The number of epoch, batch size, learning rate and its decay were fine-tuned during the development
 1585 of the CNN.

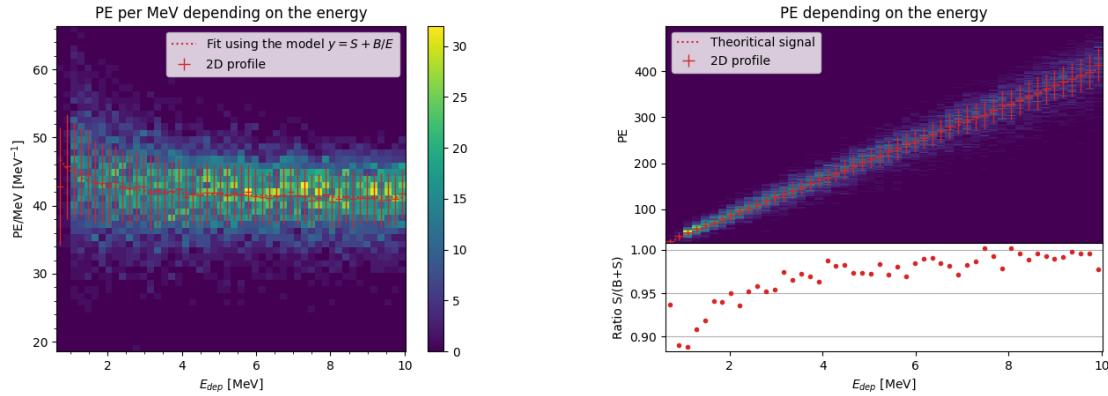
4.3 Results

1587 Before presenting the results, let's discuss the different observables.

1588 The events are considered point-like in this study. The target truth position, or vertex, is the mean
 1589 position of the energy deposits of the positron and the two annihilation gammas. This approximation
 1590 for point-like interaction is also used for the likelihood study presented in section 2.8 and in previous
 1591 ML studies presented in section 2.8.3 [42].

1592 Due to the symmetries of the detector, we mainly consider and discuss the bias and precision evolution
 1593 depending on the radius R but we will still monitor the performances depending on the spherical
 1594 angle θ and ϕ . From the detector construction and effect we expect dependency in radius due to the
 1595 TR area effect presented in section 2.8 and the possibility for the positron or the gammas to escape
 1596 from the CD for positrons interacting near the edge. We also expect dependency on θ , the top of the
 1597 experiment being non-instrumented due to the filling chimney. It is also to be noted that the events
 1598 in the dataset are uniformly distributed in the CD, and so are uniformly distributed in R^3 and ϕ . The
 1599 θ distribution is not uniform and we will have more events for $\theta \sim 90^\circ$ than $\theta \sim 0^\circ$ or $\theta \sim 180^\circ$.

1600 We define multiple energy in JUNO:



(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

- E_ν : The energy of the neutrino.
- E_k : The kinetic energy of the resulting positron from the IBD.
- E_{dep} : The deposited energy of the positron and the two annihilation gammas.
- E_{vis} : The equivalent visible energy, so E_{dep} after the detector effect such as the LS response non-linearity.
- E_{rec} : The reconstructed energy by the reconstruction algorithm. The expected value depend on the algorithm we discuss about. For example the algorithm presented in section 2.8 reconstruct E_{vis} while the ones presented in section 2.8.3 reconstruct E_{dep} .

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

4.3.1 J21 results

The best results comes from the Gen₃₀ model, meaning then 30th model generated using the table 4.1: Gen₃₀: $N_{blocks} = 3$, $N_{channels} = 32$, FCDNN configuration: $2048 * 2 + 1024 * 2$, Loss $\equiv E + V$.

The performances of its reconstruction are presented in blue in figure 4.8. Superimposed in black is the performances of the classical algorithm from [65].

Energy reconstruction

By looking at the figure 4.8a and 4.8b, the CNN has similar performances in its energy resolution. Important biases, however, appear at low and high energy.

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

The first observation is a physics phenomena: for a positron, its minimum deposited energy is the mass energy coming from its annihilation with an electron 1.022 MeV. There is a few event with

energies inferior to 1.022 MeV, in those case the annihilation gammas or even the positron escape the detector. The deposited energy in the LS is thus only a fraction of the energy of the event.

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

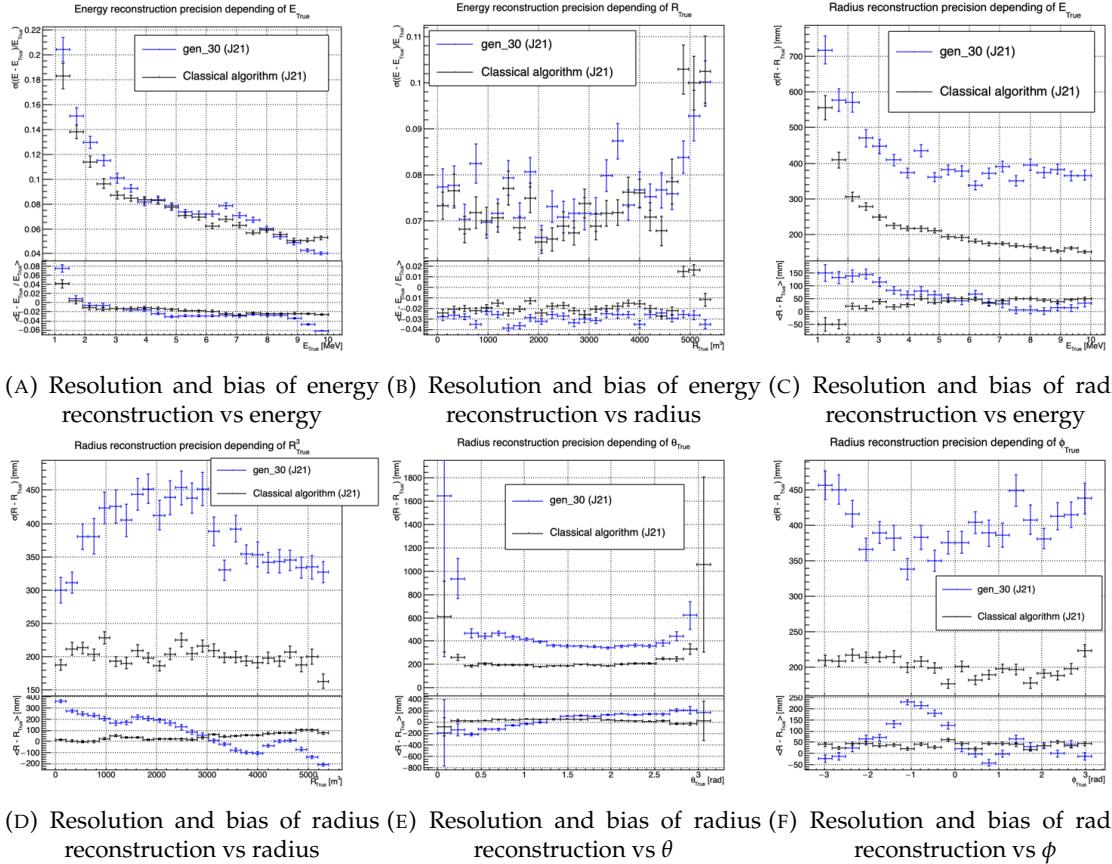


FIGURE 4.8 – Reconstruction performance of the Gen₃₀ model on J21 data and its comparison to the performances of the classic algorithm “Classical algorithm” from [65]. The top part of each plot is the resolution and the bottom part is the bias.

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

Vertex reconstruction

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

The comparison of radius reconstruction between the classical algorithm and Gen₃₀ are presented in the figures 4.8c, 4.8d, 4.8e and 4.8f. The resolution obtained by the CNN is twice worse in average,

and worse in all studied regions. 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.

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 4.3.2. Over radius, θ and ϕ 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.

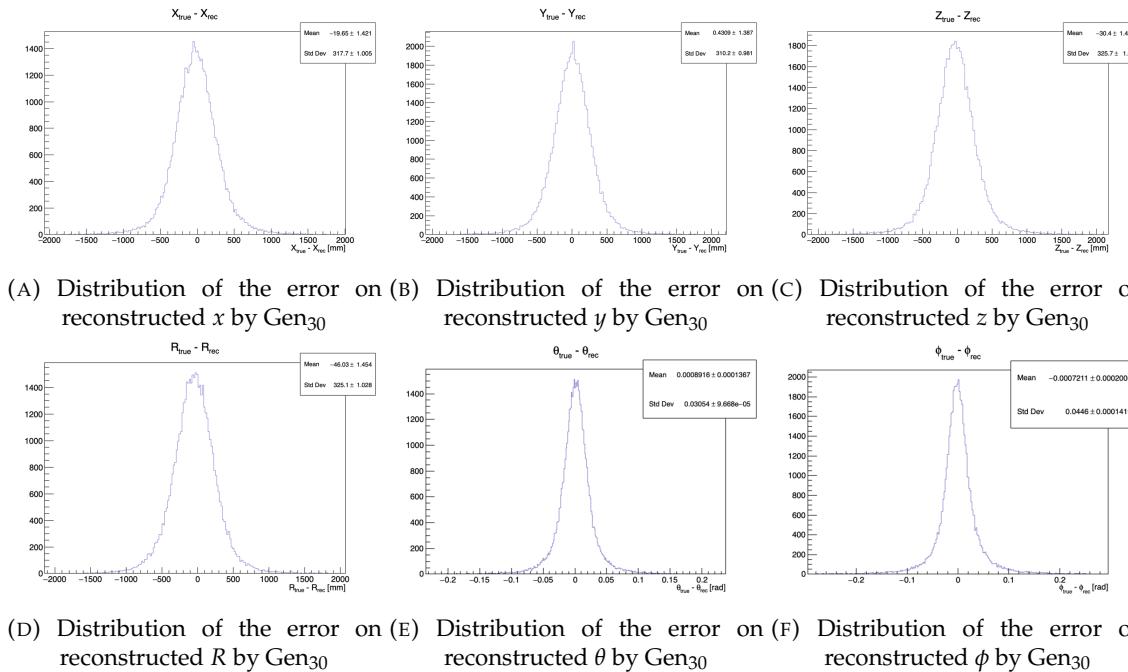
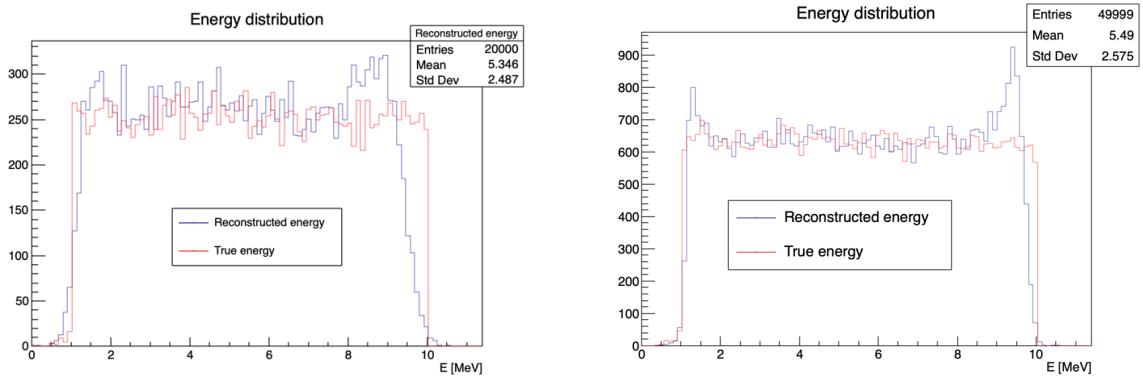


FIGURE 4.9 – Residual distribution of the different component of the vertex by Gen₃₀. The reconstructed component are x , y and z but we see similar behavior in the error of R , θ and ϕ .

As mentioned in the introduction of this chapter, this CNN initially served as a tool for learning about machine learning and JUNO’s detector and software. It eventually became necessary for use as an SPMT reconstruction tool in Chapter 7, so we made some optimizations. However, we did not invest much time in fully addressing its issues.

4.3.2 J21 Combination of classic and ML estimator

As it has been presented in previous section, there is 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



(A) Distribution of Gen₃₀ reconstructed energy and true energy of the analysis dataset (J21)

(B) Distribution of Gen₄₂ reconstructed energy and true energy of the analysis dataset (J23)

FIGURE 4.10

algorithm are respectively designated as $\theta_N(X)$ and $\theta_C(X)$.

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

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

where μ is the bias of the estimator and σ^2 its variance.

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

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

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

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

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

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

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

Where σ_{NC} is the covariance between θ_N and θ_C and ρ_{NC} their correlation.

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

By generalising the equation 4.9 to

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

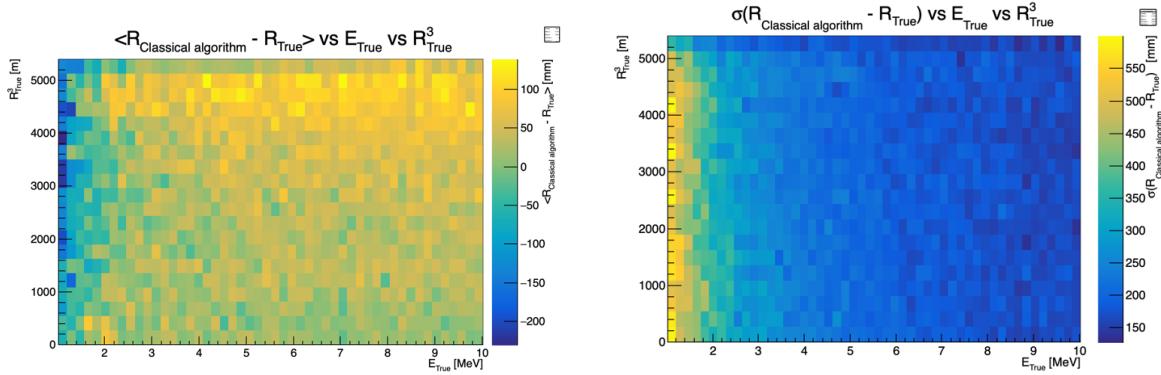


FIGURE 4.11 – Radius bias (on the left) and resolution (on the right) of the classical algorithm in a E, R^3 grid

we can determine an optimal α for two combined estimators. The estimators with the smallest variance

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

and the estimator without bias

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

See annex A for demonstration.

We present in this section the result of the estimator with the smallest variance.

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

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

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

This technique is not suited for realistic reconstruction, we rely too much on the knowledge of the resolution, bias and correlation between the two methods. While this is possible to determine using simulated data or calibration sources, the real data might differ from our model and we would need to really well understand the behavior of the two system. But this is a good tool to detect that algorithms don't all use the same information, and is a first step to identify new information that could be brought to the best algorithms, to improve their performance.

4.3.3 J23 results

We needed for Chapter 7 a SPMT reconstruction tool to run the comparison with LPMT. We thus retrained the SPMT CNN on newer, more realistic data.

The J21 simulation is fairly old and newer version, such as J23, include refined measurements of the

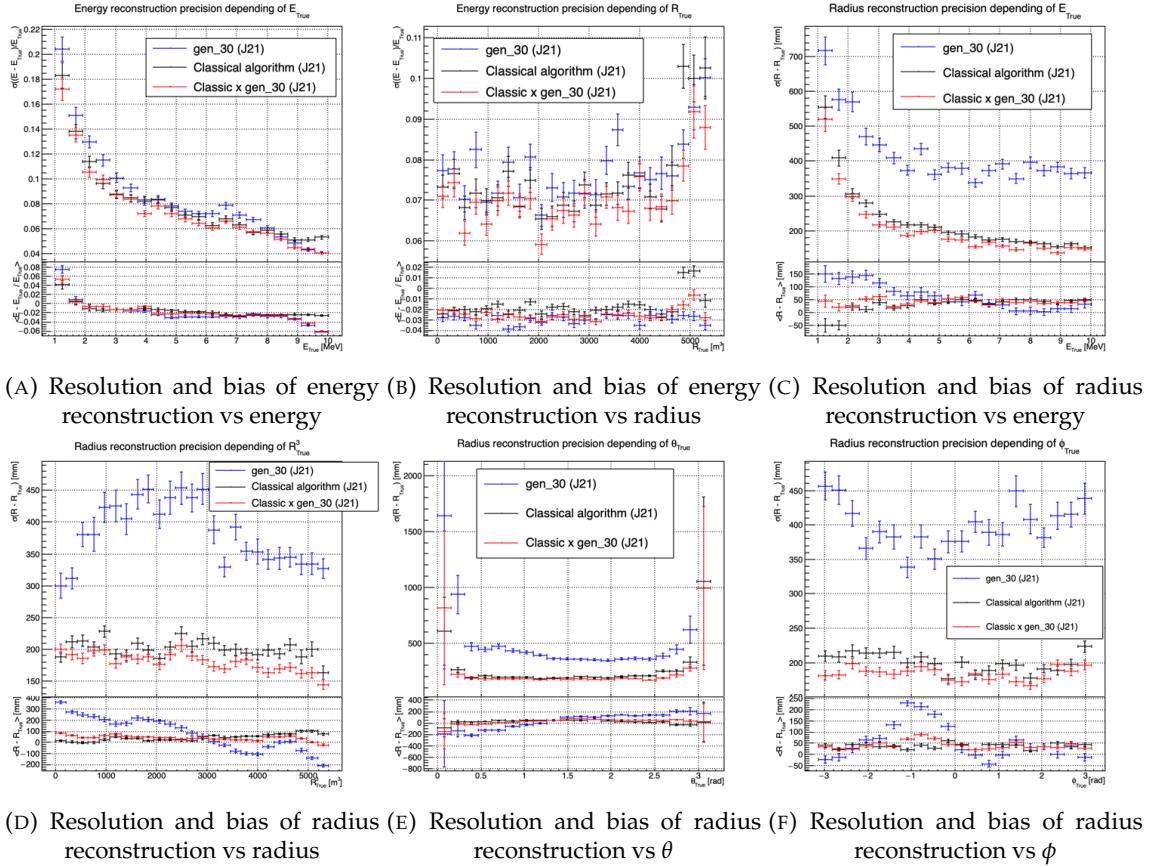


FIGURE 4.12 – Reconstruction performance of the Gen30 model on J21, the classic algorithm “Classical algorithm” from [65] and the combination of both using weighted mean. The top part of each plot is the resolution and the bottom part is the bias.

light yield, reflection indices of materials of the detector, structural elements such as the connecting structure and more realistic dark noise. Additionally, the trigger, waveform integration and time window are defined using the algorithms that will ultimately be used by the collaboration to process real physics events.

We retrained the models defined in 4.1.1 on the J23 data and used the same hyperparameter optimisation procedure. The results from the best architecture, Gen₄₂, are presented in figure 4.14. Following the table 4.1, Gen₄₂: $N_{blocks} = 3$, $N_{channels} = 64$, FCDNN configuration: $4096 * 2$, Loss $\equiv E + V$.

1701 Energy reconstruction

1702 The results of the energy reconstruction are presented in figures 4.14a and 4.14b. The resolution is
1703 close to the one of the classical algorithm with the exception of the start and end of the spectrum.
1704 This is the same effect that we saw with Gen₃₀, events are pulled from the edge of the distribution,
1705 resulting in smaller resolution but heavy biases.

1706 Vertex reconstruction

1707 The vertex reconstruction, presented in figures 4.14c, 4.14d, 4.14e and 4.14f is not yet to the level of
1708 the classical reconstruction but the degradation is smaller than for Gen₃₀ being at most a difference

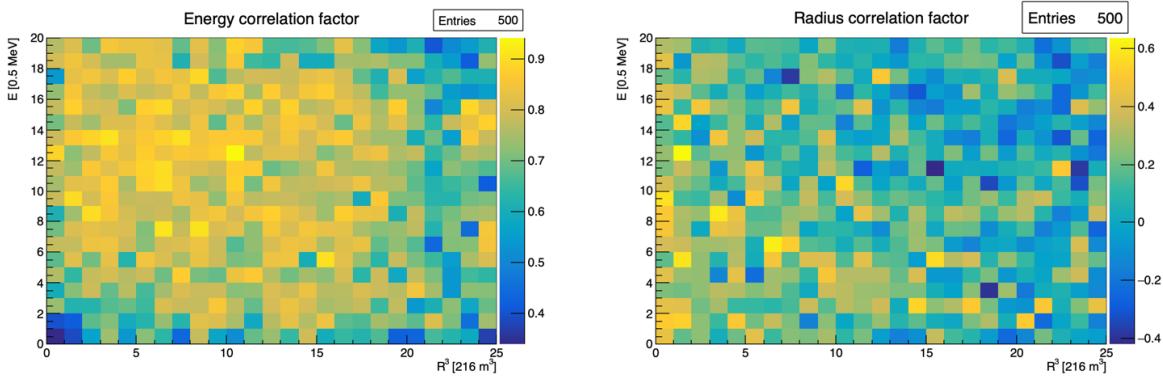


FIGURE 4.13 – Correlation between CNN and classical method reconstruction (on the left) for energy and (on the right) for radius in a E, R^3 grid

of 15cm of resolution and closing to the performance of the classical algorithm in the most favourable condition. Gen₄₂ 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.

With a more realistic description of the propagation and collection of scintillation photons, of the charge and time resolutions, of the DN and of the trigger, it seems new features can be identified by the CNN.

Unfortunately could not rerun the classical algorithm 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.

4.4 Conclusion and prospect

In this chapter we have developed a CNN for the reconstruction of IBD prompt signals. This work was the opportunity to learn about machine learning and neural networks, and familiarise ourselves with JUNO's detector and software.

This work was revisited for the needs of Chapter 7, providing a reconstruction tools for the SPMT.

The CNN we developed suffers limitations in its performance. We think one of the reasons for this lies in 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 reconstruct the θ and ϕ 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 time problematic, i.e. what is the time of a PMT that was never 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.

Another possibility is to use a kind of algorithm that does not impose a planar projection, like a GNN. It has other advantages, as will be presented in the next chapter, where we propose a GNN to reconstruct IBD's with the LPMT system.

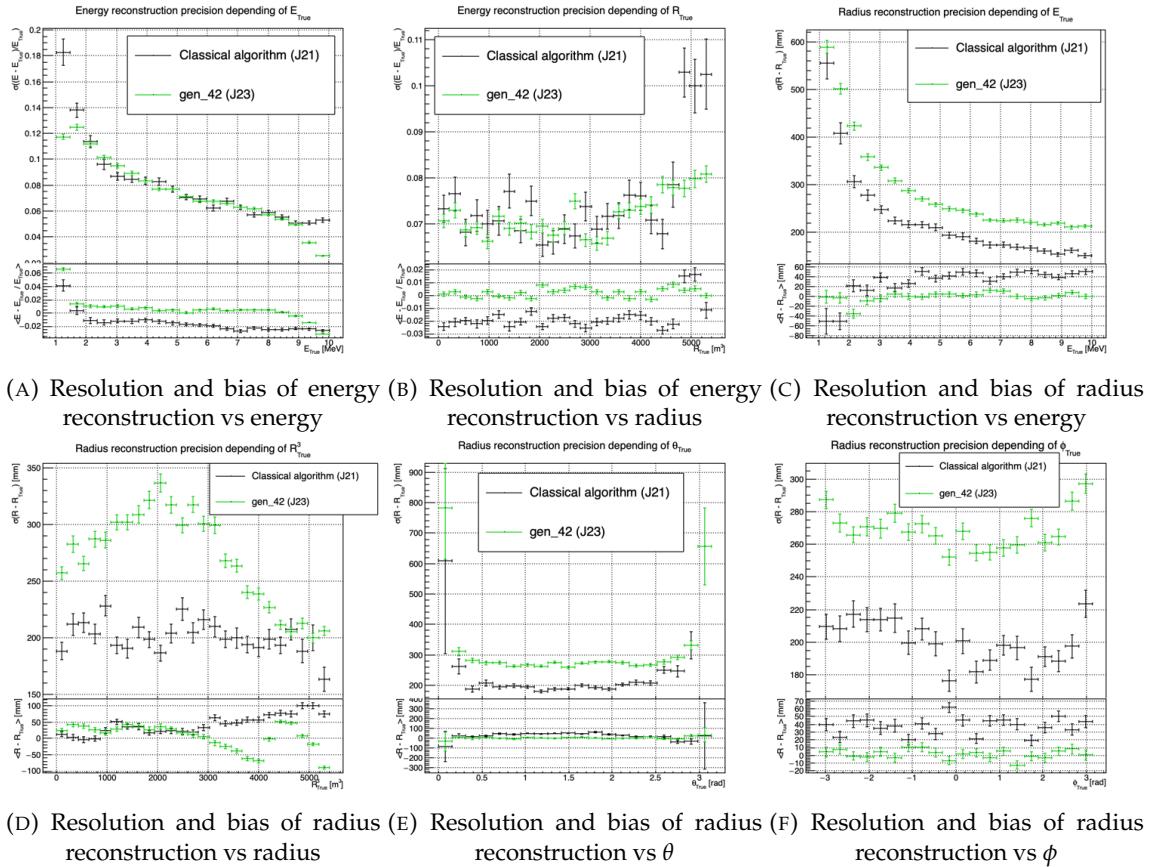


FIGURE 4.14 – Reconstruction performance of the Gen42 model on J23 data and its comparison to the performances of the classic algorithm “Classical algorithm” from [65]. The top part of each plot is the resolution and the bottom part is the bias.

¹⁷³⁸ **Chapter 5**

¹⁷³⁹ **Graph representation of JUNO for
IBD reconstruction**

¹⁷⁴¹

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

Douglas Adams, The Hitchhiker's Guide to the Galaxy

¹⁷⁴²

Contents

1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755

5.1 Data representation	78
5.2 Message passing algorithm	80
5.3 Data	82
5.4 Model	84
5.5 Training	84
5.6 Optimization	86
5.6.1 Software optimization	86
5.6.2 Hyperparameters optimization	87
5.7 performance of the final version	87
5.8 Conclusion	91

1756
1757
1758
1759
1760
1761
1762

In section 2.8.3, we showed that all ML methods developed before this thesis to reconstruct IBDs have similar results, and that their performance is very similar to that of the classical, likelihood-based algorithm. We think these similarities can reasonably be explained by this: the input data used by all these methods to compute E or \vec{X} is the same full list of PMT integrated signals $\{(Q_i, t_i); i \in 1, \dots, N_{PMTs}\}$, and by the high level of sophistication of the detector's description in the likelihood. It's probable that the likelihood method looses very little information.

1763
1764
1765
1766
1767
1768

May be some was, but that the ML algorithms were not designed well enough to recover it. It's also reasonable to think that ML algorithms will make a difference when, instead of the list of (Q_i, t_i) , a rawer information will be used in input, like the full waveform. To actually be able to learn from such a complex and high dimensional input, well designed architectures (that would guide the learning toward the solution) are necessary. In any case, it seemed welcome to us to propose an additional algorithm, with an original architecture.

1769
1770
1771
1772

For the fist stage of its development, the purpose of this part of my thesis, we considered it was enough to also take the (Q_i, t_i) list as the input. In case better of equivalent performance would be achieved, we could hope the architecture would make a difference when more complex inputs would be used. If not, we can conclude it's probably not relevant.

1773
1774
1775

The algorithm we propose is a GNN. It also has the advantage of addressing sphericity issues described in Chapter 4. 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,

1776 i.e. the energy and radius of the event do change by rotation our referential. For more details see
 1777 section 3.2.3. Graph representation also has the advantage to be able to encode global and higher
 1778 order informations.

1779 5.1 Data representation

1780 In section 2.8.3, we mentioned a GNN developed before the beginning of this thesis to reconstruct
 1781 IBD energies in JUNO [42]. In their approach: nodes of the graph correspond to 3072 pixels representing
 1782 geometric regions of the detector and the information of the ~ 6 LPMTs found in a pixel are then
 1783 aggregated on those nodes. The network then process the data using the equivalent of convolution
 1784 but on graph [49]. In the first layer, each node is connected only with its direct neighbours.

1785 To determine the energy released by an IBD in the LS, it is helpful to determine the position of
 1786 the main energy deposit. Therefore, relative Q and t's of PMTs all around the sphere is a useful
 1787 information. If in the first layer only neighbour nodes are linked, several layers are necessary to
 1788 access this detector-wide information. In an ideal world, we would develop a Graph NN where each
 1789 PMT is a node (even if it has not been hit in the event under consideration, since this is in itself an
 1790 information) and where each node is connected to all the other ones. This makes the detector-wide
 1791 information available as early as the first layer. This architecture might help the network to better
 1792 learn. Such an architecture can also be motivated this way: one of the strength of GNN's is their
 1793 capacity to encompass the characteristics of a detector. A node can be the representation of a detector
 1794 element, and the edge can represent its relationship with other elements. In the case of JUNO, any
 1795 measurement is collective : an interaction is seen by all the PMTs, with no a priori hierarchy in the
 1796 role of each. A fully connected GNN, in that respect, seems to make sense.

1797 Another advantage of a GNN is also that it is well adapted to inhomogenous detectors. We therefore
 1798 tried to build GNNs including both LPMTs and SPMTs.

1799 With 17612 LPMTs and 25600 SPMTs, the ideal fully connected Graph mentioned above is impossible:
 1800 even excluding self relation and considering the relation to be undirected (the edge from a node A
 1801 to a node B being the same from as the one from B to A) the amount of necessary edges would be
 1802 $n(n - 1)/2$ with $n = 43212$ nodes. This amounts to 933'616'866 edges. If we encode an information
 1803 with double precision (64 bits) in what we call an adjacency matrix, illustrated in figure 3.12, each
 1804 information we want to encode in the relation would consume 4 GB of data. When adding the
 1805 overhead due to gradient computation during training, this would put us over the memory capacity
 1806 of a single V100 gpu card (20 GB of memory). We could use parallel training to distribute the training
 1807 over multiple GPU but we considered that the technical challenge to deploy this solution was too
 1808 high.

1809 We finally decided of a middle ground where we define three *families* of nodes:

- 1810 — The core of the graph is composed of nodes representing geometric regions of the detector.
 1811 We call those nodes **mesh** nodes. Those mesh nodes are all connected to each other. We keep
 1812 their number low to gain in memory consumption.
- 1813 — PMTs in which Photo-Electrons (PE) are found are represented by **fired** nodes. Fired nodes
 1814 are connected to the mesh node they geometrically belong to.
- 1815 — A final node is called the input/output node (**I/O**). It is connected to every mesh node. Its
 1816 features are combinations of signals found in the whole detector.

1817 Those nodes and their relations are illustrated in figure 5.1a. From this representation, we end up
 1818 with three distinct adjacency matrices

- 1819 — A $N_{\text{fired}} \times N_{\text{mesh}}$ adjacency matrix, representing the relations between fired and mesh. Those
 1820 relations are undirected.

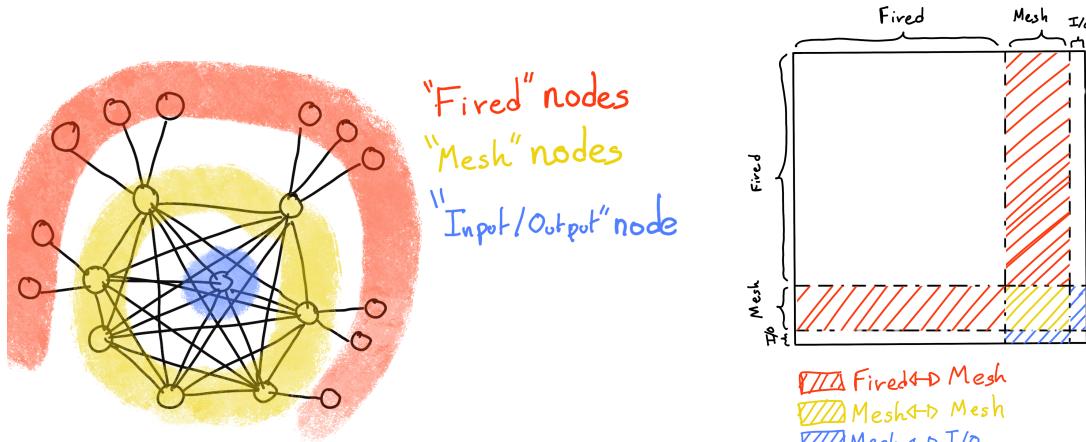


FIGURE 5.1



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

- 1822 — A $N_{mesh} \times N_{mesh}$ adjacency matrix, representing the relation between meshes. Those relation
1823 are directed.
1824 — A $N_{mesh} \times 1$ adjacency between the mesh and I/O nodes. Those relations are undirected.
1825 The adjacency matrix representing those relation is illustrated in figure 5.1b.

1826 The mesh segmentation is following the Healpix segmentation [75]. This segmentation offer the
1827 advantage that almost each mesh have the same number of direct neighbours and it guarantee that
1828 each mesh represent the same extent of the detector surface. The segmentation can be infinitely
1829 subdivided to provide smaller and smaller pixels. The number of pixel follow the order n with
1830 $N_{pix} = 12 \cdot 4^n$. This segmentation is illustrated in figure 5.2. To keep the number of mesh small, we
1831 use the segmentation of order 2, $N_{pix} = 12 \cdot 4^2 = 192$.

1832 We decided on having the different kind of nodes **mesh (M)**, **fired (F)** and **I/O** have different set of
1833 features. The features used in the graph are presented in tables 5.1 and 5.2. Most of the features
1834 are low level informations such as the charge or time information but we include some high order
1835 features such as

- 1836 1. P_l^h : Is the normalized power of the l th spherical harmonic. For more details about spherical

1837 harmonics in JUNO, see annex [B](#).

2. \mathbb{A} and \mathbb{B} are informations that are related the likeliness of the interaction vertex to be on the segment between the center of two meshes.

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

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

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

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

1838 where \vec{i} is the position vector of the mesh i , D_{ij} is the distance between the center of the meshes
 1839 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
 1840 i and n the optical index of the LS. \mathbb{A} is the vertex between center of meshes distance ratio
 1841 between i and j based on the time information. For \mathbb{B} , the charge ratio evolve with the square
 1842 of the distance, so the mesh couple with the smallest \mathbb{B} should be the one with the interaction
 1843 vertex between its two center.

Fired	Mesh	I/O
Q	$\langle Q_m \rangle$	$\langle X \rangle$
t	σQ_m	$\langle Y \rangle$
x	$\min(t_m)$	$\langle Z \rangle$
y	$\max(t_m)$	$\sum Q$
LPMT/SPMT: 1/-1	σt_m X_m Y_m Z_m	$P_l^h; l \in [0, 8]$

TABLE 5.1 – Features on the nodes of the graph. All charge are in [nPE], time in [ns] and position in [m].

Q and t are the reconstructed charge and time of the hit PMTs. (x, y, z) is the position of the PMTs and the last parameter represent the type of the PMT. It's 1 for LPMT and -1 for SPMT

Q_m and t_m is the set of charges and time of the PMT belonging the mesh m . (X_m, Y_m, Z_m) i the position of the center of the geometric region represented by the mesh m

$(\langle X \rangle, \langle Y \rangle, \langle Z \rangle)$ is the position of the charge barycenter, $\sum Q$ the sum of the collected charge in the detector and P_l^h is the relative power of the l th harmonic. See annex [B](#) for details.

1844 Since our different nodes do not have the same number of features, they exist in distinct spaces.
 1845 Traditional graph neural networks only handle homogeneous graphs, where the nodes and edges
 1846 have the same number of features at each layer. Therefore, the libraries and publicly available
 1847 algorithms we found were not suited to our needs. As a result, we had to develop and implement a
 1848 custom message-passing algorithm capable of handling our heterogeneous graph.

5.2 Message passing algorithm

1850 As introduced in previous section and in the tables [5.1](#) and [5.2](#), our graphs nodes and edges will
 1851 have different number of features depending on their nature, meaning that we cannot have a single

Fired → Mesh	Mesh ($m1$) → Mesh ($m2$)	Mesh → I/O
$x - X_m$	$X_{m1} - X_{m2}$	$\langle X \rangle - X_m$
$y - Y_m$	$Y_{m1} - Y_{m2}$	$\langle Y \rangle - Y_m$
$z - Z_m$	$Z_{m1} - Z_{m2}$	$\langle Z \rangle - Z_m$
$t - \min(t_m)$	$\min(t_{m1}) - \min(t_{m2})$	$\sum Q_m / \sum Q$
$Q / \sum Q_m$	$\frac{\langle Q_{m1} \rangle - \langle Q_{m2} \rangle}{\langle Q_{m1} \rangle + \langle Q_{m2} \rangle}$ $D_{m1 \rightarrow m2}^{-1}$ \mathbb{A} \mathbb{B}	$\langle t_m \rangle$

TABLE 5.2 – Features on the edges on the graph. It use the same notation as in table 5.1. $D_{m1 \rightarrow m2}^{-1}$ is the inverse of the distance between the mesh $m1$ and the mesh $m2$. The features \mathbb{A} and \mathbb{B} are detailed in section 5.1

1852 message passing function. We thus need to define a message passing function for each transition
1853 inside or outside a family. Using the notation presented in section 3.2.3

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

and denoting the mesh nodes M , the fired nodes F and the I/O node IO , we need to define

$$\begin{aligned} & \phi_{u;F \rightarrow M}; \phi_{m;F \rightarrow M} \\ & \phi_{u;M \rightarrow F}; \phi_{m;M \rightarrow F} \\ & \phi_{u;M \rightarrow M}; \phi_{m;M \rightarrow M} \\ & \phi_{u;M \rightarrow IO}; \phi_{m;M \rightarrow IO} \\ & \phi_{u;IO \rightarrow M}; \phi_{m;IO \rightarrow M} \end{aligned}$$

1854 to update the nodes after each layers. Following the illustration in figure 5.3, for each transition
1855 between families or inside a family we need an aggregation, a message and an update function. For
1856 the aggregation, we use the sum. We use the same, simple, formalism for every ϕ_u :

$$\phi_u \equiv I_{i'}^{n'} = I_i^n A_{i',e}^i W_n^{e,n'} + I_i^n S_n^{n'} + B^{n'} \quad (5.6)$$

1857 using the Einstein summation notation. The second order tensor, or matrix, I_i^n is holding the nodes
1858 informations with i the node index and n the feature index. n represent the features of the previous
1859 layer and n' the features of this layer.

1860 $A_{i',e}^i$ is the adjacency tensor, discussed in the previous section, representing the edges between the
1861 node i' and the node i , each edges holding the features indexed by e . If the edge does not exist, the
1862 features are set to 0. This choice is justified by the linearity of the operation in equation 5.6 : whatever
1863 the weights, when multiplied by 0 the results is 0 and the sum result is unchanged.

1864 The learnable parameters are composed of:

- 1865 — The third order tensor $W_n^{e,n'}$ which represent the passage from the previous combined feature
1866 space between the node and the edge features $n \otimes e$, the previous layer, to the current space
1867 n' , this layer.
- 1868 — The first order tensor $B^{n'}$ which is a learnable bias on the new features n' .
- 1869 — The second order tensor $S_n^{n'}$, which can be viewed as a self loop relation where the node update
1870 itself based on the previous layer informations, going from the previous space n to the current
1871 space n' .

1872 If a node have neighbours in different families, the different IAW coming from the different families

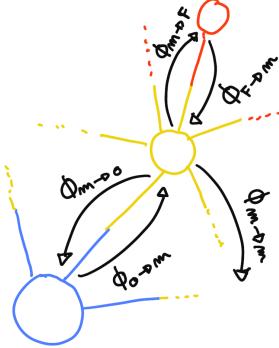


FIGURE 5.3 – Illustration of the different update function needed by our GNN

1873 are summed.

$$I' = \sum_{\mathcal{N}} [I_{\mathcal{N}} AW] + IS + B \quad (5.7)$$

where \mathcal{N} are the neighbouring family. In our case, dropping the tensor indices and indexing by family for readability, we get

$$I'_F = I_M A_{M \rightarrow F} W_{M \rightarrow F} + I_F S_F + B_F \quad (5.8)$$

$$I'_M = I_F A_{F \rightarrow M} W_{F \rightarrow M} + I_M A_{M \rightarrow M} W_{M \rightarrow M} + I_{IO} A_{IO \rightarrow M} W_{IO \rightarrow M} + I_M S_M + B_M \quad (5.9)$$

$$I'_{IO} = I_M A_{M \rightarrow IO} W_{IO \rightarrow M} + I_{IO} S_{IO} + B_{IO} \quad (5.10)$$

1874 We thus have a S , W and B for each of the ϕ_u function we defined above. The IAW sum can be
 1875 seen as the ϕ_m function and $IS + B$ as the second part of the ϕ_u function. Eq 5.5 gave the generic
 1876 form of message passing : to update a node i , one first combines informations from the surrounding
 1877 nodes and edges and then combine the result ($\square_j \phi_m$) with the current features of node i . Many
 1878 practical ways to combine can be tried. In our implementation of message passing (Eq. 5.6 and 5.7)
 1879 the latter combination is the simple sum of the former (IAW , the equivalent of $\square_j \phi_m$) with a linear
 1880 combination of the current features of node i ($IS + B$).

1881 Interestingly, the number of learnable weight in those layer is independent of the number of nodes
 1882 in each family and depends solely on the number of features on the nodes and the edges.

1883 The expression above only update the node features. We could update the edges, using the results of
 1884 ϕ_m for example, but for technical simplicity we only update the nodes and keep the edges constant.
 1885 Preserving the edges after each layers allow to share the adjacency matrix between all layers, saving
 1886 memory and computing time.

1887 This operation of message passing is the constituent of our message passing layers, designed in this
 1888 work as *JWGLayer*, each of them owning their own set of parameter W , S and B . To those layers, we
 1889 can adjoin an activation function such as *PReLU*

$$I' = PReLU \left(\sum_{\mathcal{N}} [I_{\mathcal{N}} AW] + IS + B \right) \quad (5.11)$$

1890 5.3 Data

1891 For this study we will be using a 1M positrons event dataset, uniformly distributed in energy with
 1892 $E_k \in [0, 9]$ MeV and uniformly distributed in the detector. Those events come from the JUNO

1893 official simulation version J23.0.1-rc8.dc1. All the event are *calib* level, with simulation of the physics,
 1894 electronics, digitizations and triggers. 900k events will be used for the training, 50k for validation
 1895 and loss monitoring and 50k for the results analysis in section 5.7. Each events is between 2k and
 1896 12k fired PMTS, resulting in fired nodes being the largest family in our graphs in all circumstances
 1897 as illustrated in figure 5.4c.

1898 As expected, by comparing the scale between the figure 5.4a and 5.4b we see that the LPMT system
 1899 is predominant in term of informations in our data. The number of PMT hits grow with energy but
 1900 do not reach 0 for low energy event due to the dark noise contribution which seems to be around
 1901 1000 hits per event for the LPMT system (left limit of figure 5.4a) and around 15 hits per event for the
 1902 SPMT system (left limit of figure 5.4b) which is consistent with the results show in section 4.1.2.

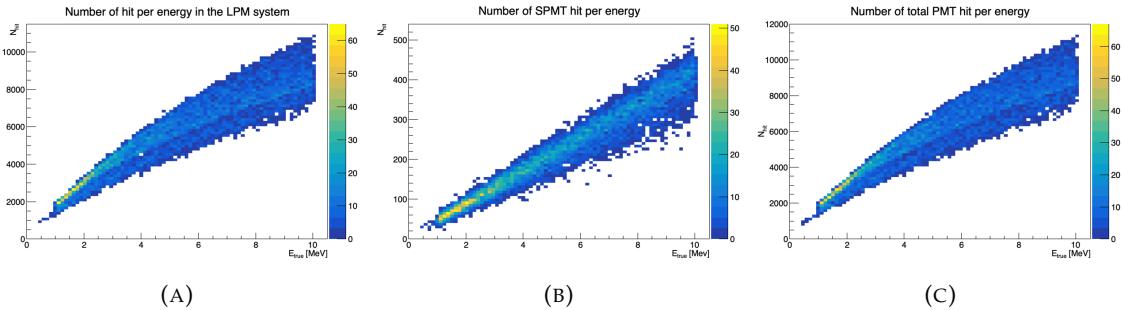


FIGURE 5.4 – Distribution of the number of hits depending on the energy. **On the right:** for the LPMT system. **In the middle :** for the SPMT system. **On the left:** For both system.

1903 The structure seen in the distribution in figure 5.4a comes from the shape of the number of hits
 1904 depending on the radius as shown in figures 5.5a and 5.5b where the number of hit decrease with
 1905 radius. It is important to understand that this is not representative of the number of PE per event
 1906 and the decrease in hits over the radius means that the PE are just more concentrated in a smaller
 1907 number of PMTs.

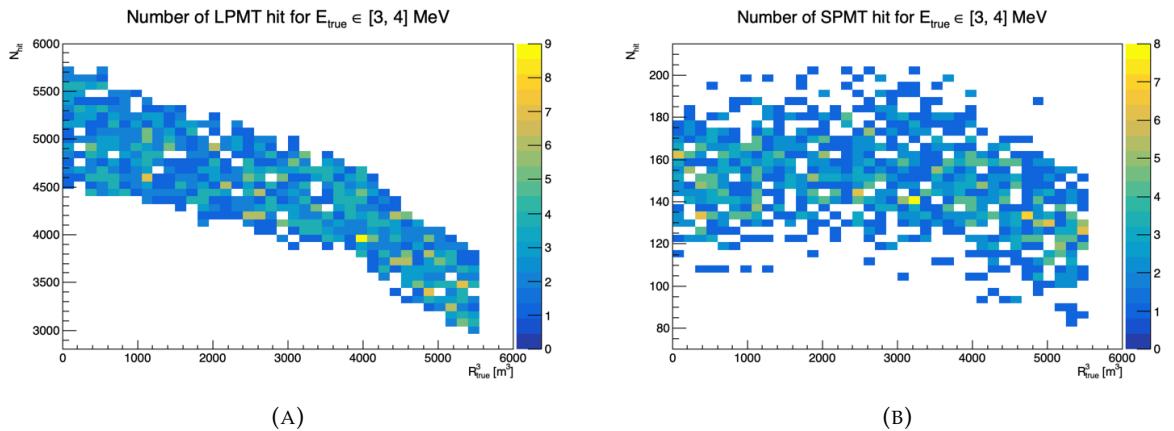


FIGURE 5.5 – Distribution of the number of hits depending on the radius. **On the right:** for the LPMT system. **On the right :** for the SPMT system. To prevent the superposition of structure of different scales we limit ourselves to the energy range $E_{true} \in [0, 9]$.

1908 No quality cut is applied here, we rely only on the trigger system. It means that event that would not
 1909 trigger are not present in the dataset but for events that triggered twice, it happens rarely, the two
 1910 trigger are considered as two separate event.

5.4 Model

In this section, we discuss the different layers that compose the final version of the model. The number of layers, their dimensions, and their arrangement were fine-tuned through multiple iterations. As mentioned earlier, 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 previous 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 tables 5.1 and 5.2 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 5.6.

We use the Mean Square Error (MSE) for the loss

$$\mathcal{L} = (E_{rec} - E_{dep})^2 + (X_{rec} - X_{true})^2 + (Y_{rec} - Y_{true})^2 + (Z_{rec} - Z_{true})^2 \quad (5.12)$$

as it was the best resulting loss in Chapter 4.

5.5 Training

The optimizer used for training is the Adam optimizer and default hyperparameters ($\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 1e-8$) with a learning rate $\lambda = 1e-8$. The training last 200 epochs of 800 steps. We use a batch size of 32, the largest we can have with 40GB of GPU ram. The learning rate is constant during the first 20 epochs then exponentially decrease with a rate of 0.99. We save two set of parameters, the set of parameters the set that yield the lowest validation loss and the set of parameters at the end of the training. The validation is computed over a single batch.

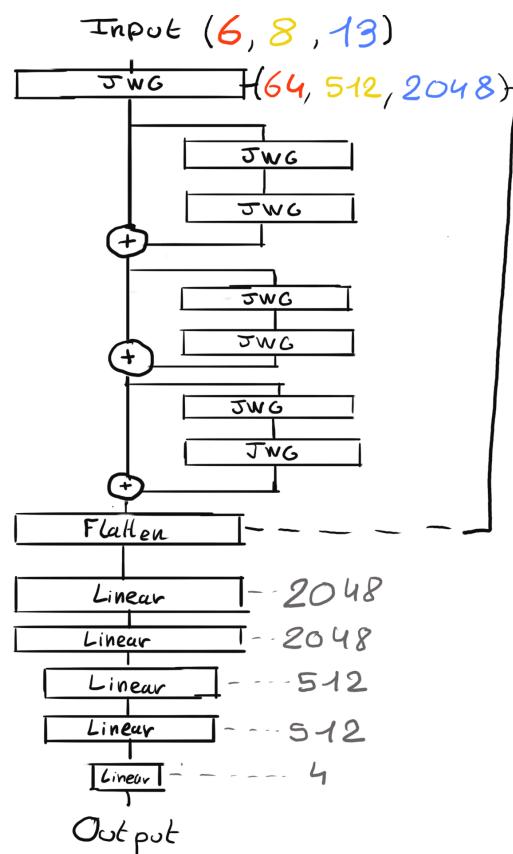


FIGURE 5.6 – Schema of the JWGv8.4.0 architecture, the colored triplet is the graph configuration after each JWG layers

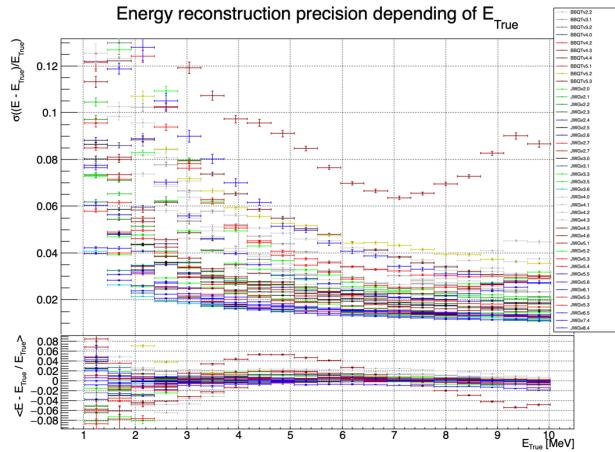


FIGURE 5.7 – Energy reconstruction depending on the true energy for samples of the different versions of the GNN

1949 5.6 Optimization

1950 The GNN model presented in previous sections is the result of a long work of optimization. Indeed,
 1951 the innovative architecture we propose left us with an infinity of possible configurations with no
 1952 guidance from prior works in literature nor in JUNO.

1953 In the end, more than 60 different configurations have been tested. This effort is illustrated on Figure
 1954 5.7¹, where the 40 configurations are compared in their ability to reconstruct the positron energy.
 1955 Although all configurations share the fundamental principles we base our innovative architecture
 1956 on (three different kinds of nodes and edges, usage of raw level features on some of them, usage of
 1957 higher level data on others, division of JUNO’s surface into regional pixels to form mesh nodes, the
 1958 very large number of edges connected to each mesh node, etc.), performances can vary a lot between
 1959 our first attempts (far beyond any acceptable energy resolution, and not even on this figure) and
 1960 recent ones. Therefore: the precise way to choose hyperparameters mattered a lot, regardless of the
 1961 relevance of the global architectural principles.

1962 The spectacular improvement between early and later configurations also explains the length of this
 1963 process : for long we hoped we would finally reach the classical performance, and it was tempting
 1964 to test yet another configuration.

1965 5.6.1 Software optimization

1966 A substantial effort was devoted to the data processing workflow. Transforming JUNO simulation
 1967 outputs into graphs is a computationally expensive task. Furthermore, due to the ever-changing
 1968 nature of the graph dimensions and features during optimization, preprocessing JUNO’s files by
 1969 precalculating the graphs and then reading them from files was not viable, as it would require a
 1970 large amount of disk space to store events for each version of the graph.

1971 Therefore, the software does not rely on preprocessed data and instead computes the observables,
 1972 adjacency matrix, etc., during training. This data processing is performed in parallel on the CPU.
 1973 The raw data comes from ROOT files produced by the collaboration software, and the Event Data
 1974 Model (EDM), used internally by the collaboration [76], had to be interfaced with our software,
 1975 an interface that had to be maintained as the collaboration’s software evolved. For the harmonic

1. Note that this figure was prepared on idealized data with no dark noise and perfect hit time determination.

1976 power calculation, we migrated from the Healpix library to Ducc0 [77] for more precise control over
1977 multithreading.

1978 5.6.2 Hyperparameters optimization

1979 The first kind of hyper-parameters that received a lot of effort concern the network's detailed architecture:
1980

- 1981 — Message passing layers where originally not JWG layers, we started by using small FCDNN
1982 in place of ϕ_u and ϕ_m . Due to low performances and memory consumption issues, we pivoted
1983 to the message passing algorithm presented in section 5.2.
- 1984 — The ResNet architecture was brought after issue with the gradient vanishing.
- 1985 — The number of layers was varied between 5 and 12.
- 1986 — The number of node features after each given message passing layer (64, 512, 2048 in the final
1987 version) was varied.
- 1988 — The Final FCDNN after the message passing layers is not present in all versions.
- 1989 — At some point, the PReLU activation function replaced the ReLU function.

1990

1991 For some of them, software work was necessary. In any case, each configuration required a training
1992 of about 90h. Adding the analysis time necessary to the verification of its performance and the
1993 comparison with other versions, one understands the number of tests had to be limited.

1994 Other hyperparameters were also tested :

- 1995 — The higher level variables described in section 5.1 (powers of various spherical harmonics, \mathbb{A} ,
1996 \mathbb{A} , $(Q_{m1} - Q_{m2})/(Q_{m1} + Q_{m2})$) were added progressively. Notice that our choice to focus
1997 our search on this kind of variables is also due to the fact that JWGLayer involves linear
1998 operations. It is therefore difficult for such a network to propose variables of this kind among
1999 the node features learned layers after layers (i.e. it's difficult for the network to understand
2000 these variables are important, or only after many layers).
- 2001 — Time allocated to training, the Learning Rate, the size of batches, etc.
- 2002 — The number of pixels (ie of mesh nodes) was varied between 192 and 768.
- 2003 — Several definitions loss functions where tried. In particular, we tried some focussed only on
2004 the E resolution, only on the vertex resolution (R) or trying to optimize both.

2005

2006 To make a long story short, each new configuration was the result of our reflections after having
2007 analysed the previous configurations, or after having thought over again about JUNO's detailed
2008 response to energy deposits – seeking for variables that could help the GNN.

2009 Another, quite common, approach was in principle possible : a random search. However, due to the
2010 extensive training time, up to 90h per training, the heavy memory consumption of the models that
2011 would often exceed the 20GB limit of the V100, this approach was not realistic in our case, though we
2012 were able to extend the memory limit to 40GB thanks to a local A100 GPU card available at Subatech.

2013 5.7 performance of the final version

2014 The reconstruction performance of "JWGv8.4" are presented in figures 5.8, 5.9, 5.10 and compared to
2015 the "Omilrec" algorithm, the official IBD reconstruction algorithm in JUNO. Omilrec is based on the
2016 QTMLE reconstruction method that was presented in section 2.8.

2017 This comparison required to use a consistent definition of E_{true} . This is not trivial since at JUNO,
2018 ML method reconstruct the true energy deposited by the positron+annihilation gammas (that's the

target implemented in the loss function), while Omilrec, which is based on probabilities to observe a given number of PE in a given PMT, reconstruct the "visible energy". It reflects the total number of radiated and detectable scintillation or Cherenkov photons (and is subject to non linear effects like quenching).

The conversion we use to obtain comparable E_{true} is explained in Appendix D.

On figures 5.8 to 5.10, we notice that the best GNN does not match the performance of the OMILREC algorithm. Generically, Energy resolution is 50% worse, while the resolution on R is three times worse. Reconstruction biases are not better either with the GNN. We have tried to understand the origin of this limited performance.

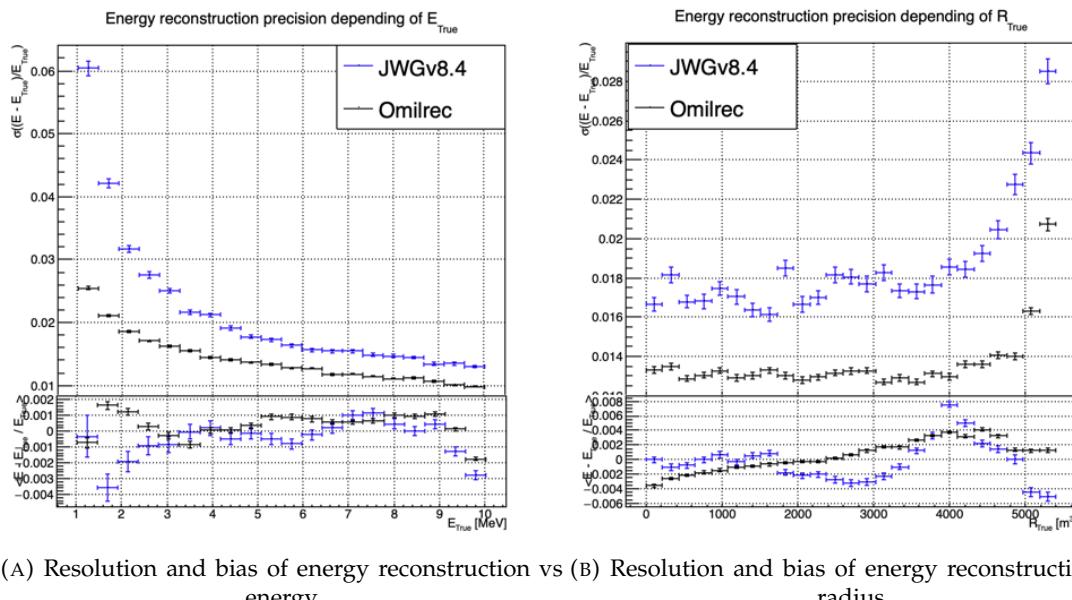


FIGURE 5.8 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.8, JWGV8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.

The first action that can be carried out in this direction was to determine if some information used by OMILREC was not used properly by JWGV8.4. For that purpose, we used again the approach presented in Chapter 4 (Sec 4.3.2 and annex A) to combine JWGV8.4 and OMILREC. We observe on figures 5.11 and 5.12 that this combination brings no sizeable improvement of the best of the two combined methods. The combination remains very close to OMILREC alone. This is an indication that JWGV8.4 does not use informations that would be overlooked by OMILREC, and that on the contrary, that's JWGV8.4 that fails to use properly important informations.

The problem described above could be inherent to our GNN's original architecture. Discussions with JUNO's colleagues when these results were presented at the collaboration pointed to the role of PMT time information (t , in the (Q, t) pairs we use as our algorithm input features). The thousands of values found in the *fired* nodes might not be aggregated well enough when transmitted to the mesh nodes, causing a loss in the redundancy of this important information.

We tested this idea in several manners, described below.

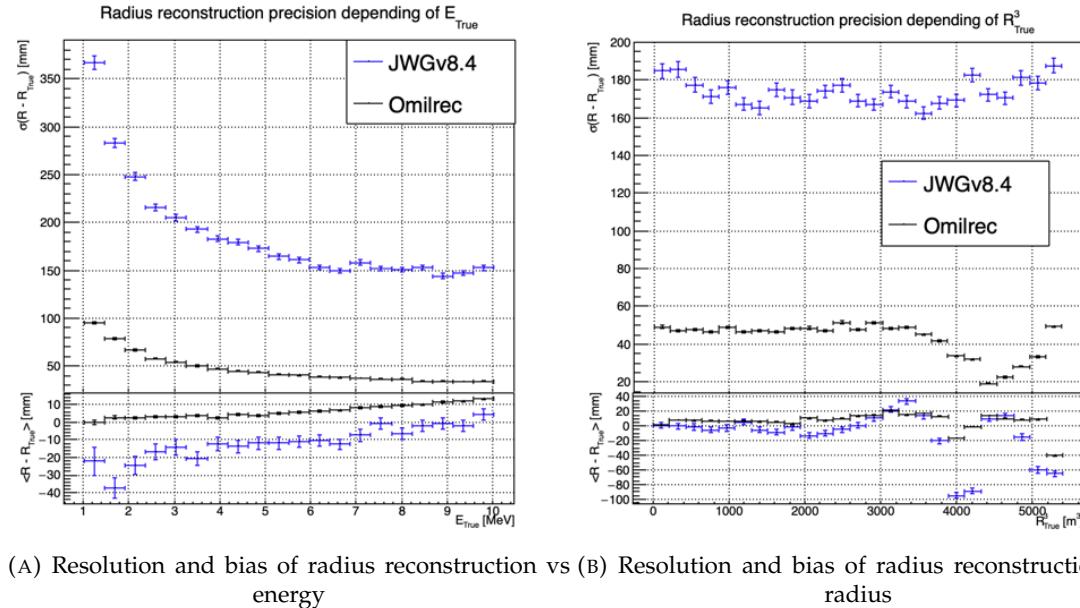


FIGURE 5.9 – Reconstruction performance of the Omilrec algorithm based on QTMLR presented in section 2.8, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.

2041 Finer granularity

2042 We tried to recover some redundancy by increasing the number of mesh nodes from 198 to 768. The
 2043 improvement we observed was small, and did not allow to get close to OMILREC's performance.

2044 To explore further in this direction, we would ideally try 3072 pixels (the next HEALPIX rank).
 2045 However, this is not possible for our GNN due to hardware limitations, mainly the available GPU
 2046 memory. Instead, we discussed the problem with Gilles Grasseau, calculus research engineer with
 2047 whom we collaborate on the subject of ML reliability (see Chapter 6). In the framework of this ac-
 2048 tivity, Gilles needs to develop reconstruction algorithms to be "attacked" by a prototype Adversarial
 2049 NN. One of them is a pseudo-spherical CNN using oriented filters, called HCNN.

2050 To produce its input image, this algorithms split the Sphere into 3072 pixels. Each channel of this
 2051 image is an aggregation of the (Q, t) values found in all the PMTs. The charge are summed and
 2052 the lowest time is kept. The performance of this algorithm can be seen on Figures 5.13 and 5.14,
 2053 compared to OMILREC. With 3072 pixels, the performance of HCNN does not match that of OMIL-
 2054 REC, but is closer to it than our GNN. The granularity of the pixels, and the way to summarize the
 2055 individual PMTs information when going from 17000 LPMTs to only 3072 pixels indeed seems to
 2056 play a role.

2057 This is consistent with the results obtained by the first GNN tried at JUNO on reactor neutrinos
 2058 (already described in section 2.8.3). It used 3072 pixels, and also obtained an uncompetitive R
 2059 reconstruction.

2060 Information reduction, from fired to Meshes

2061 The problem described above is somehow classical. ML algorithms, ideally, would start from the full
 2062 information present in the detector, and learn to reduce it optimally.

2063 In cases where only 3072 pixels can be used instead of the complete information from 17000 PMTs,

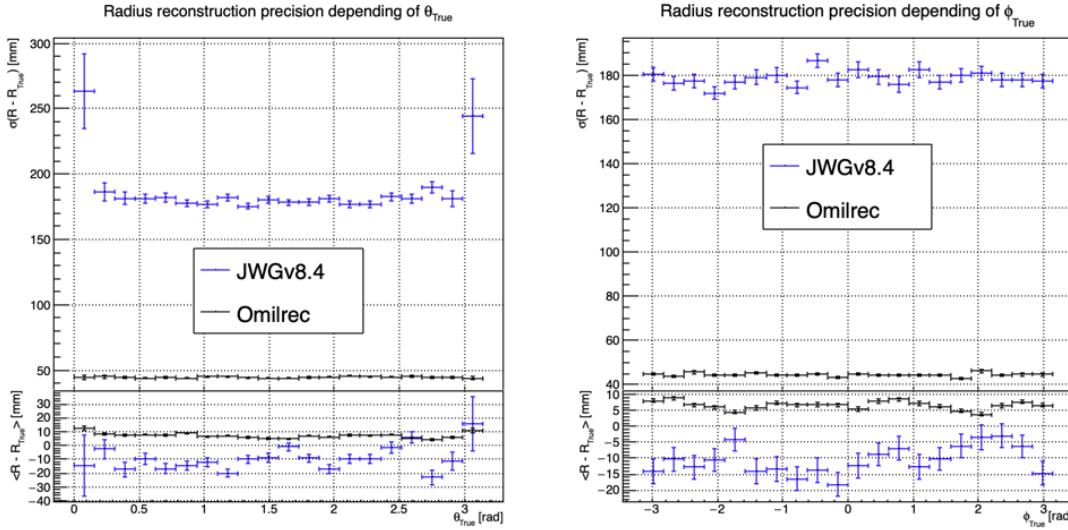
(A) Resolution and bias of radius reconstruction vs θ (B) Resolution and bias of radius reconstruction vs ϕ

FIGURE 5.10 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.8, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.

one needs to understand how to combine the individual from the 5 or 6 PMT found in each pixel into pixel-level features, without loosing important information.

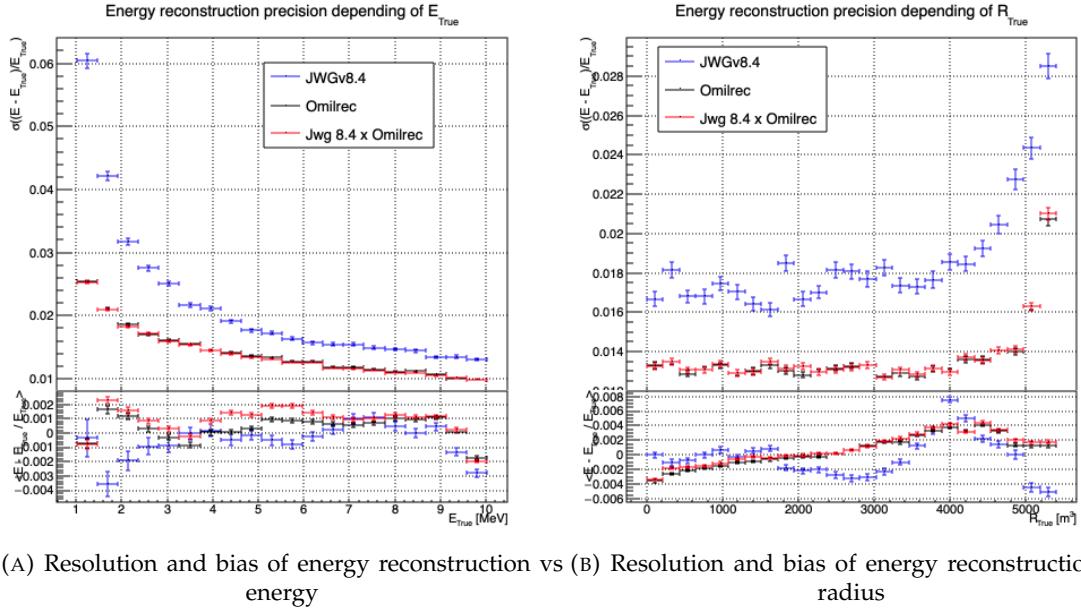
In the case of our GNN, we hoped that by connecting each mesh node to its corresponding 5 or 6 fired nodes, we could keep the full information. In reality, it seems that the message passing between fired and mesh does not work efficiently. When nodes are updated by the first (may be also by the subsequent) layer, the new mesh features might be dominated by the original features in the second column of tables 5.1, themselves a simple version of aggregation. Layer after layer, we might be limited to that level of time information, lacking time redundancy.

We have verified this by testing version of the GNN in which the link between fired and mesh was cut, or in which no time info was included among the fired nodes features. It had only a small effect which seems to confirm a problem in the way the full information, from all the individual PMTs, is used by our GNN.

2076 Possible improvements

2077 It appears that the network is unable to aggregate the timing information correctly. While this could
 2078 be addressed by using a finer segmentation, with more mesh nodes, improvements might also arise
 2079 from refining the message-passing algorithm. The algorithm presented in this thesis is still quite
 2080 basic, relying on a simple linear combination of features. We have seen through examples in CNNs,
 2081 GNNs, and other architectures, both in research and industry, that specializing the network — for
 2082 instance, by incorporating convolutional filters — can lead to improvements that were previously
 2083 unattainable with simpler FCDNNs. Applying this approach to the message-passing algorithm, by
 2084 utilizing a GNN with a more advanced message-passing, could yield better results.

2085 Regarding the timing information, we provided high-level features, assuming this would assist the
 2086 neural network in converging to the solution. However, by offering such information upfront,
 2087 the GNN might be taking the “easy” path, settling for a local and broader minimum, rather than
 2088 extracting the features that could lead to better performance.



(A) Resolution and bias of energy reconstruction vs energy (B) Resolution and bias of energy reconstruction vs radius

FIGURE 5.11 – Reconstruction performance of the Omilrec algorithm, JWGV8.4 and the combination between the two using the optimal variance estimator presented in annex A.2. The top part of each plot is the resolution and the bottom part is the bias.

If there are difficulties in transferring information between the fired and mesh nodes, it may stem from the way we connected the fired nodes to the mesh nodes. By linking the fired nodes within the same mesh, or even connecting the fired nodes of neighboring mesh nodes, the GNN might be able to construct more meaningful information.

Finally, by providing directly the PMT waveform to the GNN, in the fired nodes, we could search for even finer precision and results. An idea would be to specialise the message function $\phi_{m;F \rightarrow M}$ to be a 1D convolutional layer over the waveform. The resulting channels would be fed to the mesh nodes for their updates.

5.8 Conclusion

To achieve its scientific goals, JUNO requires a precise and well-understood reconstruction, as it needs an energy resolution of 3% at 1 MeV. Even small, unaccounted biases could make it impossible to determine the mass ordering, as explored in Chapter 7. A likelihood-based algorithm, designed to meet JUNO's requirements and referred to as the classical algorithm, was developed and is detailed in section 2.8.

Machine learning algorithms were developed to challenge this classical approach, and they are presented in Section 2.8.3. Although they achieve the precision of the classical algorithm, they do not offer significant improvements. The GNN previously developed is a convolutional GNN where nodes correspond to pixels, connected to their neighbors based on the Healpix [75] segmentation, with the (Q, t) information aggregated onto these pixels.

In this chapter, we introduce a novel and innovative architecture. In addition to the pixel segmentation represented by mesh nodes, we incorporate rawer information by directly representing the fired PMTs as nodes. We also fully connect the mesh nodes to each other, hoping to facilitate the transfer of information. Finally, we introduce a global node that holds global information about the detector.

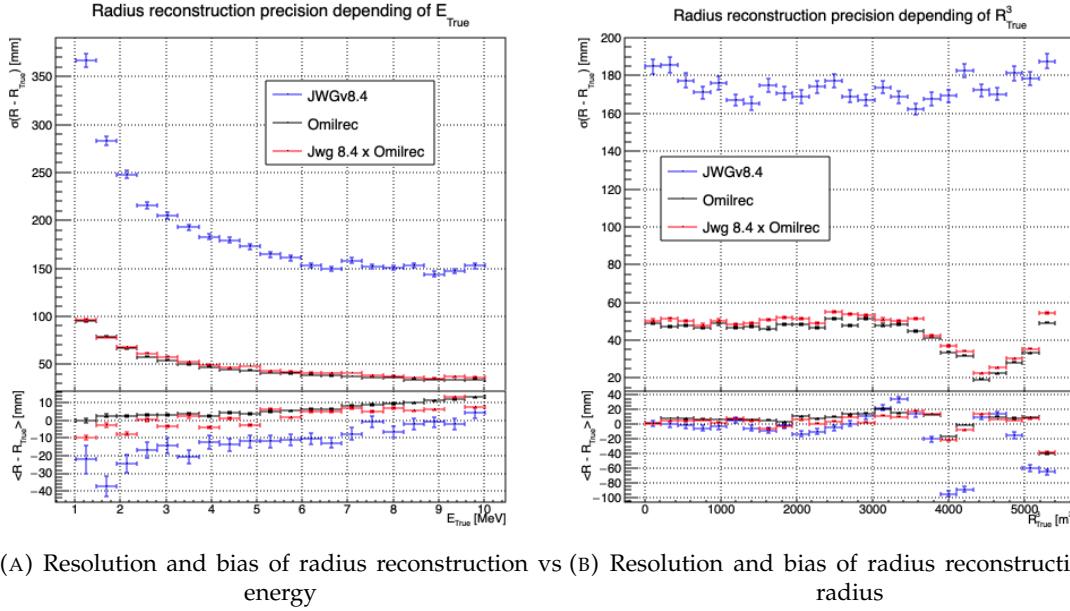
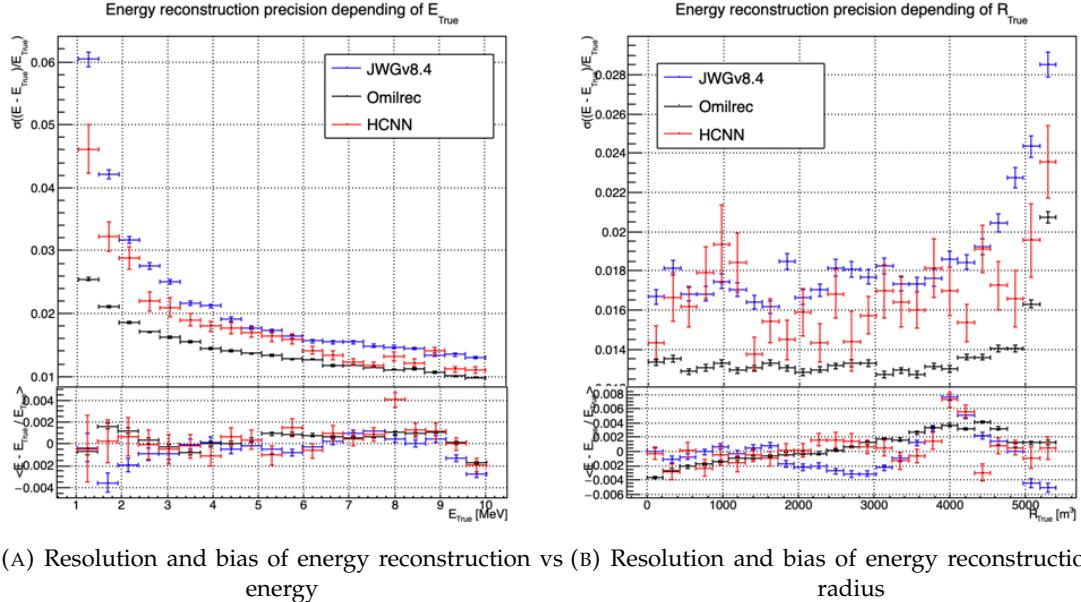


FIGURE 5.12 – Reconstruction performance of the Omilrec algorithm, JWGV8.4 and the combination between the two using the optimal variance estimator presented in annex A.2. The top part of each plot is the resolution and the bottom part is the bias.

2112 These three types, or families, of nodes do not have the same number of features, resulting in a het-
 2113 erogeneous graph. Publicly available algorithms for graph processing are designed for homogeneous
 2114 graphs, so we had to develop a custom algorithm adapted to heterogeneous graphs.

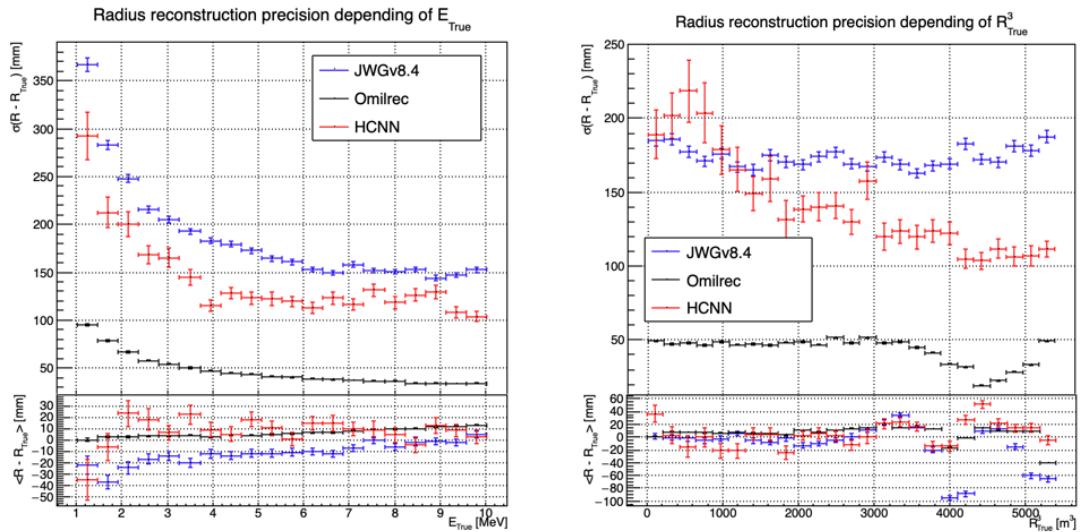
2115 This GNN required significant technical development, but the results are not at the level of the
 2116 classical algorithm. The tests we conducted suggest that the problem may lie in the aggregation
 2117 of raw information from the fired nodes onto the mesh nodes, as removing the fired nodes does
 2118 not degrade the results. Additionally, due to technical constraints, we had to reduce the number of
 2119 pixels compared to the previous GNN. Other algorithms we developed, which use a higher pixel
 2120 resolution, outperform this architecture, reinforcing our suspicion that the aggregation is the root of
 2121 the issue.

2122 Perhaps by incorporating rawer information, such as the waveform, refining the message-passing
 2123 algorithm, or adjusting the features on the different nodes, we could match the precision of the
 2124 classical algorithm. However, it is also possible that deeper, more radical changes are needed to
 2125 become competitive.



(A) Resolution and bias of energy reconstruction vs energy (B) Resolution and bias of energy reconstruction vs radius

FIGURE 5.13 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.8, JWGv8.4 presented in this chapter and the HCNN algorithm. The top part of each plot is the resolution and the bottom part is the bias.



(A) Resolution and bias of radius reconstruction vs energy (B) Resolution and bias of radius reconstruction vs radius

FIGURE 5.14 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.8, JWGv8.4 presented in this chapter and the HCNN algorithm. The top part of each plot is the resolution and the bottom part is the bias.

2126 **Chapter 6**

2127 **Reliability of machine learning
methods**

2128

2129 “*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

2130 **Contents**

<small>2131</small> 6.1 Method	<small>96</small>
<small>2132</small> 6.2 Architecture	<small>96</small>
<small>2133</small> 6.2.1 Back-propagation problematic	<small>98</small>
<small>2134</small> 6.2.2 Reconstruction Network	<small>99</small>
<small>2135</small> 6.2.3 Adversarial Neural Network	<small>99</small>
<small>2136</small> 6.2.4 Training	<small>99</small>
<small>2137</small> 6.3 Results	<small>99</small>
<small>2138</small> 6.3.1 Back to identity	<small>99</small>
<small>2139</small> 6.3.2 Breaking of the reconstruction	<small>99</small>
<small>2140</small> 6.4 Conclusion and prospect	<small>99</small>

2141 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.

2142 We already shown a large variety of reconstruction algorithms, OMILREC for LPMT reconstruction in section 2.8, numerous machine learning algorithms in section 2.8.3 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 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.

2143 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 describe the method behind the algorithm. In section 6.2 I detail the architecture of our algorithm and in section 6.3 the results of it. Finally, in section 6.4, I conclude and discuss about the prospect and possible improvements to bring to this work.

2163 **6.1 Method**

2164 As introduced above, JUNO needs a very good understanding of the biases and effects affecting its
2165 reconstruction as a small bias could wrong the mass ordering measurement. To calibrate those biases
2166 and effect, JUNO rely on multiples sources that can be located at various point in the detector. The
2167 calibration strategy was already discussed in section 2.4 and show calibrations sources of gammas,
2168 neutrons and positrons, with the catch that the positrons will annihilate inside the encapsulation and
2169 only the two 511 keV gammas will be seen.

2170 None of the calibrations sources considered are positron event. While electrons and positrons events
2171 should be pretty similar in their interaction with the electronic cloud of the LS atoms, electron
2172 events are missing the two annihilations γ and the potential of forming a positronium [78]. The
2173 topology of the event thus differ of the order of magnitude of our reconstruction performance. A
2174 few nanoseconds between the energy deposit and the positronium annihilation against a time transit
2175 spread between 3 and 6 ns depending on the PMT type [79–81]. The γ from the positron annihilation
2176 will travel distances of the order of magnitude of the typical LPMT resolution of 8 cm (see section
2177 2.8).

2178 Another natural calibration source is the ^{12}B spectrum. The ^{12}B is a cosmogenically produced isotope
2179 through the passage of muons inside the LS. The ^{12}B decays via β^- emissions with a Q value of
2180 13.5 MeV with more than 98% of the decay resulting in ground state ^{12}C . The ^{12}B event will be
2181 cleanly identified by looking for delayed high energy β events after an energetic muon. Due to its
2182 natural causes, the ^{12}B events will be uniformly distributed in the detector. The calibration strategy
2183 consist in fitting the energy spectrum of ^{12}B with the results of the simulation to adjust the simulation
2184 parameters. Both sources will be used to *control* the response of the detector.

2185 Unlike lasers and radioactive, from which the localization and energy will be well known, the in-
2186 dividual truth of ^{12}B will be unknown with only the localisation loosely constrained by the muon
2187 track. Only higher order observables such as the energy distribution will be accessible.

2188 All of those considerations could hide potential unknown or undetected effect that could lead to
2189 issue in the mass ordering analysis. But, while we have idea from where the issue could come, the
2190 production by hand of event perturbations that go unseen in the calibration would be tedious. That's
2191 why we propose to use an Adversarial Neural Network (ANN) to produce those perturbations if they
2192 exists. A schematic of the concept is presented in figure 6.1.

2193 This network should produce physically sound perturbation, that would not be seen by the calibra-
2194 tion but also by the visualisation of the event. If the ANN manage to produce such perturbations,
2195 we can derive systemic uncertainties from it. If it fail to find some, it is a proof of robustness for the
2196 attacked reconstruction method.

2197 For this study we consider a “physics” dataset composed of 1M positron events from J23, uniformly
2198 distributed in the Central Detector (CD) and in deposited energy between $E_{dep} \in [1.022; 10.022]$. This
2199 set represent the IBD events we want to *wrongly* reconstruct.

2200 We use a second “control” dataset of electron event also uniformly distributed in the detector and
2201 over the same energy range. They mimic the energy deposit of ^{12}B decay and are used as the sample
2202 to compute the control observables.

2203 **6.2 Architecture**

2204 We can describe the goal of the ANN by the following loss function:

$$\mathcal{L} = \mathcal{L}_{adv} + \mathcal{L}_{reg} \quad (6.1)$$

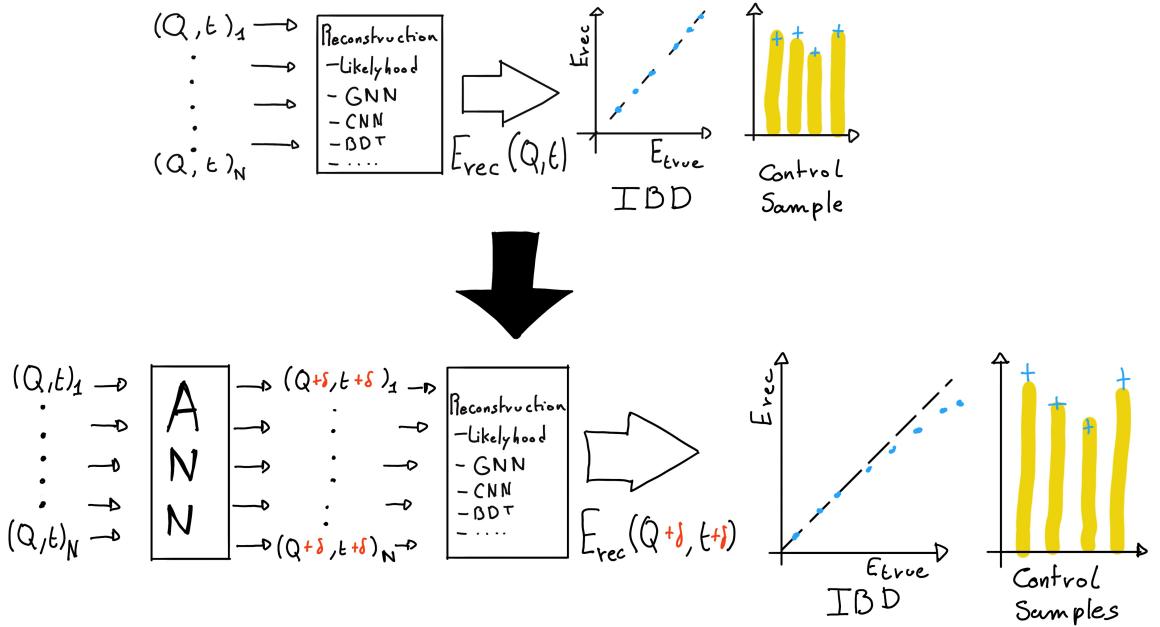


FIGURE 6.1 – Schema of the method to discover vulnerabilities in the reconstruction methods. **On the top** of the image, the standard data flow. The individual charge and times are fed to a reconstruction algorithm. From the reconstructed energies, we can produce an IBD spectrum and compute control observables from the control samples. **On the bottom**, the same data flow but we add an ANN between the input and the reconstruction. The ANN will slightly change the input charge and time so the reconstruction algorithm inaccurately reconstruct the IBD energy, but the perturbation is not visible in the control sample.

where \mathcal{L}_{adv} is the adversarial loss, which is minimal when the reconstruction is “broken”. We thus need to define what is a *wrong* reconstruction. We choose to define it via the correlation between the reconstructed and deposited energy

$$\mathcal{L}_{adv} = |\text{Corr}(E_{rec}, E_{dep})| \quad (6.2)$$

which is positive or null and is minimal when the reconstructed energy is decorrelated with the deposited energy, the reconstruction is wrong.

The term \mathcal{L}_{reg} is the regularisation term, which is minimal when the control variable are correctly reconstructed

$$\mathcal{L}_{reg} = \sum_{\lambda} (O_{\lambda}^{rec} - O_{\lambda}^{th})^2 \quad (6.3)$$

where λ index the different control observables that will be considered in this study. It’s minimal when the control observables after perturbation O_{λ}^{rec} are coherent with their expected values O_{λ}^{th} . In this exploratory work, we choose as the control observable the difference between the reconstructed position and energy and the ground truth from the Monte Carlo simulation complemented with a penalty term P

$$\mathcal{L}_{reg} = \sum_{\lambda \in \{x, y, z, E\}} (\lambda_{rec} - \lambda_{true})^2 + P \quad (6.4)$$

This penalty P is here to prevent the ANN from producing event too different from the initial event. It will be further detailed in section 6.2.4.

We see that the final loss is the equilibrium between the adversarial and regularisation loss.

2220 **6.2.1 Back-propagation problematic**

We would like this method to be applicable to any kind of reconstruction algorithm but this complicated considering standard training method through backward-propagation, discussed in details in section 3.1.3. For explanation, let's define the application of the reconstruction algorithm as \mathcal{F} on an event X , resulting in the prediction Y and the application of the ANN \mathcal{G} on X to give a perturbed event X' , we can parametrize the equation 6.1

$$Y = \mathcal{F}(X); Y' = \mathcal{F}(X') = \mathcal{F}(\mathcal{G}(X)) \quad (6.5)$$

$$\mathcal{L} \equiv \mathcal{L}(\mathcal{F}(\mathcal{G}(X)), Y_t) \quad (6.6)$$

2221 where Y_t is the reconstruction target of Y .

2223 Now if we consider a parameter θ of the ANN on which we want to optimize \mathcal{L} , in the backward-
2224 propagation optimisation framework we need to compute

$$\frac{\partial \mathcal{L}(\mathcal{F}(\mathcal{G}(X)))}{\partial \theta} \quad (6.7)$$

2225 which, when using the chain rule, become

$$\frac{\partial \mathcal{L}(\mathcal{F}(\mathcal{G}(X)))}{\partial \theta} = \frac{\partial \mathcal{G}}{\partial \theta} \cdot \frac{\partial \mathcal{F}}{\partial \mathcal{G}} \cdot \frac{\partial \mathcal{L}}{\partial \mathcal{F}} \quad (6.8)$$

2226 The terms $\frac{\partial \mathcal{G}}{\partial \theta}$ and $\frac{\partial \mathcal{L}}{\partial \mathcal{F}}$ are easily computable but $\frac{\partial \mathcal{F}}{\partial \mathcal{G}}$ depends on the nature of the reconstruction
2227 algorithm. While it comes naturally when using NN algorithms, it's not so trivial for other kind
2228 of algorithms like likelihood. Solutions exists to optimize networks that work in complex, non
2229 differentiable environments, such as *Deep Reinforcement Learning* [82, 83] but as a first prototype we
2230 will restrict ourselves to neural networks for the reconstruction algorithm.

2231 The choice to use gradient descent, and therefore neural network, also allowed us to keep all technical
2232 software development wrapped in the same language and framework, PyTorch [58].

2233 The backward-propagation introduce a second issue. At the beginning of the subsection we intro-
2234 duce $X' = \mathcal{G}(X)$, the event after perturbation. It's an input of the reconstruction \mathcal{F} , thus, let's say
2235 that the event, in its form X , is a list of tuples (id, Q, t) which are the hit on the PMT id . If \mathcal{F} require
2236 the information to be formatted in a specific way (graph, images, ...) via an algorithm $\tau(X)$, it means
2237 that

$$\frac{\partial \mathcal{L}(\mathcal{F}(\tau(\mathcal{G}(X))))}{\partial \theta} = \frac{\partial \mathcal{G}}{\partial \theta} \cdot \frac{\partial \tau}{\partial \mathcal{G}} \cdot \frac{\partial \mathcal{F}}{\partial \tau} \cdot \frac{\partial \mathcal{L}}{\partial \mathcal{F}} \quad (6.9)$$

2238 which also requires that $\frac{\partial \tau}{\partial \mathcal{G}}$ is differentiable.

2239 On the other hand, if X is already formatted as the input of \mathcal{F} , it mean that \mathcal{G} take the same format
2240 as input and we drop the requirement on τ to be differentiable. Concretely, if \mathcal{F} takes an image as
2241 input, it mean that \mathcal{G} will also takes an image as input and output an image. That also unfortunately
2242 mean that if some informations is loss before \mathcal{G} , for example during the charge and time aggregation
2243 in pixels, it cannot retrieve and modify it.

2244 A more elegant solution would that \mathcal{G} would also compute the transformation τ in addition to
2245 finding relevant perturbation, but for the simplicity of this exploratory work, we use a \mathcal{G} that process
2246 transformed data.

2247 6.2.2 Reconstruction Network

2248 As introduced just before, we need a NN algorithm for IBD reconstruction. We could have used the
 2249 GNN presented in Chapter 5 but we preferred a more simplistic approach to not be constrained by
 2250 the memory consumption of the reconstruction neural network.

2251 This network takes as input a vector containing the results of the aggregation of charge and time on
 2252 pixels. We consider JUNO composed of 3072 pixels defined by the Healpix [75] pixelisation. On each
 2253 of those pixel, we sum the charges and keep the minimal time of hit, resulting in 3072 (Q, t) tuples.
 2254 To those tuples, we adjoin the position of the center of those pixels, resulting in 3072 (Q, t, x, y, z)
 2255 tuples. The data is finally represented as a $3072 \times 5 = 15360$ vector. In the case the charge in a pixel
 2256 is 0, the time is set to 2048 ns, way after the closing of the trigger window.

2257 The simplistic neural network is simply and Fully Connected Neural Network (FCDNN) composed
 2258 of the following layer: the input layer, providing the 15360 items vector, followed by fully connected
 2259 linear layers

2260 6.2.3 Adversarial Neural Network

2261 — Décrire l'architecture de l'ANN

2262 6.2.4 Training

2263 — Présentation du dataset
 2264 — 2 étapes d'entraînement
 2265 — Retour à l'identité -> que l'ANN ne fasse pas n'importe quoi
 2266 — Cassage de la reconstruction

2267 Hyperparameter optimization

2268 — Pour les mêmes raisons que l'ANN:
 2269 — Phase exploratoire, architecture très changeante, random search n'est pas viable
 2270 — Architecture consomme beaucoup, besoin d'entrainer sur l'A100
 2271 — Possiblement que de l'optimisation permettrait de faire passer sur V100, mais développement techniques nécessaires.

2273 6.3 Results

2274 — Voir slide Gilles

2275 6.3.1 Back to identity

2276 6.3.2 Breaking of the reconstruction

2277 6.4 Conclusion and prospect

2278 — Not enough
 2279 — Probably guide the ANN

2280 **Chapter 7**

2281 **Dualcalorimetric analysis for Precision
Measurement**

2283 “We demand rigidly defined areas of doubt and uncertainty!”
Douglas Adams, The Hitchhiker’s Guide to the Galaxy

2284 **Contents**

<small>2285</small>	7.1 Motivations	<small>103</small>
<small>2286</small>	7.1.1 Discrepancies between the SPMT and LPMT results	<small>103</small>
<small>2287</small>	7.1.2 Charge Non-Linearity (QNL)	<small>104</small>
<small>2288</small>	7.2 Approach	<small>105</small>
<small>2289</small>	7.2.1 Data production	<small>106</small>
<small>2290</small>	7.2.2 Individual fits	<small>107</small>
<small>2291</small>	7.2.3 Joint fit	<small>108</small>
<small>2292</small>	7.2.4 Data and theoretical spectrum generation	<small>109</small>
<small>2293</small>	7.2.5 Limitations	<small>110</small>
<small>2294</small>	7.3 Fit software	<small>110</small>
<small>2295</small>	7.3.1 IBD generator	<small>111</small>
<small>2296</small>	7.3.2 Fit	<small>112</small>
<small>2297</small>	7.4 Technical challenges and development	<small>113</small>
<small>2298</small>	7.5 Results	<small>113</small>
<small>2299</small>	7.5.1 Validation	<small>113</small>
<small>2300</small>	7.5.2 Covariance matrix	<small>117</small>
<small>2301</small>	7.5.3 Statistical tests	<small>122</small>
<small>2302</small>	7.6 Conclusion and perspectives	<small>124</small>

2303 JUNO is precision measurement experiment. To determine the NMO with the aimed significance,
2304 JUNO must be sensitive to the tiny spectral phase shift shown on figure 7.1. Once detection effects
2305 are accounted for, the difference between IO and NO spectra is further reduced, as can be seen on
2306 figure 7.2.

2307 Among other condition, a precise and complete understanding of the reconstruction and detector
2308 effects is crucial. The challenge reside in the technology used in the detector, which, while based
2309 on well known technology: scintillator observed by PMT, is being deployed on a scale never seen
2310 before, in term of scintillator volume and PMT size. Understanding every effects that goes in the
2311 detector can become extremely complicated. Any method to help detecting problems is therefore
2312 welcome. Comparing the data and results obtained by two systems measuring the same events, but
2313 subject to different sources of error, is therefore precious. This is the purpose of the dual calorimetry
2314 techniques used in JUNO thanks to the existence of 2 PMT systems: the LPMT and SPMT systems.

The reconstruction of the IBD positron energy must be very performant: an unprecedented resolution of 3% at 1 MeV [50] is necessary to determine the NMO with the aimed significance. Moreover, it is necessary to know the energy scale with an uncertainty below 1% to correctly evaluate in our data the likelihood of the NO and IO hypotheses. Beyond that value, the risk progressively appears to exclude the NI(IO) hypothesis with the significance with which one should actually have excluded the IO(NO) if the energy scale was precisely known, as can be seen in the introduction of Chapter 4 of [24].

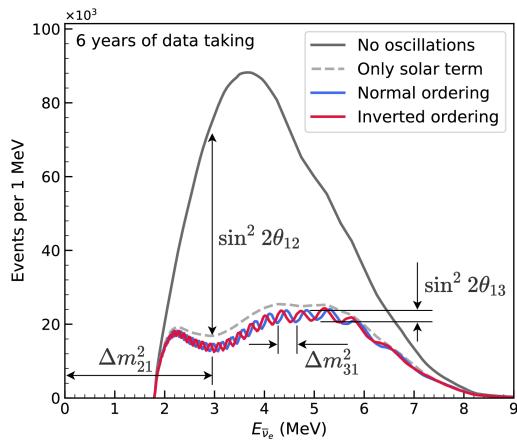


FIGURE 7.1 – Expected number of neutrinos event per MeV in JUNO after 6 years of data taking. The black curve shows the flux if there was no oscillation. The light gray curve shows the oscillation if only the solar terms are taken in account (θ_{12} , Δm_{21}^2). The blue and red curve shows the spectrum 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.

One of the possible source of non-linearity, which will be used as a reference in this chapter, is the charge non-linearity (QNL) that will be discussed in next section. Several dual calorimetry techniques can address this issue. Some are calibration techniques, that are also described in section 4.3 of [24]. More generally, comparing the results of the two systems will allow for the detection of potential issues on the calibration or reconstruction. This is done in this thesis by comparing directly the spectra and oscillation parameters measurements of the two PMT systems. We call this kind of dual calorimetry "Dual calorimetry with neutrino oscillation", since it is based on the visible energy spectra used by the oscillation analysis of reactor antineutrinos.

In this chapter, we explore several ways to perform this comparison. One of them relies on the difference between the values of Δm_{21}^2 , $\sin^2(2\theta_{12})$ measured with the LPMT and the SPMT systems. Both systems measure them with similar uncertainties. For reasonable values of the QNL, we expect these differences to be smaller than the individual uncertainties. However, the significance of these differences might still be high. Indeed, both systems reconstruct the same events, therefore the same distribution of the true positron energy, as well as the same scintillation photon emission. Therefore, the energy spectra reconstructed by the two systems share a part of their fluctuations. This translates into correlated reconstructed spectra and consequently lead to correlations between the measurements of Δm_{21}^2 and $\sin^2(2\theta_{12})$. The uncertainty on the SPMT-LPMT difference is largely decreased by this correlation. Other ways to perform the comparison (see next sections) all rely on the reconstructed spectra, therefore on the evaluation of the correlation between the LPMT and SPMT spectra.

In the next section we will discuss the motivations behind this study. In section 7.2, I present the methods we propose to implement Dual calorimetry with neutrino oscillation, and of the way we

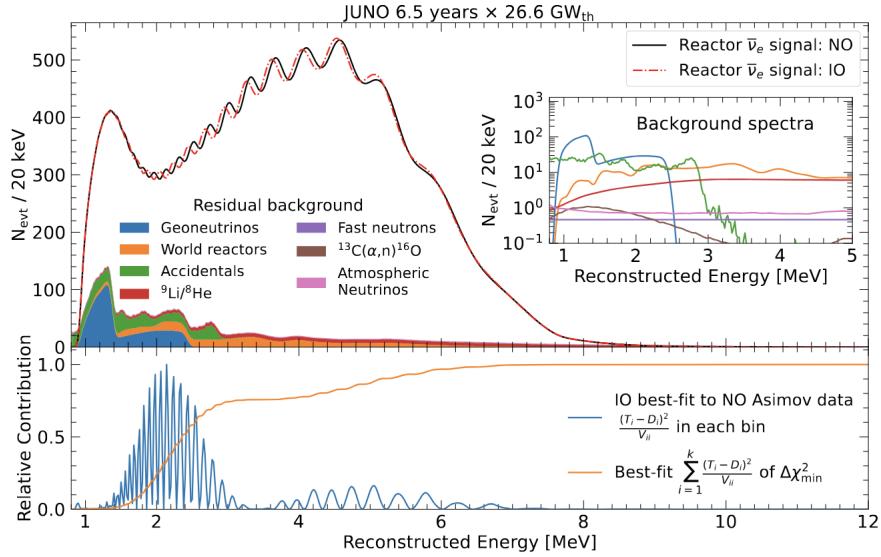


FIGURE 7.2 – Oscillated reactor $\bar{\nu}_e$ spectra for the Normal Ordering (Black) and Inverted Ordering (Red) for 6,5 years data taking and a resolution of 3% without any statistical or systematic fluctuation. Figure from [32].

2348 estimate their sensitivity. In section 7.3, I present the fit framework used, and then, in section 7.4
 2349 the technical improvement brought and the difficulties faced during the development. To end this
 2350 chapter I present the results in 7.5 and discuss the conclusions and perspectives in 7.6.

2351 7.1 Motivations

2352 7.1.1 Discrepancies between the SPMT and LPMT results

2353 As discussed in the introduction of this chapter, the SPMT and LPMT systems will observe the same
 2354 events. This mean that, after calibration, if the two system show significant differences in their results
 2355 this is the signal of potential overlook of an effect or problem. Being able to detect such differences
 2356 is thus crucial, as discussed above, even the smallest deviation from our model could lead to the
 2357 impossibility to measure the Mass Ordering (MO) or even worse, wrong our measurement.

2358 The two systems are expected to have the same sensitivity to the oscillation parameters θ_{12} and Δm_{21}^2
 2359 [3]. We will thus rely on the measurement of those two parameters to detect potential discrepancies.

2360 We could just look at the value and compare them to the estimated independent error of the two
 2361 system, but we believe and will demonstrate in this chapter that the independent study of the two
 2362 system is missing a lot of informations, and that, by taking into account the statistic and systematic
 2363 correlations between the two systems, we can produce much more powerful statistical tests.

2364 Our work in this chapter is to develop such tools, which in practice implies to define test statistics. A
 2365 first step will be to determine the distribution of these test statistics in the case when no unexpected
 2366 problem affects the LPMT nor the SPMT problem. This will give us the distribution of those statistical
 2367 test in absence of discrepancies. Later, the value of the test statistics that we will measure in real data
 2368 can be compared to these distributions to produce p-values, to judge of the potential present of an
 2369 unexpected effect.

To evaluate the power of our methods, we need to simulate a concrete difference between the two spectra. We have decided to study a plausible effect, the Charge Non-Linearity (QNL) that is detailed next section. Note that these tests should in principle be able to detect unexpected effects whatever their source (calibration issues, insufficient simulation tuning, etc.), provided that the distortion caused to the energy spectrum is important enough.

7.1.2 Charge Non-Linearity (QNL)

The CD energy response is subject to two kinds of non-linearity, the first one is the LS response non-linearity, where the LS photo-production is not linear with the deposited energy as illustrated in figure 2.12a. The LS response is composed of physical non-linearity. Particle interactions in the LS will produce mainly scintillation light, as discussed in section 2.3.2, but will also produce some Cherenkov light (< 10% of the collected light). Both mechanisms possess intrinsic non-linearity, for the Cherenkov emission it depends on the velocity of charged particle velocity while the scintillation photon-yield follows a so-called Birk's law with a "quenching" effect depending on the energy and type of particle [34]. This result in a event-wise non-linearity.

The second type of non-linearity comes from the LPMT charge measurements. When photons hit a PMT and give rise to PEs, a current pulse is formed. In the photon counting regime, simply exceeding a certain threshold allows to conclude that a single photon hit the PMT. When several photons hit the PMT simultaneously, one enters the photon integration regime : the pulse is sampled and integrated over a certain time window to produce a reconstructed charge Q. Calibration methods are applied to determine the relationship between the charge Q and the number of PEs (which is the quantity proportional to the energy deposit one wants to measure). Several effects impact this procedure: the signal pulse can fluctuate and be distorted between two events where the same number PEs occurred; the PMT gain might not be linear as a function of the number of photons that hit the PMT; the charge reconstruction algorithm is not supposed to be perfect, and its results are further affected by electronic noise and inter-channel cross-talk. The impact of these effects grows with the number of PEs.

Precedent studies [24] suggest a model for the channel-wise QNL:

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

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

Studies at previous experiments, like Daya Bay, concluded that the best reachable control of QNL in the 1-10 PEs range was $\gamma_{qnl} = 0.01$ [84]. As already mentionned in section 2.3.2, JUNO LPMTs operate in a larger range : 1-100 PEs (See also table 7.1). In such a case, a realistic value of γ_{qnl} is not known.

	1PE	2~5PE	5~10PE	10~20PE	20~50PE	50~100PE	>100PE
LPMT	42.56%	40.54%	8.74%	5.12%	2.80%	0.24%	0.003%
SPMT	95.19%	4.80%	0.01%	0%	0%	0%	0%

TABLE 7.1 – The charge fraction in terms of the number of PE collected at the single PMT for the reactor $\bar{\nu}_e$ IBD events. Table taken from [24]

The event-wise impact resulting from the channel-wise QNL can be parameterised this way :

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

In JUNO, the visible energy is proportional to the number of emitted photons per unit energy deposit. It includes the physical non linearities. In the equation above E_{vis}^{true} is this visible energy, while E_{vis}^{rec} is what it becomes when the reconstructed charges found in an event are modified according to Eq. 7.1.

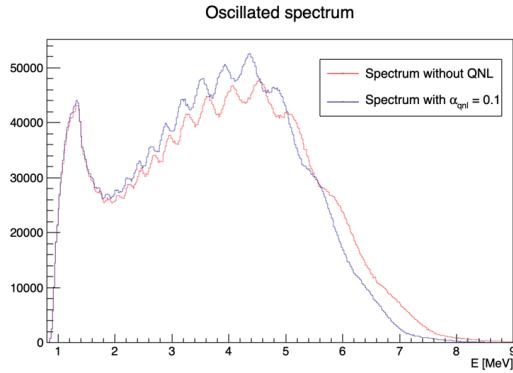


FIGURE 7.3 – Two oscillated spectra of $1e7$ event expected in JUNO. In red the spectrum without supplementary QNL. In blue the same spectrum but where an event-wise QNL $\alpha_{qnl} = 10\%$ is introduced.

An example is shown on Fig. 2.14, where we show the $E_{vis}^{rec}/E_{vis}^{true}$ ratio for several samples of uniformly distributed electron events, generated with various values of E_{vis}^{true} . Here, an extreme value $\gamma_{qnl} = 0.05$ was assumed. On can see on Fig. 2.14 that it corresponds to a 2% effect at 8 MeV, equivalent to $\alpha_{qnl} = 0.025$.

This example is from references [24], which aimed at demonstrating the potential of the dual calorimetry calibration method mentioned in section 2.4.3. If it works as hoped, the residual event-wise QNL effect will be below 0.3%. In this chapter, we propose methods to detect residuals higher than this.

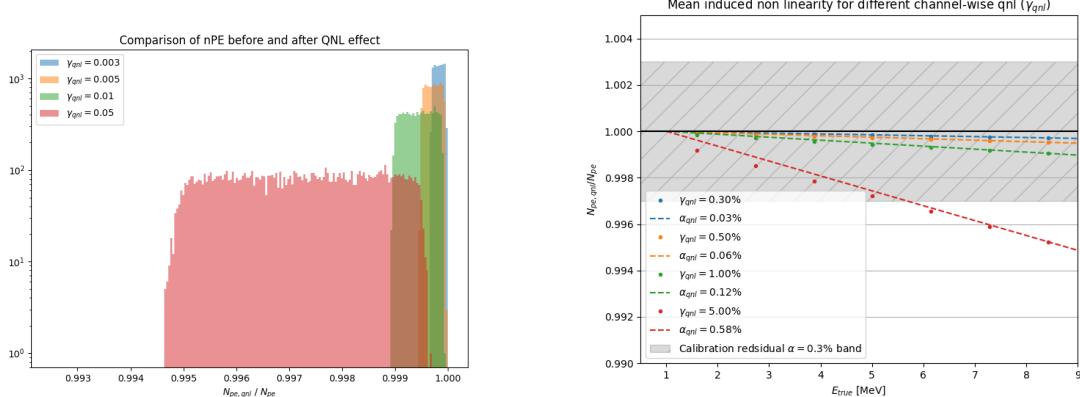
Fig. 7.4b show several other examples with varying γ_{qnl} values, and the corresponding values of α_{qnl} . Using 1M events from the JUNO official simulation J23.0.1-rc8.dc1 (released on 7th January 2024), we simulated events up to the photon collection in LPMTs and introduced an additional channel-wise QNL by using the equation 7.1 to modify the number of collected photons.

In figure 7.4a we show the distribution of the ratio $\frac{Q_{rec}}{Q_{true}}$ for central events ($R < 4m$) and different values of γ_{qnl} . In figure 7.4a, we show the mean of this distribution as a function of the energy. We also present the effective α_{qnl} for each value of γ_{qnl} . We observe that using the event-wise QNL is equivalent to the mean behavior of using channel-wise QNL.

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

7.2 Approach

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



(A) Distribution of ratio of collected nPE after the additional QNL over the number of nPE that would be collected for different γ_{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 7.4a 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 [26].

FIGURE 7.4

2430 7.2.1 Data production

2431 IBD spectra

2432 The first step involves generating the data on which our tools will be tested. In this study we
 2433 use Monte-Carlo toys. For each toy we generate a $\bar{\nu}_e$ energy spectrum from the Taishan, Yangjiang
 2434 and Dayabay nuclear power plants, the reactors used as source for the NMO analysis. The reactors
 2435 parameters comes from JUNO official database, which shared among all physics analysis, the JUNO
 2436 common inputs. This provides the initial spectra for the LPMT and SPMT systems. We then incorpo-
 2437 rate physic effects such as the LS non-linearity etc... (more details in section 7.3.1). Finally, we apply
 2438 the reconstruction resolution for each system to their respective spectra, resulting in the final LPMT
 2439 and SPMT spectra.

2440 We will study the effect of exposure on our methods at different threshold: 100 days, 1 year, 2 year
 2441 and finally 6 years which is the nominal data taking period for the NMO analysis.

2442 These spectra are generated for different QNL, $\alpha_{qnl} = 0$ (no spectrum distortion) and for $\alpha_{qnl} \in$
 2443 $\{0.01, 0.005, 0.003, 0.002, 0.001\}$. As a reminder, the calibration guarantees a residual event-wise non-
 2444 linearity of $\alpha_{qnl} \leq 0.003$ [26].

2445 The first test does not require any fitting, we are just comparing the LPMT and SPMT spectra using
 2446 the expected statistical correlation matrix in the case $\alpha_{qnl} = 0$. For details about the generation of this
 2447 correlation matrix, refer to section 7.5.2. This test is the spectrum χ^2 or χ^2_{spe} . In this test we compute

a χ^2 representing the compatibility between the LPMT and SPMT spectra:

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

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

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

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 the covariance matrix of the LPMT + SPMT spectra. A is a transformation matrix defined as:

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

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

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

Background spectra

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

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

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

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 χ^2 representing the compatibility between the measurements of the systems. Because the SPMT system is not sensible to the oscillation parameters Δm_{31}^2 and θ_{13} , the test is only done on the oscillation parameters θ_{12} and Δm_{21}^2 . We can thus produce the individual chi square χ_{ind}^2

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

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

$$U = AVA^T \quad (7.10)$$

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

where λ_L and λ_S are the measured parameters by the LPMT and SPMT systems respectively. The different λ considered are θ_{12} and Δm_{21}^2 . V here is the 4×4 covariance matrix between the parameters

$\sin^2(2\theta_{12})$	Δm_{21}^2	Δm_{31}^2	$\sin^2(2\theta_{13})$
$0.851^{+0.020}_{-0.018}$	$7.53 \pm 0.18 \times 10^{-5} \text{ eV}^2$	$2.5283 \pm 0.034 \times 10^{-3} \text{ eV}^2$	0.8523 ± 0.00268

TABLE 7.2 – Nominal PDG2020 value [34]. All value are reported assuming Normal Ordering.

²⁴⁶⁴ $\theta_{12,L}$, $\Delta m_{21,L}^2$, $\theta_{12,S}$ and $\Delta m_{21,S}^2$. A is the transformation matrix that allow us to compute the covariance
²⁴⁶⁵ matrix de $\vec{\Delta}$ from V following

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

²⁴⁶⁶ Same as described above, by comparing the distribution of this χ^2_{ind} when $\alpha_{qnl} = 0$ and $\alpha_{qnl} \neq 0$ we
²⁴⁶⁷ 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 as the same time. We minimize a χ^2_{joint} defined over the two spectra, the LPMT and SPMT one

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

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

²⁴⁷⁰ where D_i is the content of the i th bin measured, from the data, and T_i is the theoretical number of
²⁴⁷¹ event in this bin. V is the covariance matrix of our spectrum.

²⁴⁷² T is the fitted function and depend on multiple parameters

- ²⁴⁷³ — The oscillation parameters θ_{12} , Δm_{21}^2 , θ_{13} and Δm_{31}^2 . Those parameters can be free, have a pull
²⁴⁷⁴ term or be fixed during the fit.
- ²⁴⁷⁵ — We take into account in the data production the matter effect and parametrize it by the pa-
²⁴⁷⁶ rameter ρ , the effective rock density between the reactors and the experiment. Same as the
²⁴⁷⁷ oscillation parameters, this parameter can be free, pulled or fixed.
- ²⁴⁷⁸ — The exposure of the considered data which is just a normalization factor in front of the theo-
²⁴⁷⁹ retical spectrum. This parameter is fixed at the start of the fit.

²⁴⁸⁰ In the standard joint fit, the free parameters are $\sin^2(2\theta_{12})$, Δm_{21}^2 and Δm_{31}^2 . $\sin^2(2\theta_{13})$ is fixed to the
²⁴⁸¹ PDG nominal value. For simplicity, we refer to $\sin^2(2\theta_{12})$ and $\sin^2(2\theta_{13})$ as θ_{12} and θ_{13} respectively.

²⁴⁸² Both of the LPMT and SPMT systems are sensitive to θ_{12} and Δm_{21}^2 , thus these parameters are totally
²⁴⁸³ free and start at the PDG nominal value. Only the LPMT system is sensitive to Δm_{31}^2 , we let it
²⁴⁸⁴ free so we can observe the effect of the deformation on it while the solar parameters θ_{12} , Δm_{21}^2 are
²⁴⁸⁵ constrained by the SPMT system. To prevent Δm_{31}^2 to take absurd value, we add a pull term using
²⁴⁸⁶ the PDG nominal value and errors. The PDG nominal values used in this study can be found in table
²⁴⁸⁷ 7.2.

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

²⁴⁸⁸ θ_{13} is the parameter on which we are least accurate. It's fixed to nominal value to prevent degeneracy
²⁴⁸⁹ (table 7.2).

²⁴⁹⁰ The covariance matrix is produced from a correlation matrix C

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

²⁴⁹¹ where σ_i is the uncertainty on the number of event in the i th bin. We consider in this study that the
²⁴⁹² content of each bin follow a Poisson statistic, thus the uncertainty is $\sigma_i = \sqrt{N_i}$ where N_i is the content
²⁴⁹³ of the i th bin. The bin content used for the uncertainty can come from two sources: the data and the
²⁴⁹⁴ theoretical spectra $\sigma_i = \sqrt{D_i}$ (Pearson test) and $\sigma_i = \sqrt{T_i}$ (Neyman test). Precedent studies have
²⁴⁹⁵ show that both Pearson and Neyman tests show bias at low statistic, we thus use the Pearson V test
²⁴⁹⁶ where

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

²⁴⁹⁷ and the covariance matrix V is computed using the data spectrum for the uncertainty.

²⁴⁹⁸ The estimation of the covariance is crucial in this study as the strength of this test rely on the sys-
²⁴⁹⁹ tematic and statistical correlations between the LPMT and SPMT spectrum. The generation methods
²⁵⁰⁰ and results of this matrix is detailed in section 7.5.2.

²⁵⁰¹ Delta joint fit

²⁵⁰² Using the same structure we define a second joint fit, the Delta joint fit where, in addition to every-
²⁵⁰³ thing that was discussed above, we add two other parameters $\delta\theta_{12}$ and $\delta\Delta m_{21}^2$ and split the theoretical
²⁵⁰⁴ $T(\theta_{12}, \Delta m_{21}^2, \dots)$ spectrum in two

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

²⁵⁰⁵ If the there is no additional distortion between the LPMT and the SPMT spectra, the fit should
²⁵⁰⁶ converge to $\delta\theta_{12} = \delta\Delta m_{21}^2 = 0$. By observing the dispersion of those parameters we can define
²⁵⁰⁷ 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 $\alpha_{qnl} \neq 0$
²⁵⁰⁸ to define a p-value.

²⁵⁰⁹ The last test we explore in this thesis is to fit the same spectrum with the Standard Joint fit, that
²⁵¹⁰ we consider as the hypothesis without distortion H_0 , and the Delta Joint fit, designated as the H_1
²⁵¹¹ hypothesis. By looking at the dispersion of $\chi^2_{joint, H_0} - \chi^2_{joint, H_1}$ we can extract a sensitivity to potential
²⁵¹² distortion.

²⁵¹³ 7.2.4 Data and theoretical spectrum generation

²⁵¹⁴ To implement the joint fit, we have technically two data spectra and two theoretical spectra. The data
²⁵¹⁵ in this study are produced using an IBD generator *IBD gen*, see section 7.3.1. The theoretical spectrum
²⁵¹⁶ are produced the same way as data spectrum but with much higher statistics, 10^7 events to compare
²⁵¹⁷ with the $\approx 10^5$ events for 6 years statistic. The two spectrum, that we get as a collection of events,
²⁵¹⁸ are binned in two histograms from 0.8 to 9 MeV of reconstructed energy with bins of 0.02 MeV each,
²⁵¹⁹ resulting in 410 bins per spectrum. An illustration of the theoretical spectrum can be found in figure
²⁵²⁰ 7.5. The low number of events in the tail of the spectrum can cause instability due to the low statistic,
²⁵²¹ we thus cut the spectrum at 7.5 MeV / 335 bins for the fit.

²⁵²² All the IBD spectra presented and used in this study are produced assuming Normal Ordering using
²⁵²³ the PDG nominal value [34] for the oscillation parameters. Those values are reported in table 7.2.

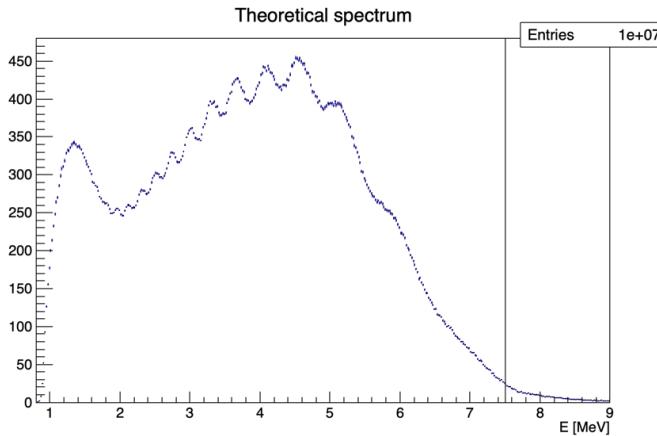


FIGURE 7.5 – Theoretical LPMT spectrum at nominal oscillation values binned using 410 bins from 0.8 to 9 MeV. It is rescaled to 6 years statistic. The black line represent the 335 bin cut

2524 7.2.5 Limitations

2525 In this work we are only working considering the statistical errors. We can ignore systematic effects,
 2526 such as effects that would affect the neutrino spectrum or the background spectrum, as they are
 2527 entirely correlated between the two systems. The details of those systematic effects can be found in
 2528 [3].

2529 Most of our results assume decorrelated detection effects between the SPMT and LPMT systems.
 2530 Their respective reconstruction effects are simulated using simple gaussian drawing on the resolution,
 2531 independently from the event position. This approach was used in previous sensitivity and
 2532 precision studies [3, 32]. The potential effect of those reconstruction effects and a first attempt to take
 2533 them into account are explored in section 7.5.2.

2534 Even if the goal of this work is to propose deformation agnostic tools, the QNL we use in this study is
 2535 simplistic as we consider event-wise, position uniform deformation. We show in figure 7.4a and 7.4b
 2536 that event-wise QNL is equivalent to the mean behaviour of channel-wise QNL but a more complete
 2537 study would simulate channel-wise deformation for each event.

2538 7.3 Fit software

2539 In this section, I describe the ft framework that was used in this study. The software is composed
 2540 of two parts as illustrated in figure 7.6: A standalone part composed of ROOT [85] macros, and the
 2541 Avenue framework.

2542 The Avenue framework is responsible for the spectrum and configuration reading, transforming
 2543 the raw collection of events into spectra, managing the physics effect such as the oscillation and
 2544 computing and minimizing the χ^2 with the help of the RooFit library. The macros are invoking, if
 2545 necessary, the Avenue framework and are the entry point for fitting, generating the necessary inputs
 2546 quantity such as the spectra and correlation matrix, analysing the fit results and managing jobs for
 2547 distributed computing.

2548 In this section we will focus on the IBD generator in section 7.3.1 and the fit macro in itself in section
 2549 7.3.2.

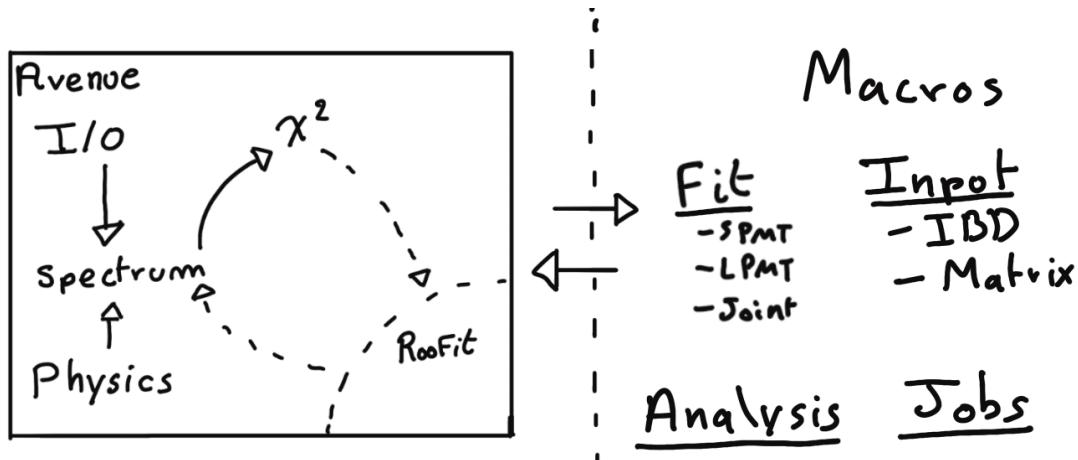


FIGURE 7.6 – Schematic description of the fit framework

2550 7.3.1 IBD generator

2551 The IBD generator is a standalone generator used to produce oscillated and non oscillated spectra
 2552 as the one seen by the JUNO experiment. It takes as inputs physics parameters and a collection
 2553 of histograms, values and function provided by JUNO to its analysis groups, referred as the JUNO
 2554 common inputs.

2555 Options allow to enable or disable effects such as non-uniformity and non-linearity. It finally take as
 2556 an argument the number of events to generate N_{evt} . Optionally, we generate an effective number of
 2557 events N by drawing in a Poisson distribution of mean N_{evt} .

2558 Then for each event we

- 2559 1. Choose randomly, following the reactor power fraction, the source reactor of the neutrino.
- 2560 2. Generate a random interaction position in the detector following a uniform distribution over
the detector volume.
- 2562 3. Draw a random neutrino energy E_ν from the expected neutrino emission spectrum of every
reactor. This spectrum is computed by:
 - 2564 (a) Computing the power spectrum of each isotopes ^{235}U , ^{238}U , ^{239}Pu , ^{241}Pu using the Huber-
Mueller model [5, 8].
 - 2566 (b) Summing the contribution of each isotopes following the respective fission fraction [0.58,
0.07, 0.30, 0.05] as reported in [86].
 - 2568 (c) The power of each reactor is then adjusted by their distances from the detector, the detector
efficiency and their mean duty cycle (11 of 12 month).
 - 2570 (d) The total spectrum is then finally adjusted by taking into account the correction of the Day
Bay bump [11], adjustment due to spent nuclear fuel and due to the non-equilibrium.
- 2572 4. (Optional) Compute the survival probability due to oscillation at nominal oscillation param-
eters value. If the neutrino does not survive, the event is rejected and the algorithm restart
from step (1).
- 2575 5. Compute the emitted positron energy E_{pos} from the mass difference. If the neutrino does not
have enough energy reject the event and start from step (1).
- 2577 6. Compute the deposited energy E_{dep} by incrementing E_{pos} by 511 keV to account for the positron
annihilation. We do not consider cases where some of the energy leak outside of the detector
(positron or annihilation gammas escaping the CD).

- 2580 7. Correct the deposited energy with the expected event-wise non-linearity from [26] to obtain
 2581 the visible energy E_{vis} .
- 2582 8. (Optional) Add a custom non-linearity as described in section 7.1.2. This non linearity is
 2583 characterized by α_{qnl} to obtain E_α .
- 2584 9. Finally, using the expected resolution of the LPMT and SPMT systems, provided in the JUNO
 2585 common inputs, we draw from a gaussian characterized by those resolution the reconstructed
 2586 energy E_{rec} or E_{lpmt} and E_{spmt} for each systems. The resolutions are provided as ABC param-
 2587 eters using

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

2588 where A is the term driven by the Poisson statistics of the total number of detected photoelec-
 2589 trons, C is dominated by the PMT dark noise, and B is dominated by the detector's spatial
 2590 non-uniformity. The relative and absolute resolutions of the LPMT and SPMT systems are
 2591 illustrated in figure 7.7.

2592 The events are stored as n-tuples and are not yet binned at the end of the generator.

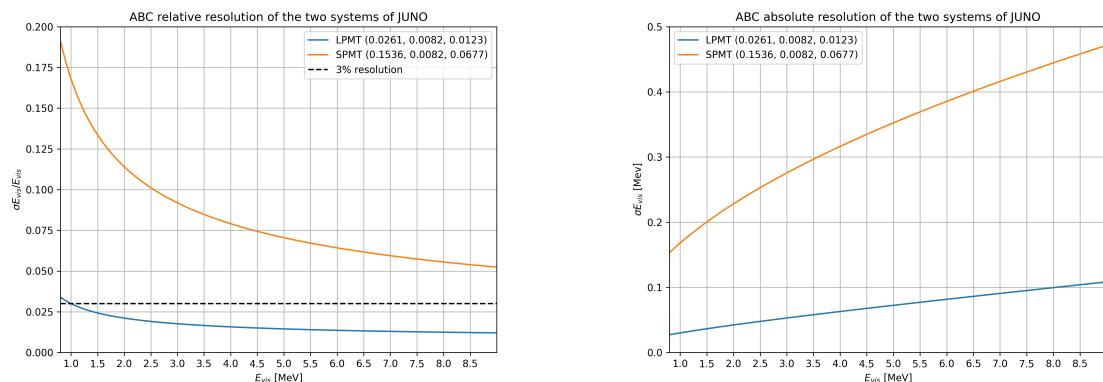


FIGURE 7.7 – Relative (On the left) and absolute (On the right) resolutions of the LPMT and SPMT systems used in this study. The number in parenthesis are the parameter A, B and C respectively for each systems.

2593 7.3.2 Fit

2594 The fit macro is the core of this fitting procedure. This macro is responsible for loading the fit
 2595 configuration and setup the Avenue framework. Using Avenue, it will setup the data files, theoretical
 2596 spectrum, choose the binning, χ^2 , etc... It also have the possibility to generate toys on the fly based
 2597 on the theoretical spectrum. Given this theoretical spectrum we can randomize the bin content either
 2598 by:

- 2599 1. Drawing the bin content in a Poisson distribution with the bin content as parameter.
 2600 2. Drawing the bin content in a Gaussian distribution with the bin content as mean and variance.
 2601 The bin content is then rounded to the nearest integer.
 3. Drawing the bin difference following a given covariance matrix using the Choleski decomposi-
 2602 tion. This matrix is at least the statistical covariance matrix but can also contain systematic

uncertainties.

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

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

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

$$(7.23)$$

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

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

7.4 Technical challenges and development

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

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

While the IBD generator was already able to generate LPMT and SPMT spectrum, it was not designed for generating correlated spectrum. As detailed in section 7.3.1, up to the reconstruction effect, the two spectrum need to share the same generation else the two spectrum would be decorrelated and it would be like we would run two different experiment.

7.5 Results

7.5.1 Validation

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

- 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 [3].

- The standard joint fit with an independent covariance matrix (*Indep Standard joint*), meaning that the covariance between the LPMT and SPMT spectra is 0, believe to have twice as much informations, and thus believe to have a grater precision than the individual fits.
- The standard joint (*Standard joint*) fit with a correlated covariance matrix has errors similar to the LPMT individual fit as the LPMT drive the precision on θ_{13} and Δm_{31}^2 and that the LPMT as SPMT are expected to have close precision on θ_{12} and Δm_{21}^2 .
- The delta joint (*Delta joint*) fit with covariance matrix have the same resolution as the standard joint fit. The supplementary parameter $\delta\theta_{12}$ and $\delta\Delta m_{21}^2$ should not bring supplementary precision.

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

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

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

Asimov studies

We ran Asimov studies on the tests presented above on the updated framework, the results are reported in table 7.3. All those test are ran considering statistics error only, 6 years exposure with all backgrounds, Pearson χ^2 (covariance is estimated using data spectrum) and θ_{13} fixed to nominal value. For the SPMT fit Δm_{31}^2 is fixed at nominal value as the SPMT system is net expected to be sensitive to this parameter.

	Δm_{31}^2 error	$\delta\Delta m_{21}^2$ error	θ_{12} error	$\delta\theta_{12}$ error	Δm_{31}^2 error	χ^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 Δm_{31}^2 and Δm_{21}^2						
Indep Standard joint			9.33082e-04			4.82955e-26
LPMT			1.27032e-03			2.58849e-26
SMPT			1.31070e-03			2.24106e-26
Weighted			9.12193e-04			
Fixed Δm_{31}^2 and θ_{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 7.3 – Results of the Asimov studies on the updated framework. All results are Asimov fit, considering 6 years exposure, θ_{13} is fixed to nominal value, χ^2 is pearson meaning that he error is estimated using the data spectrum

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

We observe, as expected, that $\sigma_{Weighted} \approx \sigma_{Indep Standard joint}$ with the exception of $\sigma\theta_{12}$. This could from the slight difference in statistic between the SPMT and LPMT spectra. Indeed, due to a larger smearing in energy resolution, events that would be inside the spectrum range [0.8, 7.5] MeV are

smeared outside it. This deficit is partially compensated by event outside the spectrum coming back in it but we expect very few event outside the spectrum in comparison to event at the edges of it. Thus the event deficit is not totally compensated. θ_{12} being mainly driven by the amplitude of the spectrum (see illustration 2.2), that's why we think this the origin of the difference.

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

Finally for the *Delta* fit, the error on $\delta\theta_{12}$ and $\delta\Delta m^2_{21}$ are of the same order of magnitude than the errors on θ_{12} and Δm^2_{21} 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 χ^2 (Eq. 7.13) without pull parameter. We show in figure 7.8 the effect of using a simple Pearson χ^2 . We see that $\sin^2(2\theta_{12})$ (reported as θ_{12} for simplicity) is biased of about 0.5σ and Δm^2_{21} biased of about 0.1σ . When introducing the PearsonV χ^2 (Eq. 7.17) the bias disappear as reported in figure 7.9.

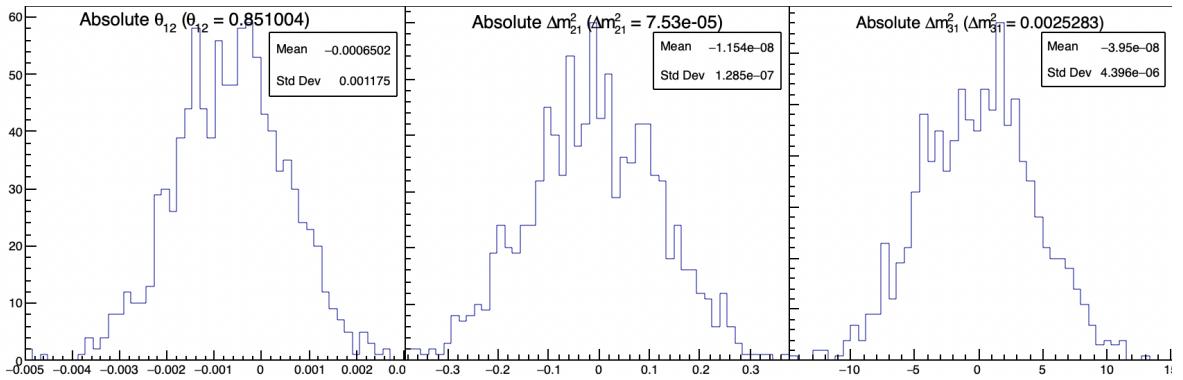


FIGURE 7.8 – Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, Pearson χ^2 , θ_{13} fixed.

When the supplementary parameters are introduced in the Delta Joint fit, the fit is stable as shown in the results figure 7.10. 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 δ 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.

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

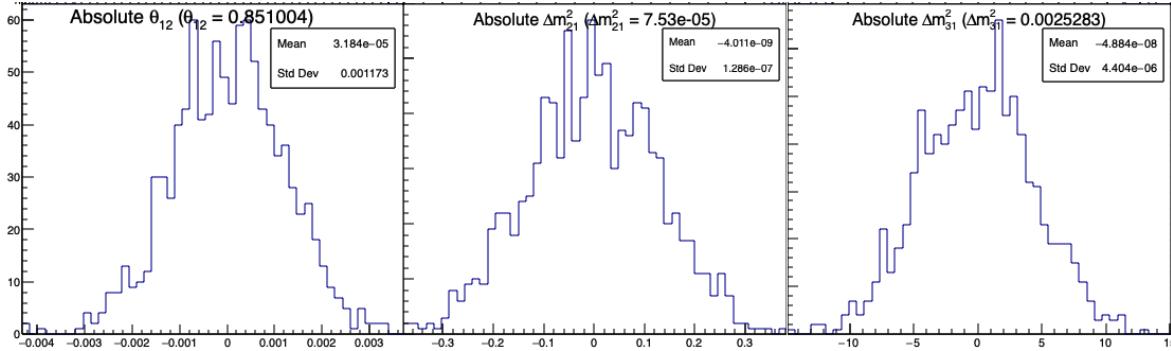


FIGURE 7.9 – Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, PearsonV χ^2 , θ_{13} fixed.

2692 calibration ($\alpha_{qnl} = 0.3\%$ [26]). The background had to be removed as JUNO provide them already
 2693 smeared, thus the introduction of supplementary QNL is not trivial, the resolution being dependent
 2694 of E_{vis} which is affected by the QNL. We use a covariance matrix assuming no QNL. The effect of this
 2695 QNL on the spectrum is illustrated in figure 7.11. In table 7.4 we report the results of the different
 2696 scenarios.

Mean (std dev)	$\theta_{12} [10^{-3}]$	$\Delta m^2_{21} [10^{-7}\text{eV}^2]$	$\Delta m^2_{31} [10^{-6}\text{eV}^2]$	$\delta\theta_{12} [10^{-3}]$	$\delta\Delta m^2_{21} [10^{-7}\text{eV}^2]$
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 7.4 – 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 Δm^2_{31} is fixed at nominal value. The χ^2 is PearsonV. The correlation matrix used to fit assume no QNL in the spectrum.

2697 The results in table 7.4 are subtracted from their nominal value, themselves reported in table 7.2.
 2698 We clearly see the bias induced by $\alpha_{qnl} = 1\%$ when comparing the SPMT and LPMT results. The
 2699 Indep Standard is, as expected, the mean value between the SPMT and LPMT: the fit having no
 2700 informations about the correlation between the spectrum think it have two uncorrelated experiments
 2701 thus report an in between value. When introducing the relationship between the LPMT and SPMT
 2702 spectra in the Standard fit, the joint fit cannot find a clean minima, it thus converge to a completely
 2703 incorrect value.

2704 Introducing the δ without the correlation in Delta Indep remove the bias and converge to the SPMT
 2705 minima, the δ absorbing the deformation of the LPMT spectra.

2706 Finally, with the δ and the covariance matrix, θ_{12} is unbiased, $\delta\theta_{12}$ absorbing the deformation. $\delta\Delta m^2_{21}$
 2707 is still heavily biased, even more than LPMT only, for the same reason than the Standard fit: the
 2708 correlation make it difficult to converge to the nominal value.

2709 Overall Δm^2_{31} bias is unchanged as the SPMT spectrum bring no information about the parameter.
 2710 The δ are significant, naively up to 7.46σ for $\delta\theta_{12}$ in the Delta fit.

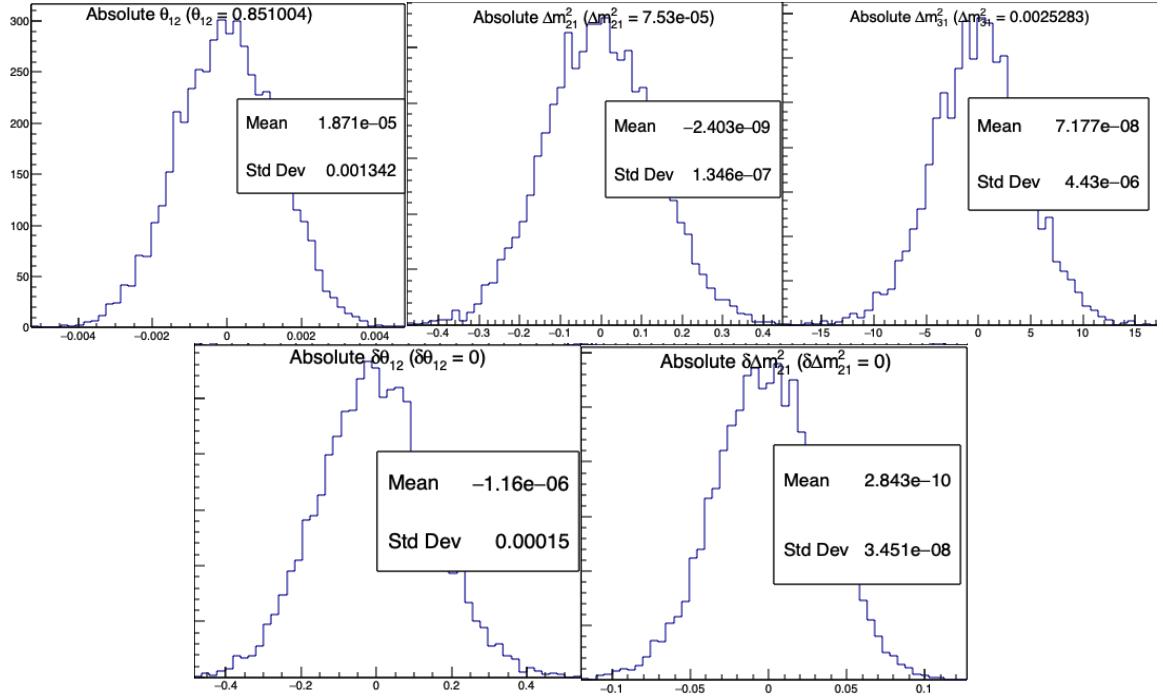


FIGURE 7.10 – Distribution of BFP - nominal value for 5000 toy Delta joint fit. 6 years exposure, all background, PearsonV χ^2 , θ_{13} fixed.

7.5.2 Covariance matrix

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

Analytical method

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

The first step is to compute the statistical uncertainty of the input spectrum while taking into account the smearing, considering no uncertainty on the smearing. For this, using the notation of section 39.2.5 *Propagation of errors* of PDG2020 [34] and considering an extended spectrum of 820 bins following the binning scheme introduced in 7.2.4, the first 410 for the LPMT and the last 410, we consider

- $\boldsymbol{\theta} = (\theta_0, \dots, \theta_n); n = 820$ the content of the spectrum bins.
- $\boldsymbol{\eta}(\boldsymbol{\theta}) = (\eta_0(\boldsymbol{\theta}), \dots, \eta_m(\boldsymbol{\theta})); m = 820$ the set of smearing functions representing the PMT resolutions.

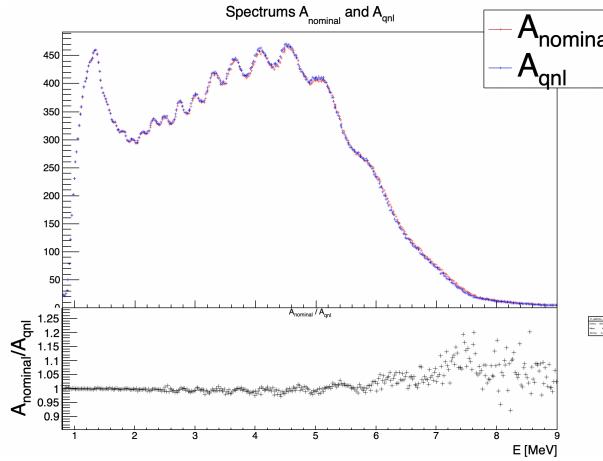


FIGURE 7.11 – **Top:** Theoretical spectrum without QNL (in red) and with $\alpha_{qnl} = 1\%$ (in blue). **Bottom:** Ratio between the theoretical spectrum with and without QNL.

²⁷³¹ η_m can thus be defined as

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

²⁷³² where $G(i, \sigma(E_i))(j)$ is the smearing function defined as

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

²⁷³³ 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 ²⁷³⁴ the i th bin respectively.

²⁷³⁵ We can then construct the transfer matrix A as

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

²⁷³⁶ and then compute the first part of our covariance matrix

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

²⁷³⁷ where V is the uncorrelated covariance matrix simply defined, under the assumption of poissonian ²⁷³⁸ statistic for the bin content,

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

²⁷³⁹ Now we just need to consider the uncertainty on the smearing $\sigma \eta_i$, considering no uncertainty on ²⁷⁴⁰ the unsmeared spectrum. From Eq. 7.25, the $G(i, j) \equiv G(i, \sigma(E_i))(j)$ are considered independents ²⁷⁴¹ from each other $\forall i, j$. This mean that this covariance matrix is diagonal, we only need $\sigma G(i, j)$. We ²⁷⁴² can derive this term from two equation:

- ²⁷⁴³ — The term $G(i, j) \theta_j$ represent the number of event smeared from the bin j that end up in the bin ²⁷⁴⁴ i . This is a number, we thus assume poissonian statistic so that $\sigma[G(i, j) \theta_j] = \sqrt{G(i, j) \theta_j}$.
- ²⁷⁴⁵ — Using basic error propagation we can say that $\sigma^2[G(i, j) \theta_j] = \theta_j^2 \sigma^2 G(i, j) + G(i, j)^2 \sigma^2 \theta_j$.

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

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

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

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

By summing the two covariance matrix, we can extract a correlation matrix presented in figure 7.12. The correlation between the SPMT and LPMT spectra is greater at the start of the spectrum, where the absolute smearing is the smallest, up to 5% correlation, and diffuse as the bins are further from each other and the absolute resolution grow.

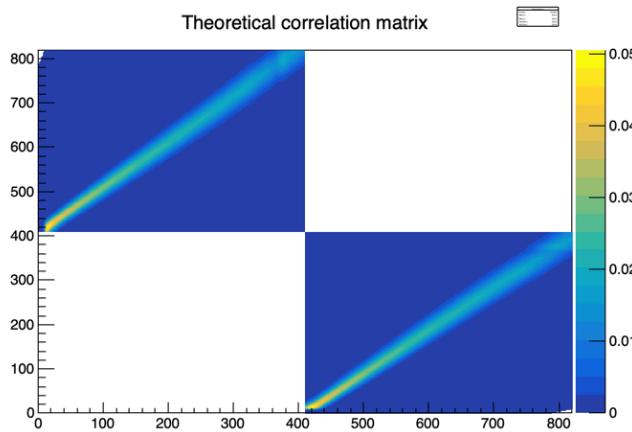


FIGURE 7.12 – Theoretical correlation matrix between the LPMT spectrum (bins 0-409) and the SPMT spectrum (410-819). The diagonal has been set to 0 (it was 1) for readability purpose.

Empiric method

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

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

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

Due to the nature of our estimator, the estimated correlation coefficient is subject to statistical fluctuation as any estimator. There is no definite formula to compute the standard deviation of the correlation coefficient as suggested in this study [87] but all cited formula depend solely on the 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 7.13 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 7.13a that if the toy number N is too low, the statistical noise make the correlation pattern almost completely disappear, in figure 7.13b we see clearly the same correlation pattern as in the theoretical matrix in figure 7.12. On the final matrix in figure 7.13c 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.

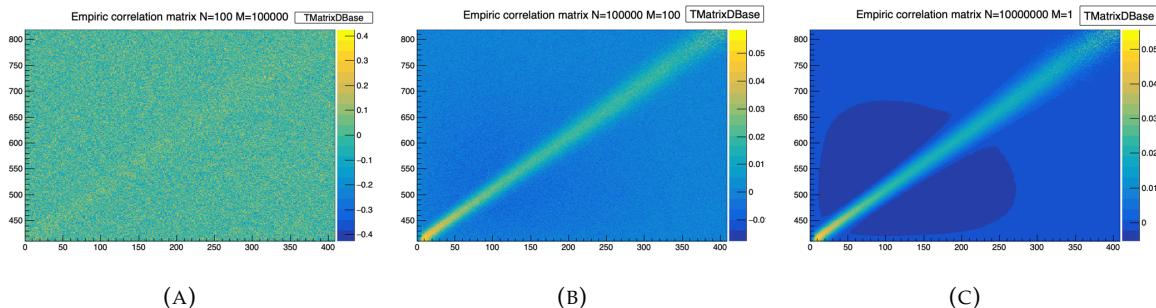


FIGURE 7.13 – Upper left corner of the estimated correlation matrix between the LPMT and SPMT spectrum for different configuration of N toy with different number of M events per toy

The difference between the element of the theoretical and the empiric correlation matrices are presented in figure 7.14a. We see 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 7.14b, 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 χ^2 calculation. The framework use the Cholesky decomposition [88] 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 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 that was actually deposited.

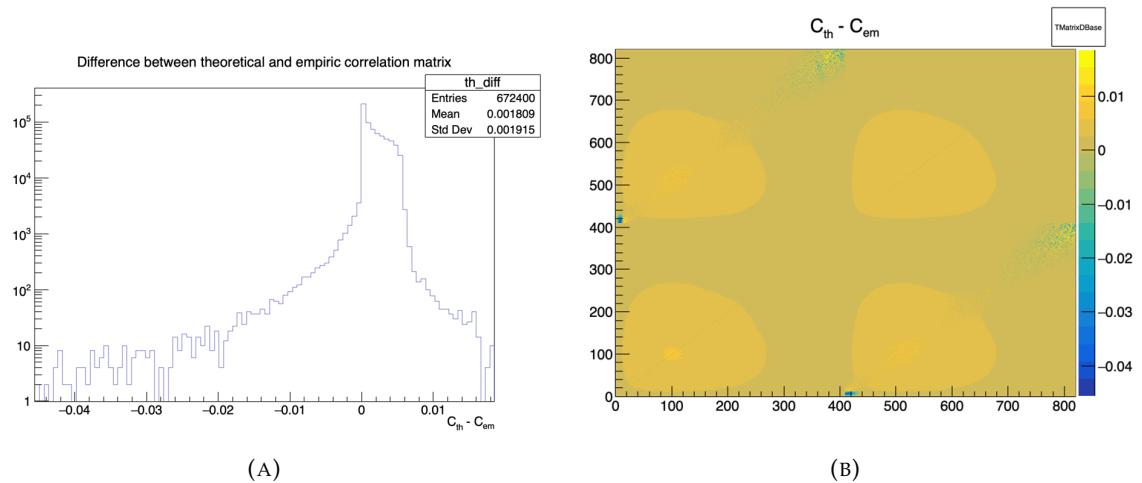


FIGURE 7.14 – Difference between the element of the theoretical and empiric correlation matrix

2794 On a more subtle scale, the randomness in the production of scintillation photons is common for the
 2795 two systems, if the liquid scintillator produces fewer scintillation photons for an event, both systems
 2796 are likely to underestimate the energy.

2797 We study those effects by computing from a dataset of IBD events, uniformly distributed in the CD,
 2798 the correlation between the reconstruction errors on the energy

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

2799 where E_{lpmt} and E_{spmt} are the reconstructed energies from both systems and E_{dep} is the deposited
 2800 energy in the detector.

2801 With this observable, the bias difference between the two reconstructions at fixed R and E is irrele-
 2802 vant. However, since we compute the correlation in E and R^3 bins, we need to account for the
 2803 potential spurious relationship between the errors and their respective biases. If the bias is small
 2804 relative to the resolution, it can be ignored; but if the bias variation is on the same order of magnitude
 2805 as the error, it may introduce false correlations. For this reason, based on the CNN results shown in
 2806 figure 4.8, we restrict our analysis to the $1 < E_{dep} < 9$ MeV range.

2807 The results of those correlations are presented in figure 7.15 for the single energy and radius depen-
 2808 dency and figure 7.16 for the dual energy and radius dependency.

2809 We see correlation increase with respect to the energy which can be attributed to the signal over dark
 2810 noise ratio. As more PMTs hits come from the signal, the reconstruction becomes more signal related.
 2811 Regarding the R^3 distribution, we see almost no dependency until the total reflection area. After this
 2812 point the correlation rises as the event are exposed to the optical effect of the total reflection area.

2813 By looking at figure 7.16, we can see that the rising in correlation with respect to the energy is mostly
 2814 due to the radius dependency.

2815 The exploitation of those correlations in the fit and the data production, without generating and
 2816 reconstructing full spectra from SNIPER, is a bit more complicated. As seen in section 7.3.1, we
 2817 characterize the resolution of both systems by the ABC parameters. The correlation shown here take
 2818 into account all of the ABC terms, as they are the complete correlation between the two systems, but
 2819 the generation and the modeling this correlation needs to be very well understood as, as seen before,
 2820 the mass ordering and parameters measurements are very sensitive to even small correlations.

2821 We consider the binned approach that we used here, knowing that the CNN reconstruction was

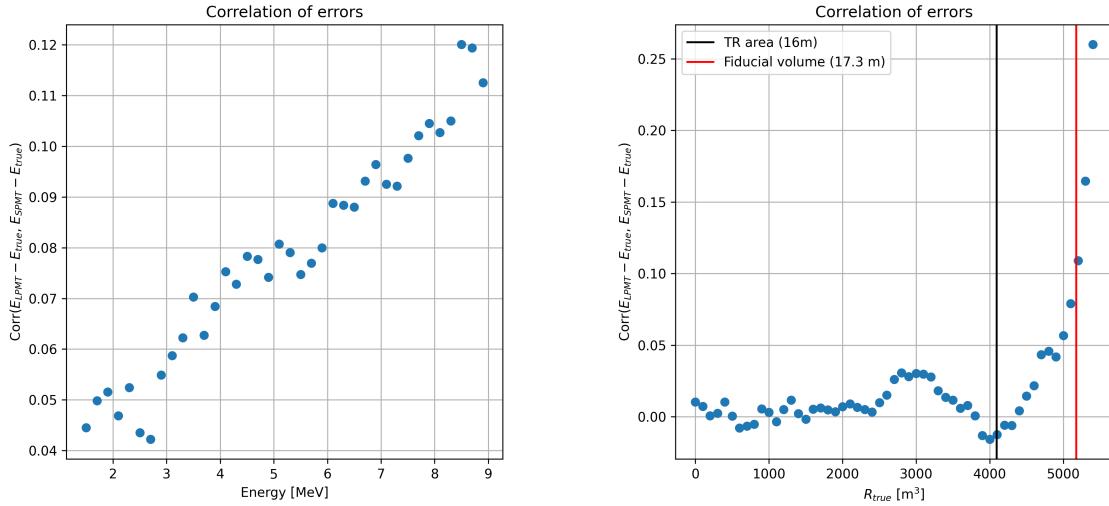


FIGURE 7.15 – Correlation on the reconstruction error between the LPMT and SPMT system as a function of (On the left) the energy, (On the right) the radius. The SPMT reconstruction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.8. To prevent effect due to the CNN bad reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.

2822 deemed efficient but flawed, to be insufficient for the complete study of those effects on the fit.

2823 7.5.3 Statistical tests

2824 In this part, I present the results of the statistical tests presented in section 7.2.

2825 Test χ_{spe}^2

2826 The χ_{spe}^2 is a chi-square representing the compatibility between the LPMT ans SPMT spectra under
2827 constraints of the correlation matrix between the two.

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

2828 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
2829 calculation of V_{spe} , see section 7.2.

2830 The results for different exposures can be found in figure 7.17. To give an idea of the significance of
2831 this test, we provide the median p-value for each test $\alpha_{qnl} \neq 0$. As expected, the power of this test
2832 rises as the exposure does. We see significant discrimination at 6 years for $\alpha_{qnl} \geq 0.3\%$ where the
2833 p-value for $\alpha_{qnl} = 3\%$ is 0.005 ± 0.0022 .

2834 This test relies solely on the estimated covariance matrix between the two spectra, requiring no
2835 fitting. As a result, it is a very lightweight test that can still provide valuable indications of potential
2836 unknown distortions between the two spectra.

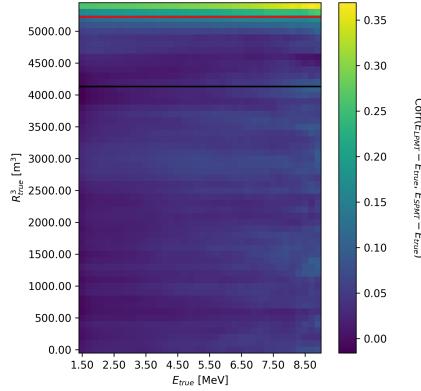


FIGURE 7.16 – Correlation on the reconstruction error between the LPMT and SPMT system as a function of the energy and the radius. The SPMT reconstruction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.8. To prevent effect due to the CNN bad reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.

2837 **Test χ_{ind}^2**

2838 The χ_{ind}^2 is the chi-square that represent the agreement between the measured oscillation parameters
 2839 θ_{12} and Δm_{21}^2 . This test is defined as

$$\chi_{ind}^2 = \Delta\lambda V_{ind} \Delta\lambda^T; \Delta\lambda = \{\theta_{12}^L - \theta_{12}^S, (\Delta m_{21}^2)^L - (\Delta m_{21}^2)^S\} \quad (7.36)$$

2840 where θ_{12}^L and $(\Delta m_{21}^2)^L$ are the oscillation parameters measured by the LPMT system. Same for θ_{12}^S
 2841 and $(\Delta m_{21}^2)^S$ for the SPMT system. We use V_{ind} computed for $\alpha_{qnl} = 0$. For more details about the
 2842 calculation of V_{ind} see section 7.2.

2843 The results are presented in figure 7.18. This test does not require any joint fit or covariance matrix
 2844 estimation between the two spectrum, it just need the estimated covariance matrix between the four
 2845 parameters. We see that the p-value are much less significant than the other tests, this is because this
 2846 test possess much less information about the relation between the LPMT and SPMT systems.

2847 This test is the most straightforward as it require only the fit of the two spectra and the estimation
 2848 of the parameters covariances, but is also the less powerful with a p value for $\alpha_{qnl} = 0.3\%$ of 0.09 ± 0.009 .

2850 **δ parameters significance**

2851 This test involves observing the values of the δ parameters in the Delta Joint fit and comparing them
 2852 tho their dispersion in the case where $\alpha_{qnl} = 0$. The results are shown in figures 7.19 and 7.20.

2853 We can see that the $\delta\Delta m_{21}^2$ has a very small discriminative power (figure 7.20) even at 6 years
 2854 exposure with a p-value of 0.34 ± 0.01 for $\alpha_{qnl} = 0.3\%$. On the other hand $\delta\theta_{12}$ (figure 7.19) has
 2855 much more discriminative power with a p-value for $\alpha_{qnl} = 0.3\%$ of 0.025 ± 0.005 . This test with a
 2856 single joint fit seems to be still less powerful than the χ_{spe}^2 . This can be explained as this method
 2857 only get information through the oscillation parameters θ_{12} and Δm_{21}^2 missing potential informations
 2858 contained in Δm_{31}^2 .

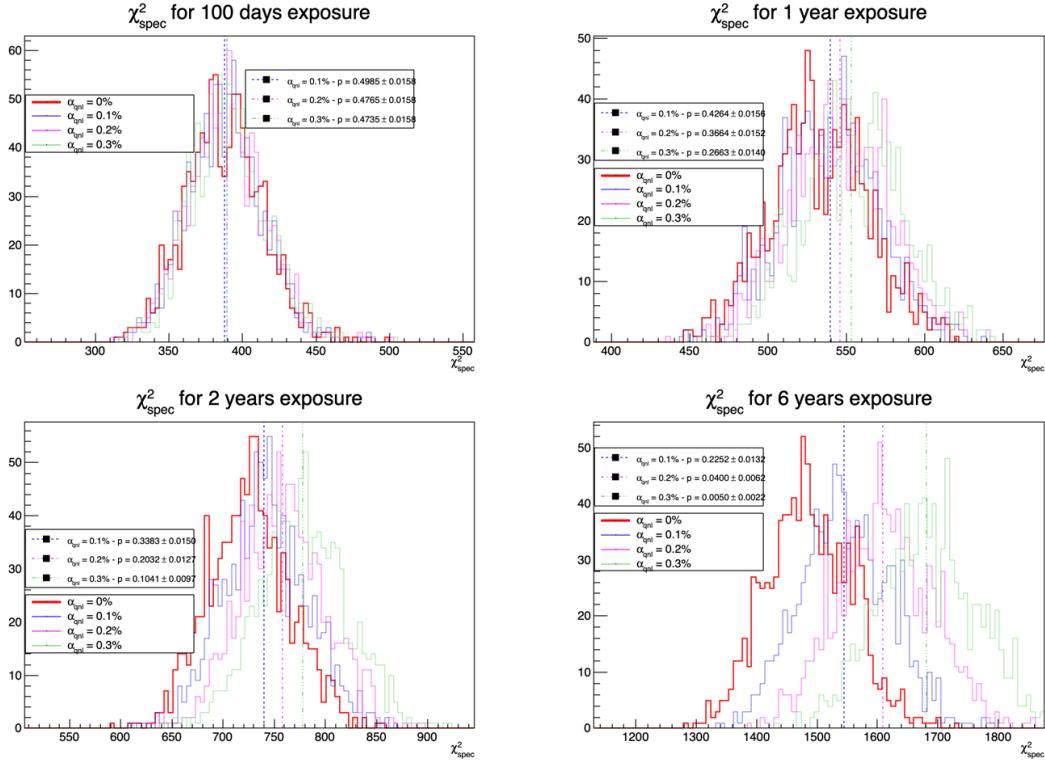


FIGURE 7.17 – Distribution of the χ^2_{spe} for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

2859 Hypothesis test

2860 In this last test we consider the two fit Standard Joint and Delta Joint as two hypothesis. The first
 2861 one, Standard Joint, is the H_0 hypothesis: we do not need supplementary parameters to describe the
 2862 energy spectrum. The second one, Delta Joint, is the H_1 hypothesis: we do need those supplementary
 2863 δ parameters to, if not correctly, approach the energy spectrum. If the δ parameter are unnecessary
 2864 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
 2865 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
 2866 exposures and α_{qnl} . The results are presented in figure 7.21.

2867 This test is the most complex, requiring two fit and the covariance matrix between the LPMT and
 2868 SPMT spectra. The results are good, close to the χ^2_{spe} , one with a p-value at 6 years for $\alpha_{qnl} = 0.3\%$ of
 2869 0.01 ± 0.003 .

2870 As explained in section 7.2.4, the spectra used for the fit are cut at 335 bins / 7.5 MeV to prevent
 2871 instability, while in χ^2_{spe} we use full 410 bins spectra. The χ^2_{spe} thus has more informations that the
 2872 hypothesis test leading to this difference in power.

2873 7.6 Conclusion and perspectives

2874 In this chapter, we present the development of a fit framework that allows us to fit multiple spectra
 2875 simultaneously. We also introduce a set of tools that enable us to detect potential distortions in one of
 2876 the two spectra. As an illustration of the capability of these tools, we use supplementary event-wise

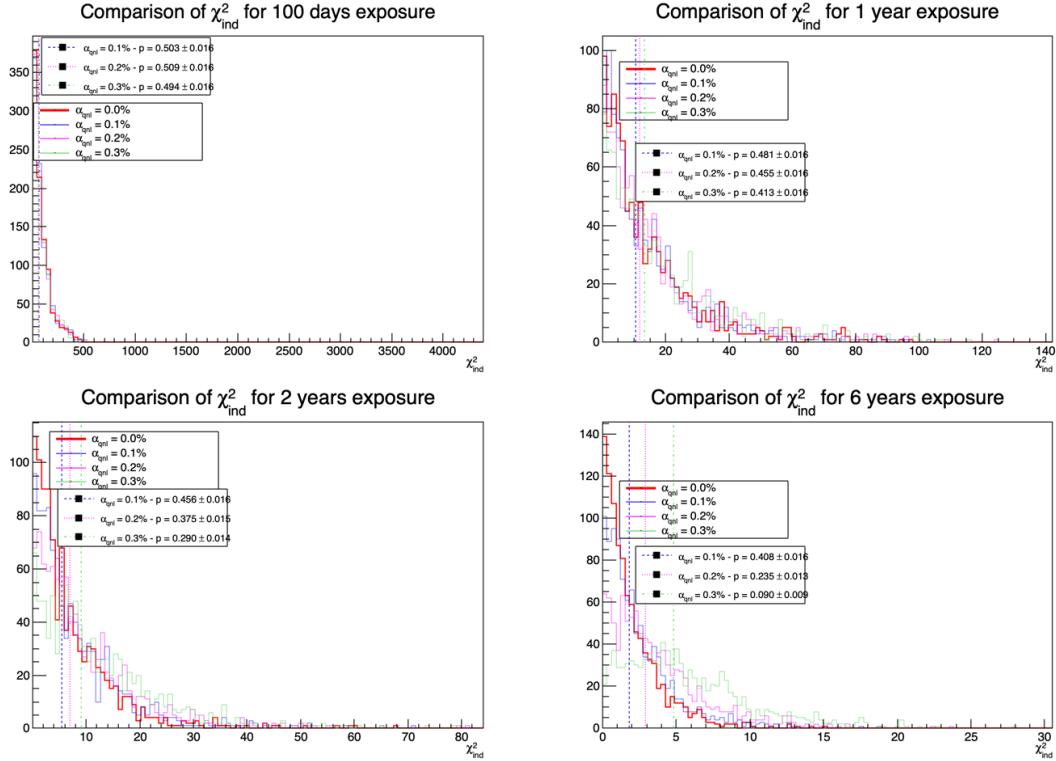


FIGURE 7.18 – Distribution of the χ^2_{Ind} for 1000 toys for different exposures. The dashed lines represent the median of the distributions and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

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.

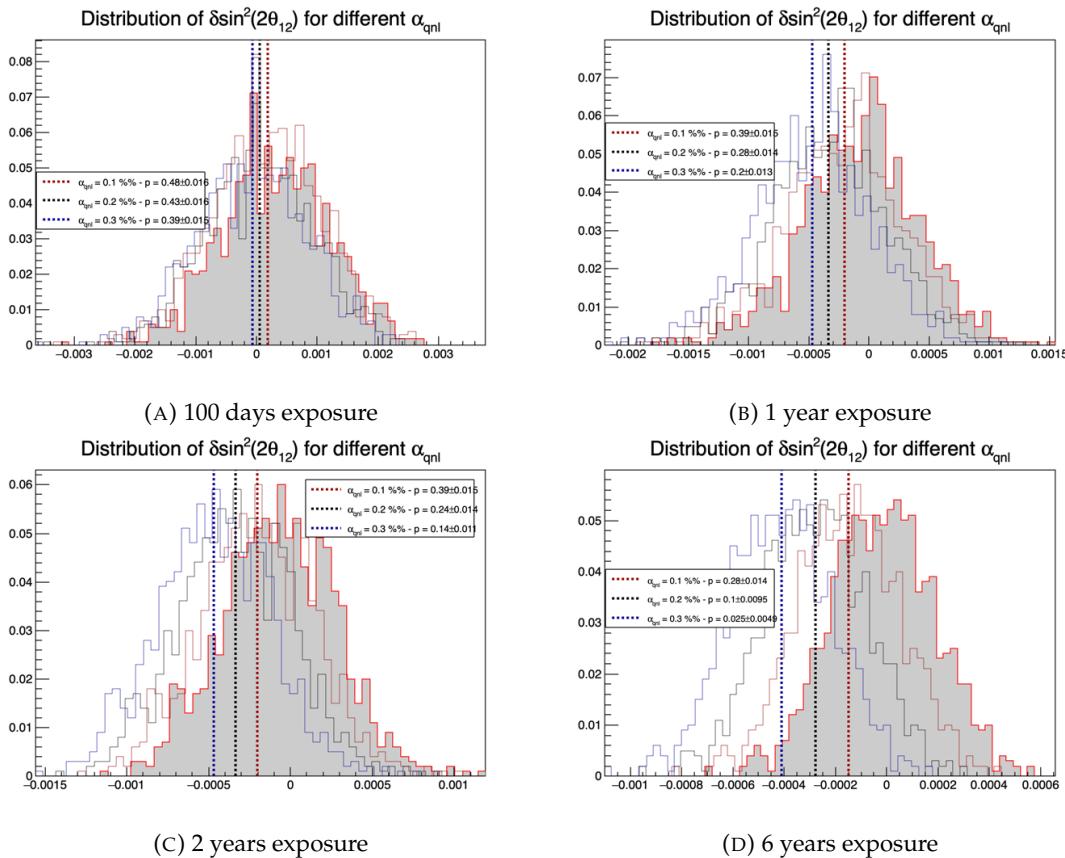


FIGURE 7.19 – Distribution of the $\delta \sin^2(2\theta_{12})$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

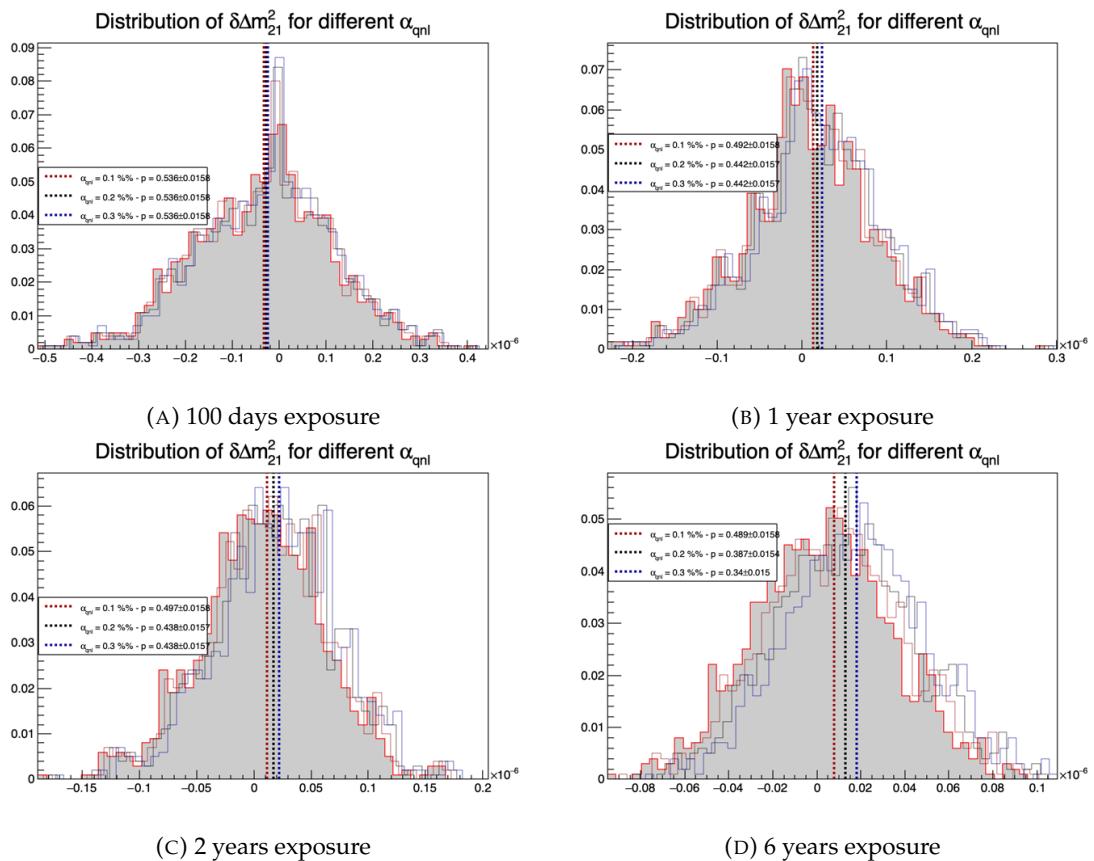


FIGURE 7.20 – Distribution of the $\delta\Delta m_{21}^2$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

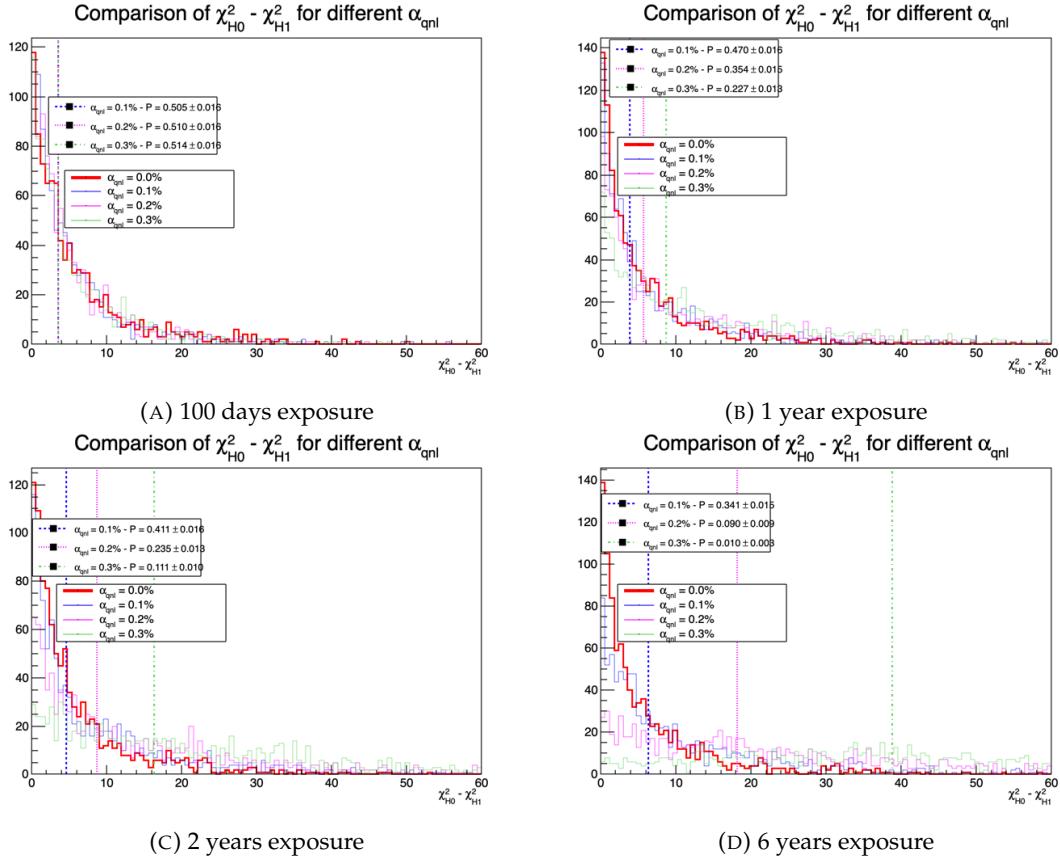


FIGURE 7.21 – Distribution of $\chi^2_{H_0} - \chi^2_{H_1}$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

²⁸⁸⁶ Chapter 8

²⁸⁸⁷ Conclusion

²⁸⁸⁸ **Appendix A**

²⁸⁸⁹ **Calculation of optimal α for estimator
combination**

²⁸⁹¹ This annex the details of the determination of the optimal α for estimator combination presented in
²⁸⁹² section 4.3.2.

²⁸⁹³ As a reminder, the combined estimator $\hat{\theta}$ of X is defined as

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

²⁸⁹⁴ where θ_N and θ_C are both estimator of X .

²⁸⁹⁵ **A.1 Unbiased estimator**

For the unbiased estimator, it is straight-forward. We search α such as $E[\hat{\theta}] = X$

$$E[\hat{\theta}] = E[\alpha\theta_N + (1 - \alpha)\theta_C] \quad (\text{A.2})$$

$$= E[\alpha\theta_N] + E[(1 - \alpha)\theta_C] \quad (\text{A.3})$$

$$= \alpha E[\theta_N] + (1 - \alpha)E[\theta_C] \quad (\text{A.4})$$

$$= \alpha(\mu_N + X) + (1 - \alpha)(\mu_C + X) \quad (\text{A.5})$$

$$X = \alpha\mu_N + \mu_C - \alpha\mu_C + X \quad (\text{A.6})$$

$$0 = \alpha(\mu_N - \mu_C) + \mu_C \quad (\text{A.7})$$

$$(A.8)$$

$$\Rightarrow \alpha = \frac{\mu_C}{\mu_C - \mu_N} \quad (\text{A.9})$$

²⁸⁹⁶ **A.2 Optimal variance estimator**

The α for this estimator is a bit more tricky. By expanding the variance we get

$$\text{Var}[\hat{\theta}] = \text{Var}[\alpha\theta_N + (1 - \alpha)\theta_C] \quad (\text{A.10})$$

$$= \text{Var}[\alpha\theta_N] + \text{Var}[(1 - \alpha)\theta_C] + \text{Cov}[\alpha(1 - \alpha)\theta_N\theta_C] \quad (\text{A.11})$$

$$= \alpha^2\sigma_N^2 + (1 - \alpha)^2\sigma_C^2 + 2\alpha(1 - \alpha)\sigma_N\sigma_C\rho_{NC} \quad (\text{A.12})$$

²⁸⁹⁷ where, as a reminder, ρ_{NC} is the correlation factor between θ_C and θ_N .

Now we try to find the minima of $\text{Var}[\hat{\theta}]$ with respect to α . 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

$$\frac{d}{d\alpha} \text{Var}[\hat{\theta}] = 0 \quad (\text{A.15})$$

$$2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) - 2\sigma_C^2 + 2\sigma_N\sigma_C\rho_{NC} = 0 \quad (\text{A.16})$$

$$2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) = 2\sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC} \quad (\text{A.17})$$

$$\alpha = \frac{\sigma_C^2 - \sigma_N\sigma_C\rho_{NC}}{\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}} \quad (\text{A.18})$$

2898 This equation shows only one solution which is a minima. From Eq. A.18 arise two singularities:

- 2899 — $\sigma_N = \sigma_C = 0$. This is not a problem because as physicists we never measure with an absolute
2900 precision, neither us or our detectors are perfect.
- 2901 — $\sigma_N = \sigma_C$ and $\rho_{CN} = 1$. In this case θ_C and θ_N are the same estimator in term of variance thus
2902 any value for α yield the same result: an estimator with the same variance as the original ones.

²⁹⁰³ **Appendix B**

²⁹⁰⁴ **Charge spherical harmonics analysis**

²⁹⁰⁵ When looking at JUNO events we can clearly see some pattern in the charge repartition based on
²⁹⁰⁶ the event radius as illustrated in figure B.4. When dealing with identifying features and pattern on a
²⁹⁰⁷ spherical plane, the astrophysics community have been using, with success, the spherical harmonic
²⁹⁰⁸ decomposition. The principle is similar to a frequency analysis via Fourier transform. It comes to
²⁹⁰⁹ saying that a function $f(r, \theta, \phi)$, here our charge repartition of the spherical plane constructed by our
²⁹¹⁰ PMTs, can be expressed

$$f(r, \theta, \phi) = \sum_{l=0}^{\infty} \sum_{m=-l}^l a_l^m r^l Y_l^m(\theta, \phi) \quad (\text{B.1})$$

²⁹¹¹ 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
²⁹¹² degree l and order m and P_l^m their associated Legendre Polynomials. Those harmonics are illustrated
²⁹¹³ in figure B.1. By reducing the problem to the unit sphere $r = 1$, we get rid of the term r^l . The Healpix
²⁹¹⁴ library [75] offer function to efficiently find the a_l^m factor from a given Healpix map.

²⁹¹⁵ For the above decomposition, we will define the *Power* of an harmonic as

$$S_{ff}(l) = \frac{1}{2l+1} \sum_{m=-l}^l |a_l^m|^2 \quad (\text{B.2})$$

²⁹¹⁶ and the *Relative Power* as:

$$P_l^h = \frac{S_{ff}(l)}{\sum_l S_{ff}(l)} \quad (\text{B.3})$$

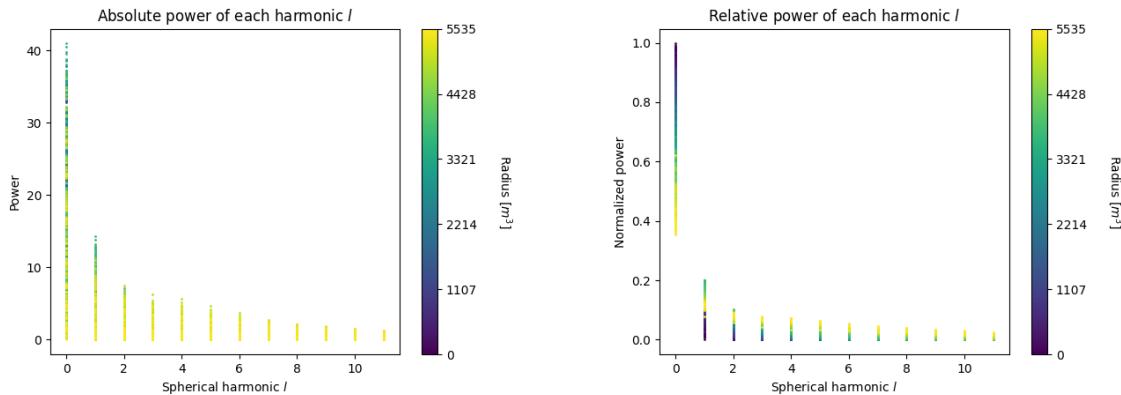
²⁹¹⁷ For this study we will use 10k positron events with $E_{kin} \in [0; 9]$ MeV uniformly distributed in the
²⁹¹⁸ CD from the JUNO official simulation version J23.0.1-rc8.dc1 (released the 7th January 2024). All the
²⁹¹⁹ event are *calib* level, with simulation of the physics, electronics, digitizations and triggers. We first
²⁹²⁰ take a sub-set of 1k events and look at the power and relative power distribution depending on the
²⁹²¹ radius and harmonic degree l . The results are shown in figure B.2. While don't see any pattern in
²⁹²² absolute power, it is pretty clear that there is a correlation between the relative power of $l = 0$ and
²⁹²³ the radius of the event.

²⁹²⁴ When applying the same study but dependent on the energy, no clear correlation appear. The results
²⁹²⁵ for the $l = 0$ harmonic are presented in the figure B.5. Thus, in this study we will focus on the radial
²⁹²⁶ dependency of the relative power of each harmonic.

²⁹²⁷ In figures B.6 and B.7 are presented the distribution of the relative power of each harmonic for $l \in$
²⁹²⁸ $[0, 11]$. The relation between the radius and the relative power become even more clear, especially
²⁹²⁹ for the first harmonics $l \in [0, 4]$. After that for $l > 4$ their relative power is close to 0 for central event,
²⁹³⁰ thus loosing power. It also interesting to note the change of behavior in the TR area, clearly visible
²⁹³¹ 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 $l = 0$.

2943 The minimization algorithm used are the Bent algorithm for $l = 0$ and the Bounded algorithm for
 2944 $l > 0$ provided by the Scipy library [89]. We then do the mean of the reconstructed radius from
 2945 the different harmonics. The reconstruction results are shown in figure B.3. The performance seems
 2946 correct but we see heavy fluctuation in the bias. To really be used as a reconstruction algorithm, the
 2947 method needs to be refined as discussed in the next section.

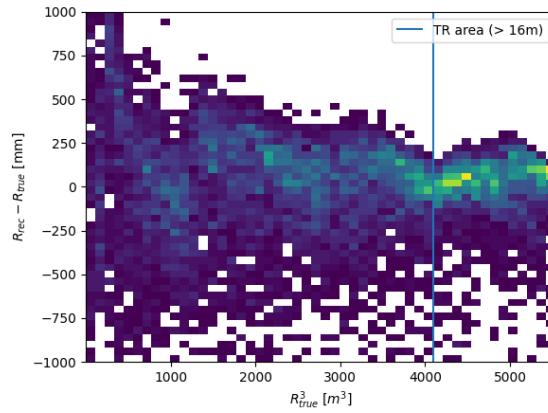


FIGURE B.3 – Error on the reconstructed radius vs the true radius by the harmonic method

Conclusion

2948 We have clearly shown in this analysis the relevance the of relative harmonic power for radius
 2949 reconstruction, and provided an erzats of a reconstruction algorithm. We will not delve further in
 2950 this thesis but if we wanted to refine this algorithm multiple paths can be explored:
 2951

- 2952 — No energy signature in the harmonics: This is surprising that there is no correlation between
 2953 the energy and the amplitude of the harmonics. We know that the energy is heavily correlated
 2954 with the total number of photoelectrons collected, it would be unintuitive that we see no
 2955 relation.
- 2956 — Localization of the event: We shown here the relation between the relative power of the har-
 2957 monic and the radius but don't get any information about the θ and ϕ spherical coordinates.
 2958 This information is probably hidden in the individual power of each order m of the degree l .
 2959 This intuition comes from the figure B.1 where in the higher degree l we see that the order m
 2960 are oriented. Intuitively, the order should be able to indicate a direction where the signal is
 2961 more powerful.
- 2962 — Combination of the degree power: Here we combined the radius reconstructed by the dif-
 2963 ferent degree via a simple mean but we shown in section 4.3.2 and annex A that this is note
 2964 the optimal way to combine estimator. A more refined algorithm probably exist to take into
 2965 account the predicting power of each order.

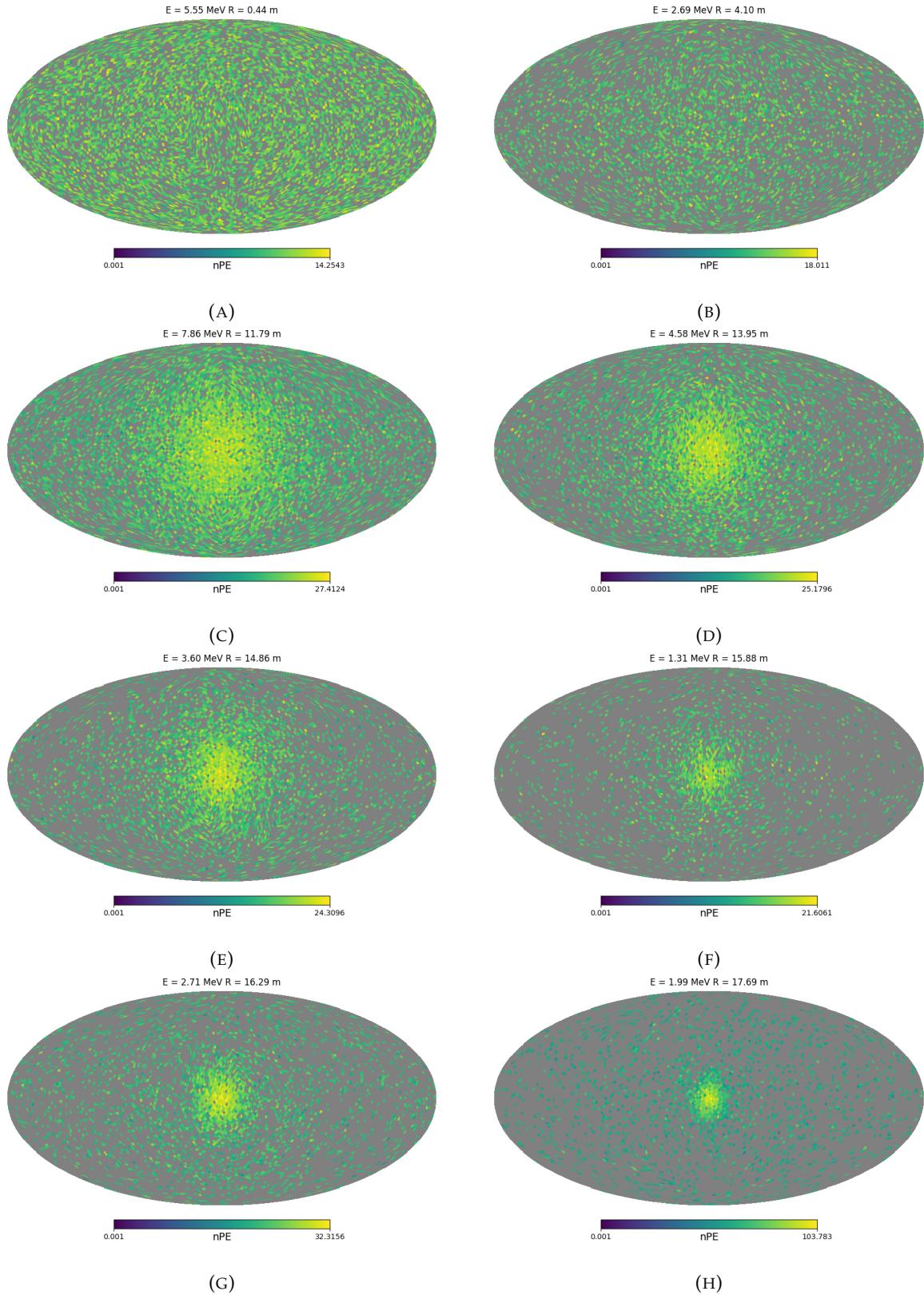


FIGURE B.4 – Charge repartition in JUNO as seen by the Healpix segmentation. Those are Healpix map of order 5 (i.e. 12288 pixels). The color represent the summed charge of the PMTs in each pixels. The color scale is logarithmic. The view have been centered to prevent event deformations.

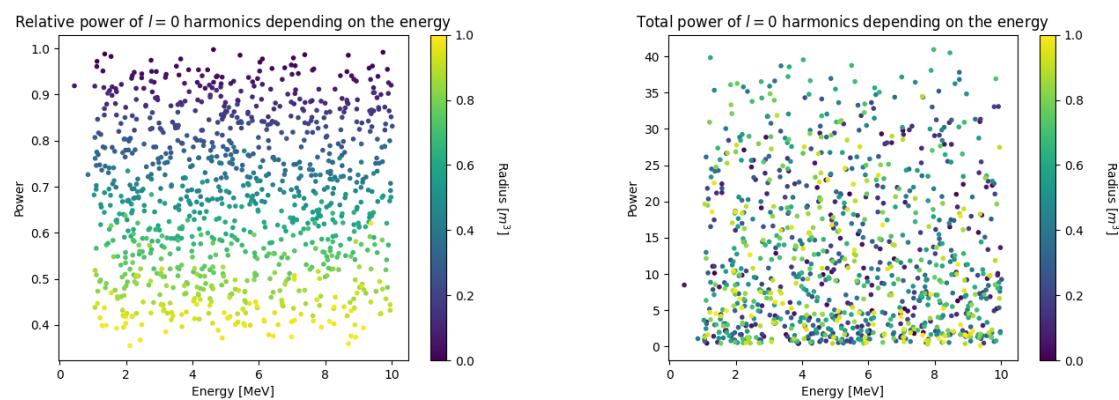


FIGURE B.5 – Scatter plot of the absolute and relative power, respectively on the left and right plot, of the $l = 0$ harmonic. The color indicate the radius of the event.

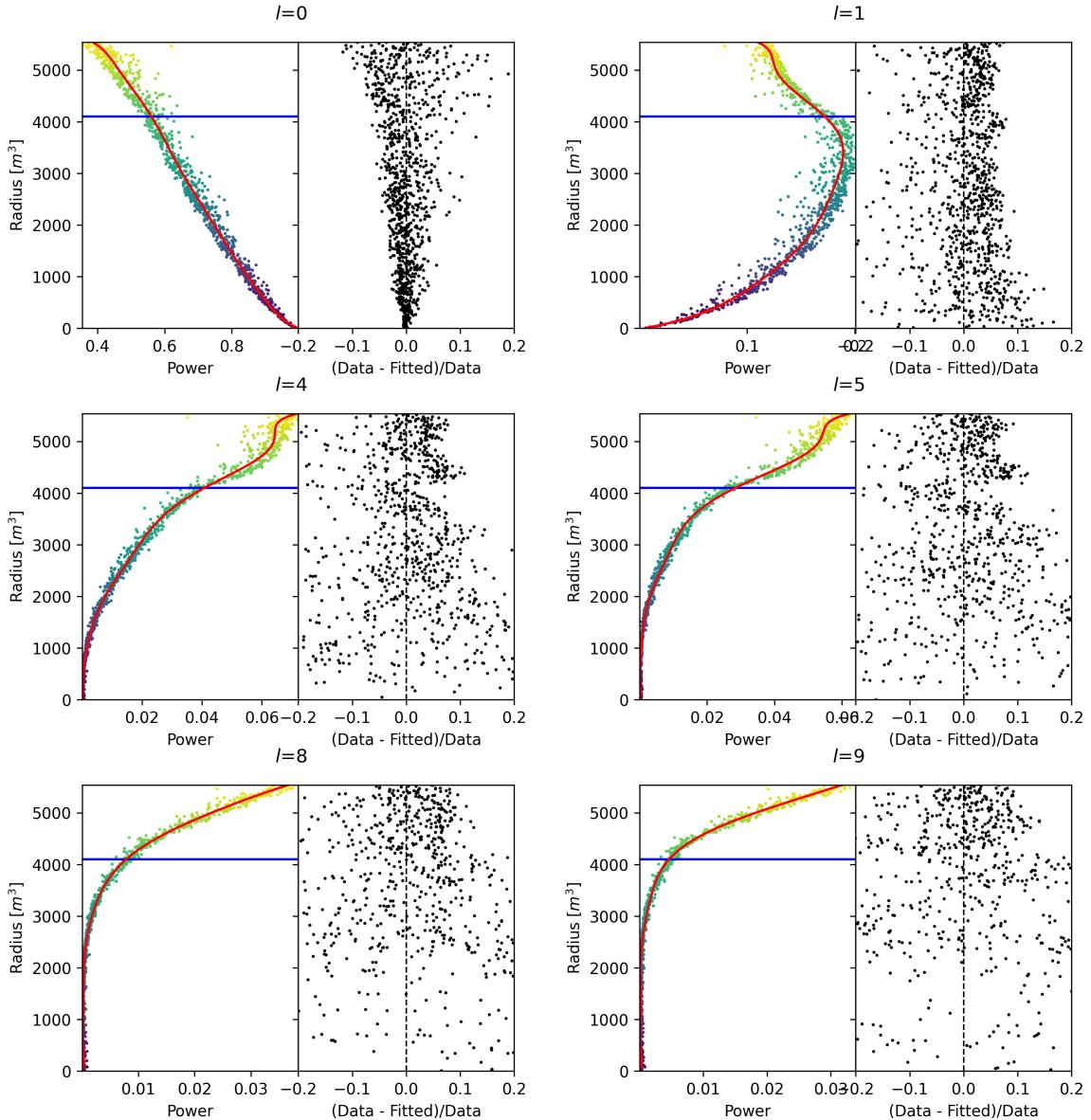


FIGURE B.6 – Plot of the distribution of the relative power of each harmonic dependent on R^3 (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. **Part 1**

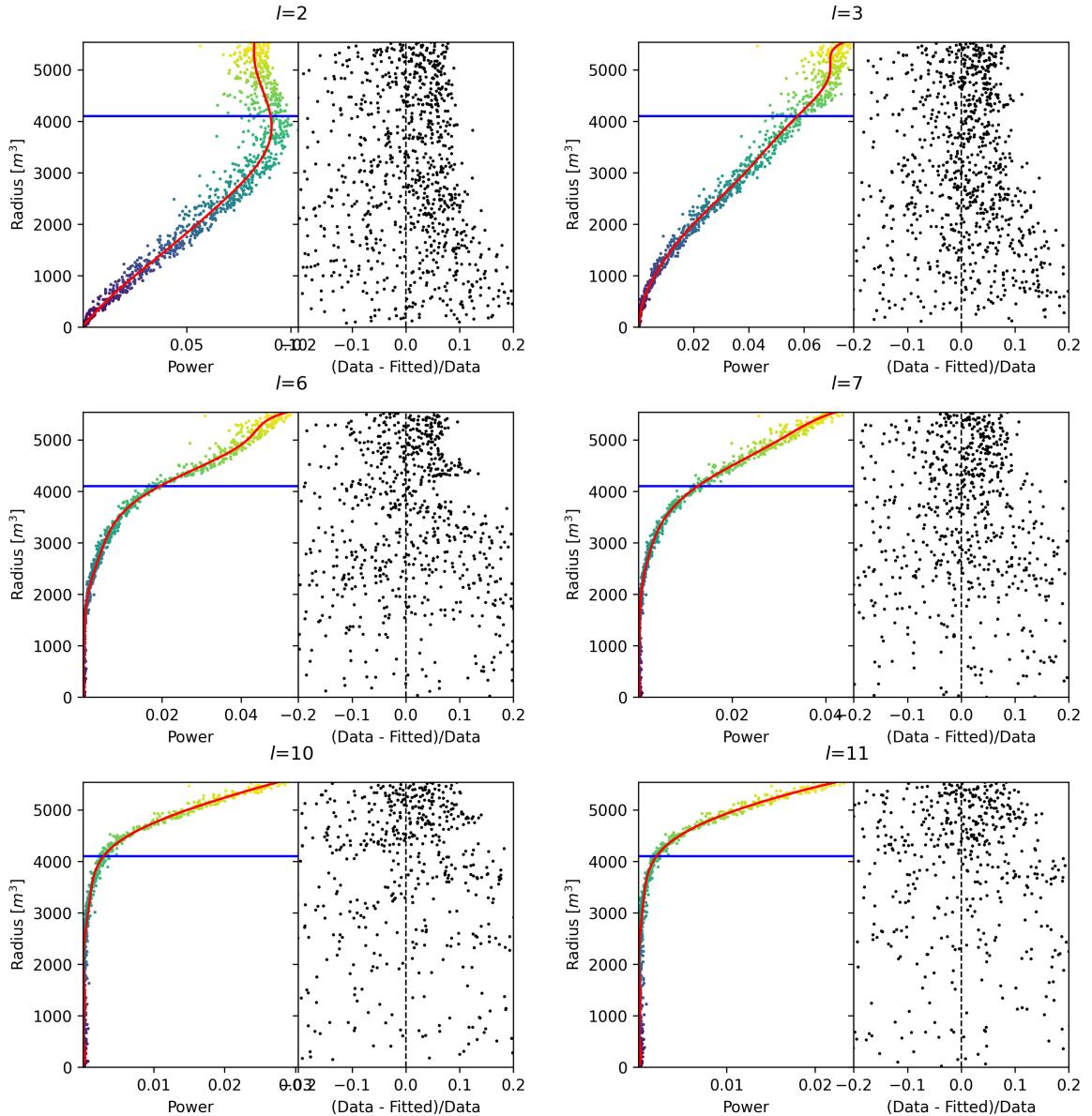


FIGURE B.7 – Plot of the distribution of the relative power of each harmonic dependent on R^3 (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. **Part 2**

²⁹⁶⁶ **Appendix C**

²⁹⁶⁷ **Additional spectrum smearing**

²⁹⁶⁸ In this section we demonstrate that a spectrum S smeared by a gaussian G parametrized by its
²⁹⁶⁹ varianse σ_1^2 can be smeared by a gaussian parametrized by the variance σ_2^2 from the the smeared
²⁹⁷⁰ 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

$$K'(E, \sigma_2) = S(E) \star G(E, \sigma_1) \star G^{-1}(E, \sigma_1) \star G(E, \sigma_2) \quad (\text{C.1})$$

$$= K(E, \sigma_1) \star G^{-1}(E, \sigma_1) \star G(E, \sigma_2) \quad (\text{C.2})$$

²⁹⁷¹ 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

$$G(E, \sigma_1) \star G^{-1}(E, \sigma_1) = \delta(E) \quad (\text{C.3})$$

$$F[G(E, \sigma_1)](\nu) \times F[G^{-1}(E, \sigma_1)](\nu) = 1 \quad (\text{C.4})$$

²⁹⁷² 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

$$F[G^{-1}(E, \sigma_1) \star G(E, \sigma_2)] = e^{\frac{\sigma_1^2(2\pi)^2}{2}\nu^2} \times e^{-\frac{\sigma_2^2(2\pi)^2}{2}\nu^2} \quad (\text{C.8})$$

$$= e^{\frac{(2\pi)^2}{2}(\sigma_1^2 - \sigma_2^2)\nu^2} \quad (\text{C.9})$$

$$= e^{\frac{(2\pi)^2}{2}\Delta\sigma^2\nu^2}; \Delta\sigma^2 = (\sigma_1^2 - \sigma_2^2) \quad (\text{C.10})$$

²⁹⁷³ 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

$$G^{-1}(E, \sigma_1) \star G(E, \sigma_2) = \frac{1}{\sqrt{|\Delta\sigma^2|}\sqrt{2\pi}} e^{-\frac{E^2}{2|\Delta\sigma^2|}} \quad (\text{C.11})$$

²⁹⁷⁴ **Appendix D**

²⁹⁷⁵ **Correction of E_{vis} bias**

²⁹⁷⁶ The reconstruction algorithms that are presented in this thesis in Chapters 4 and 5 do not reconstruct
²⁹⁷⁷ the same energy as the classical algorithms presented in section 2.8. Our algorithms reconstruct the
²⁹⁷⁸ deposited energy E_{dep} while the classical algorithms reconstruct a visible energy E_{vis} .

To understand this phenomena, let's look at the equation 2.23:

$$\hat{\mu}(r, \theta, \theta_{pmt}, E_{vis}) = \frac{1}{E_{vis}} \frac{1}{M} \sum_i^M \frac{\frac{\bar{Q}_i}{\bar{Q}_i} - \mu_i^D}{DE_i}, \quad \mu_i^D = DNR_i \cdot L$$

²⁹⁷⁹ which define the expected N_{pe}/E . This define a linear relation between the number of photoelectrons
²⁹⁸⁰ and the energy. However we discussed in sections 2.3.2 and 2.4 that the number of photoelectrons
²⁹⁸¹ collected by the LPMT system do not follow a linear relationship. Thus this visible energy is not
²⁹⁸² linear with the deposited energy. This effect is corrected in physics analysis and in Chapter 7 by
²⁹⁸³ applying the calibrated non-linearity profile the energy spectrum.

²⁹⁸⁴ When we need to compare our algorithm that reconstruct the deposited energy to the classical
²⁹⁸⁵ algorithms we need to correct this non-linearity. For this we fit the systematic bias of the classical
²⁹⁸⁶ algorithm using a 5th degree polynomial

$$\frac{E_{dep}}{E_{vis}} = \sum_{i=0}^5 P_i E_{dep}^i \quad (D.1)$$

²⁹⁸⁷ The fitted distribution and the corresponding fit is presented in figure D.1. The value fitted for this
²⁹⁸⁸ correction are presented in table D.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 D.1 – Parameters of the 5th degree polynomial used to correct Omilrec reconstructed energy.

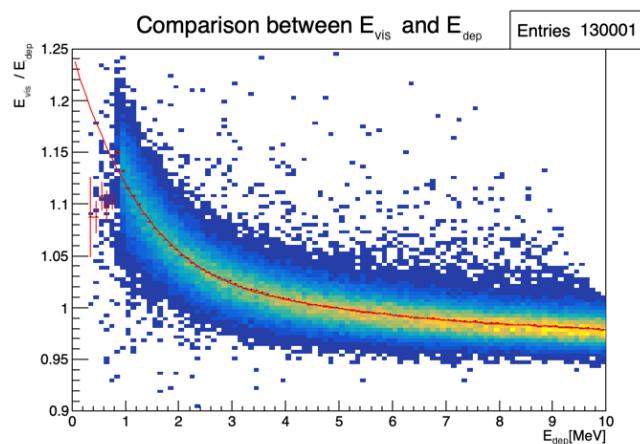


FIGURE D.1 – Comparison between Omilrec reconstructed E_{vis} and the deposited energy E_{dep} . The profile of the distribution E_{vis}/E_{dep} vs E_{dep} is fitted with a 5th degree polynomial.

List of Tables

2990	2.1	Characteristics of the nuclear power plants observed by JUNO.	14
2991	2.2	Detectable neutrino signal in JUNO and the expected signal rates and major background sources	15
2992	2.3	List of sources and their process considered for the energy scale calibration	24
2993	2.4	Calibration program of the JUNO experiment	25
2994	2.5	Summary of cumulative reactor antineutrino selection efficiencies. The reported IBD rates (with baselines <300 km) refer to the expected events per day after the selection criteria are progressively applied. Table taken from [32]	29
2995	2.6	Expected background rates, background to signal ratio (B/S), and rate and shape uncertainties. The B/S ratio is calculated by using the IBD signal rate of 47.1/day. Table taken from [32]	29
2996	2.7	A summary of precision levels for the oscillation parameters. The reference value (PDG 2020 [34]) is compared with 100 days, 6 years and 20 years of JUNO data taking.	32
2997	2.8	Features used by the BDT for vertex reconstruction	42
2998	2.9	Features used by the BDTE algorithm. <i>pe</i> and <i>ht</i> reference the charge and hit-time distribution respectively and the percentages are the quantiles of those distributions. <i>cht</i> and <i>cc</i> reference the barycenters of hit time and charge respectively	42
3000	4.1	Sets of hyperparameters values considered in this study	62
3001	5.1	Features on the nodes of the graph. All charge are in [nPE], time in [ns] and position in [m]. <i>Q</i> and <i>t</i> are the reconstructed charge and time of the hit PMTs. (<i>x</i> , <i>y</i> , <i>z</i>) is the position of the PMTs and the last parameter represent the type of the PMT. It's 1 for LPMT and -1 for SPMT <i>Q_m</i> and <i>t_m</i> is the set of charges and time of the PMT belonging the mesh <i>m</i> . (<i>X_m</i> , <i>Y_m</i> , <i>Z_m</i>) i the position of the center of the geometric region represented by the mesh <i>m</i> ($\langle X \rangle$, $\langle Y \rangle$, $\langle Z \rangle$) is the position of the charge barycenter, $\sum Q$ the sum of the collected charge in the detector and <i>P_l^h</i> is the relative power of the <i>l</i> th harmonic. See annex B for details.	80
3002	5.2	Features on the edges on the graph. It use the same notation as in table 5.1. <i>D_{m1→m2}⁻¹</i> is the inverse of the distance between the mesh <i>m1</i> and the mesh <i>m2</i> . The features A and B are detailed in section 5.1	81
3003	7.1	The charge fraction in terms of the number of PE collected at the single PMT for the reactor $\bar{\nu}_e$ IBD events. Table taken from [24]	104
3004	7.2	Nominal PDG2020 value [34]. All value are reported assuming Normal Ordering.	108
3005	7.3	Results of the Asimov studies on the updated framework. All results are Asimov fit, considering 6 years exposure, θ_{13} is fixed to nominal value, χ^2 is pearson meaning that he error is estimated using the data spectrum	114
3006	7.4	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 Δm_{31}^2 is fixed at nominal value. The χ^2 is PearsonV. The correlation matrix used to fit assume no QNL in the spectrum.	116
3029	D.1	Parameters of the 5th degree polynomial used to correct Omilrec reconstructed energy.	143

List of Figures

3031	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	12
3032		
3033		
3034		
3035		
3036		
3037		
3038		
3039		
3040	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 (θ_{12} , Δm_{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
3041		
3042		
3043		
3044	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 [3]	14
3045		
3046		
3047		
3048	2.4	17
3049	a Schematics view of the JUNO detector.	17
3050	b Top down view of the JUNO detector under construction	17
3051	2.5 Schematics of an IBD interaction in the central detector of JUNO	18
3052	2.6 Schematics of the supporting node for the acrylic vessel	19
3053	2.7 On the left: Quantum efficiency (QE) and emission spectrum of the LAB and the bis-MSB [17]. On the right: Sensitivity of the Hamamatsu LPMT depending on the wavelength of the incident photons [19].	20
3054		
3055	2.8 Schematic of a PMT	20
3056	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	21
3057		
3058	2.10 Schematic of the JUNO SPMT electronic system (left), and exploded view of the main component of the UWB (right)	22
3059	2.11 The JUNO top tracker	23
3060	2.12 Fitted and simulated non linearity of gamma, electron sources and from the ^{12}B spectrum. Black points are simulated data. Red curves are the best fits. Figures taken from [26].	24
3061	a Gamma non-linearity	24
3062	b Boron spectrum	24
3063	c Electron non-linearity	24
3064		
3065	2.13 Overview of the calibration system	25
3066	2.14 Event-level instrumental non-linearity, defined as the ratio of the total measured LPMT charge to the true charge for events at the center of the detector. The solid red line represents event-level non-linearity without the channel-level correction 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 [26].	26
3067		
3068		
3069		
3070		
3071		
3072	2.15	27
3073	a Schematic of the TAO satellite detector	27

3074	b Schematic of the OSIRIS satellite detector	27
3075	2.16 Illustration of the spectrum considered when joint fitting	33
3076	2.17	35
3077	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.	35
3081	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	35
3083	2.18	36
3084	a Δt distribution at different iterations step j	36
3085	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	36
3087	2.19 Bias of the reconstructed radius R (left), θ (middle) and ϕ (right) for multiple energies by the time likelihood algorithm	37
3089	2.20 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	38
3091	2.21 Radial resolution of the different vertex reconstruction algorithms as a function of the energy	38
3092	2.22	39
3096	a Spherical coordinate system used in JUNO for reconstruction	39
3097	b Definition of the variables used in the energy reconstruction	39
3098	2.23	41
3099	a Radial resolutions of the likelihood-based algorithm TMLE, QMLE and QTMLE	41
3100	b Energy resolution of QMLE and QTMLE using different vertex resolutions	41
3101	2.24 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	43
3102	2.25 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.	43
3106	3.1 Example of a BDT that determine if the given object is a duck	46
3107	3.2 Schema of a simple neural network	47
3108	3.3 Illustration of the training lifecycle	49
3109	3.4	50
3110	a Illustration of SGD falling into a local minima	50
3111	b Illustration of the Adam momentum allowing it to overcome local minima	50
3112	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.	51
3113	a Illustration of the SGD optimizer on one parameter θ on the MAE Loss. We see here that it has trouble reaching the minima due to the gradient being constant.	51
3117	b Illustration of the SGD optimizer on one parameter θ 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.	51
3120	3.6	52
3121	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.	52
3122	b Illustration of a very simple NN	52

3126	3.7 Illustration of the ResNet framework	53
3127	3.8 Illustration of the gradient explosion. Here it can be solved with a lower learning rate but its not always the case.	53
3128		
3129	3.9	54
3130	a Schema of a FCDNN	54
3131	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	54
3132		
3133	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].	55
3134		
3135	3.11	56
3136	a Example of images in the MNIST dataset	56
3137	b Schema of the CNN used in Pytorch example to process the MNIST dataset	56
3138		
3139	3.12 Illustration of a graph and its tensor representation.	57
3140		
3141	3.13 Illustration of the message passing algorithm. The detailed explanation can be found in section 3.2.3	57
3142		
3143	4.1 Graphic representation of the VGG-16 architecture, presenting the different kind of layer composing the architecture.	60
3144		
3145	4.2	65
3146	a Spherical coordinate system used in JUNO for reconstruction	65
3147	b Repartition of SPMTs in the image projection. The color scale is the number of SPMTs per pixel	65
3148		
3149	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.	65
3150		
3151		
3152		
3153	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.	66
3154		
3155		
3156		
3157	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.	66
3158		
3159		
3160		
3161	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.	67
3162		
3163		
3164		
3165	4.7	68
3166	a Distribution of PE/MeV in the J23 Dataset. This distribution is profiled and fitted using equation 4.6	68
3167		
3168	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	68
3169		
3170	4.8 Reconstruction performance of the Gen ₃₀ model on J21 data and it's comparison to the performances of the classic algorithm “Classical algorithm” from [65]. The top part of each plot is the resolution and the bottom part is the bias.	69
3171		
3172		
3173	a Resolution and bias of energy reconstruction vs energy	69
3174	b Resolution and bias of energy reconstruction vs radius	69
3175	c Resolution and bias of radius reconstruction vs energy	69
3176	d Resolution and bias of radius reconstruction vs radius	69
3177	e Resolution and bias of radius reconstruction vs θ	69

3178	f	Resolution and bias of radius reconstruction vs ϕ	69
3179	4.9	Residual 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 , θ and ϕ .	70
3180	a	Distribution of the error on reconstructed x by Gen ₃₀	70
3181	b	Distribution of the error on reconstructed y by Gen ₃₀	70
3182	c	Distribution of the error on reconstructed z by Gen ₃₀	70
3183	d	Distribution of the error on reconstructed R by Gen ₃₀	70
3184	e	Distribution of the error on reconstructed θ by Gen ₃₀	70
3185	f	Distribution of the error on reconstructed ϕ by Gen ₃₀	70
3186	4.10	...	71
3187	a	Distribution of Gen ₃₀ reconstructed energy and true energy of the analysis dataset (J21)	71
3188	b	Distribution of Gen ₄₂ reconstructed energy and true energy of the analysis dataset (J23)	71
3189	4.11	Radius bias (on the left) and resolution (on the right) of the classical algorithm in a E , R^3 grid	72
3190	4.12	Reconstruction performance of the Gen ₃₀ model on J21, the classic algorithm "Classical algorithm" from [65] and the combination of both using weighted mean. The top part of each plot is the resolution and the bottom part is the bias.	73
3191	a	Resolution and bias of energy reconstruction vs energy	73
3192	b	Resolution and bias of energy reconstruction vs radius	73
3193	c	Resolution and bias of radius reconstruction vs energy	73
3194	d	Resolution and bias of radius reconstruction vs radius	73
3195	e	Resolution and bias of radius reconstruction vs θ	73
3196	f	Resolution and bias of radius reconstruction vs ϕ	73
3197	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	74
3198	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 [65]. The top part of each plot is the resolution and the bottom part is the bias.	75
3199	a	Resolution and bias of energy reconstruction vs energy	75
3200	b	Resolution and bias of energy reconstruction vs radius	75
3201	c	Resolution and bias of radius reconstruction vs energy	75
3202	d	Resolution and bias of radius reconstruction vs radius	75
3203	e	Resolution and bias of radius reconstruction vs θ	75
3204	f	Resolution and bias of radius reconstruction vs ϕ	75
3205	5.1	...	79
3206	a	Illustration of the different nodes in our graphs and their relations	79
3207	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	79
3208	5.2	Illustration of the Healpix segmentation. On the left: A segmentation of order 0. On the right: A segmentation of order 1	79
3209	5.3	Illustration of the different update function needed by our GNN	82
3210	5.4	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	83
3211	a	...	83
3212	b	...	83
3213	c	...	83

3229	5.5	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 $E_{true} \in [0, 9]$	83
3230	a	83
3231	b	83
3232			
3233	5.6	Schema of the JWGv8.4.0 architecture, the colored triplet is the graph configuration after each JWG layers	85
3234			
3235	5.7	Energy reconstruction depending on the true energy for samples of the different versions of the GNN	86
3236			
3237	5.8	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.8, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.	88
3238	a	Resolution and bias of energy reconstruction vs energy	88
3239	b	Resolution and bias of energy reconstruction vs radius	88
3240			
3241	5.9	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.8, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.	89
3242	a	Resolution and bias of radius reconstruction vs energy	89
3243	b	Resolution and bias of radius reconstruction vs radius	89
3244			
3245	5.10	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.8, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.	90
3246	a	Resolution and bias of radius reconstruction vs θ	90
3247	b	Resolution and bias of radius reconstruction vs ϕ	90
3248			
3249	5.11	Reconstruction performance of the Omilrec algorithm, JWGv8.4 and the combination between the two using the optimal variance estimator presented in annex A.2. The top part of each plot is the resolution and the bottom part is the bias.	91
3250	a	Resolution and bias of energy reconstruction vs energy	91
3251	b	Resolution and bias of energy reconstruction vs radius	91
3252			
3253	5.12	Reconstruction performance of the Omilrec algorithm, JWGv8.4 and the combination between the two using the optimal variance estimator presented in annex A.2. The top part of each plot is the resolution and the bottom part is the bias.	92
3254	a	Resolution and bias of radius reconstruction vs energy	92
3255	b	Resolution and bias of radius reconstruction vs radius	92
3256			
3257	5.13	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.8, JWGv8.4 presented in this chapter and the HCNN algorithm. The top part of each plot is the resolution and the bottom part is the bias.	93
3258	a	Resolution and bias of energy reconstruction vs energy	93
3259	b	Resolution and bias of energy reconstruction vs radius	93
3260			
3261	5.14	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.8, JWGv8.4 presented in this chapter and the HCNN algorithm. The top part of each plot is the resolution and the bottom part is the bias.	93
3262	a	Resolution and bias of radius reconstruction vs energy	93
3263	b	Resolution and bias of radius reconstruction vs radius	93
3264			
3265	6.1	Schema of the method to discover vulnerabilities in the reconstruction methods. On the top of the image , the standard data flow. The individual charge and times are fed to a reconstruction algorithm. From the reconstructed energies, we can produce an IBD spectrum and compute control observables from the control samples. On the bottom , the same data flow but we add an ANN between the input and the reconstruction. The ANN will slightly change the input charge and time so the reconstruction algorithm inaccurately reconstruct the IBD energy, but the perturbation is not visible in the control sample.	97
3266			
3267			
3268			
3269			
3270			
3271			
3272			
3273			
3274			
3275			
3276			
3277			
3278			
3279			
3280			

3281 7.1	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.	102
3288 7.2	Oscillated reactor $\bar{\nu}_e$ spectra for the Normal Ordering (Black) and Inverted Ordering (Red) for 6,5 years data taking and a resolution of 3% without any statistical or sys- tematic fluctuation. Figure from [32].	103
3291 7.3	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.	105
3294 7.4	106
3295 a	Distribution of ratio of collected nPE after the additional QNL over the number of nPE that would be collected for different γ_{qnl} . We select event with an interaction radius $R < 4m$ to not be affected by the non-uniformity.	106
3296 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 7.4a 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 [26].	106
3304 7.5	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 .	110
3306 7.6	Schematic description of the fit framework	111
3308 7.7	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.	112
3310 7.8	Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, Pearson χ^2, θ_{13} fixed.	115
3311 7.9	Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, PearsonV χ^2, θ_{13} fixed.	116
3312 7.10	Distribution of BFP - nominal value for 5000 toy Delta joint fit. 6 years exposure, all background, PearsonV χ^2, θ_{13} fixed.	117
3314 7.11	Top: Theoretical spectrum without QNL (in red) and with $\alpha_{qnl} = 1\%$ (in blue). Bottom: Ratio between the theoretical spectrum with and without QNL.	118
3317 7.12	Theoretical correlation matrix between the LPMT spectrum (bins 0-409) ans the SPMT spectrum (410-819). The diagonal has been set to 0 (it was 1) for readability purpose. .	119
3319 7.13	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	120
3322 a	120
3323 b	120
3324 c	120
3325 7.14	Difference between the element of the theoretical and empiric correlation matrix . .	121
3326 a	121
3327 b	121
3328 7.15	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 recon- struction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.8. To prevent effect due to the CNN bad reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.	122

3333	7.16 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 2.8. To prevent effect due to the CNN bad reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.	123
3334		
3335		
3336		
3337		
3338	7.17 Distribution of the χ^2_{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.	124
3339		
3340		
3341	7.18 Distribution of the χ^2_{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.	125
3342		
3343		
3344	7.19 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.	126
3345		
3346	a 100 days exposure	126
3347		
3348	b 1 year exposure	126
3349		
3350	c 2 years exposure	126
3351		
3352	d 6 years exposure	126
3353		
3354	7.20 Distribution of the $\delta \Delta m^2_{21}$ 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.	127
3355		
3356	a 100 days exposure	127
3357		
3358	b 1 year exposure	127
3359		
3360	c 2 years exposure	127
3361		
3362	d 6 years exposure	127
3363		
3364	7.21 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.	128
3365		
3366	a 100 days exposure	128
3367		
3368	b 1 year exposure	128
3369		
3370	c 2 years exposure	128
3371		
3372	d 6 years exposure	128
3373		
3374	B.1 Illustration of the real part of the spherical harmonics	134
3375		
3376	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.	134
3377		
3378	B.3 Error on the reconstructed radius vs the true radius by the harmonic method	135
3379		
3380	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.	136
3381		
3382	a	136
	b	136
	c	136
	d	136
	e	136
	f	136
	g	136
	h	136
3383	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.	137
3384		

3383	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	138
3384			
3385			
3386	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	139
3387			
3388			
3389			
3390	D.1	Comparison between Omilrec reconstructed E_{vis} and the deposited energy E_{dep} . The profile of the distribution E_{vis}/E_{dep} vs E_{dep} is fitted with a 5th degree polynomial.	144
3391			
3392			

List of Abbreviations

ACU	Automatic Calibration Unit
BDT	Boosted Decision Tree
BFP	Best Fit Point
CD	Central Detector
CLS	Cable Loop System
CNN	Convolutional NN
DNN	Deep NN
DN	Dark Noise
EDM	Event Data Model
FCDNN	Fully Connected Deep NN
GNN	Graph NN
GT	Guiding Tube
IBD	Inverse Beta Decay
IO	Inverse Ordering
JUNO	Jiangmen Underground Neutrino Observatory
LPMT	Large PMT
LR	Learning Rate
LS	Liquid Scintillator
MC	Monte Carlo simulation
ML	Machine Learning
MSE	Mean Squared Error
NMO	Neutrino Mass Ordering
NN	Neural Network
NO	Normal Ordering
NPE	Number of Photo Electron
OSIRIS	Online Scintillator Internal Radioactivity Investigation System
PE	Photo Electron
PMT	Photo-Multipliers Tubes
PRelu	Parametrized Rectified Linear Unit
QNL	Charge (Q) Non Linearity
ROV	Remotely Operated under-LS Vehicle
ReLU	Rectified Linear Unit
ResNet	Residual Network
SGD	Stochastic Gradient Descent
SPMT	Small PMT
TAO	Taishan Antineutrino Oservatory
TR Area	Total Reflexion Area
TTS	Time Transit Spread
TT	Top Tracker
UWB	Under Water Boxes
WCD	Water Cherenkov Detector

Bibliography

- [1] Liang Zhan, Yifang Wang, Jun Cao, and Liangjian Wen. "Determination of the Neutrino Mass Hierarchy at an Intermediate Baseline". *Physical Review D* 78.11 (Dec. 10, 2008), 111103. ISSN: 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).
- [2] Fengpeng An et al. "Neutrino Physics with JUNO". *Journal of Physics G: Nuclear and Particle 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).
- [3] 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](https://arxiv.org/abs/2204.13249) (visited on 08/11/2023).
- [4] A. A. Hahn, K. Schreckenbach, W. Gelletly, F. von Feilitzsch, G. Colvin, and B. Krusche. "Antineutrino spectra from ^{241}Pu and ^{239}Pu thermal neutron fission products". *Physics Letters B* 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: <https://www.sciencedirect.com/science/article/pii/0370269389915980> (visited on 01/16/2024).
- [5] Th A. Mueller et al. "Improved Predictions of Reactor Antineutrino Spectra". *Physical Review C* 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). 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) (visited on 01/16/2024).
- [6] F. von Feilitzsch, A. A. Hahn, and K. Schreckenbach. "Experimental beta-spectra from ^{239}Pu and ^{235}U thermal neutron fission products and their correlated antineutrino spectra". *Physics 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). URL: <https://www.sciencedirect.com/science/article/pii/0370269382906220> (visited on 01/16/2024).
- [7] K. Schreckenbach, G. Colvin, W. Gelletly, and F. Von Feilitzsch. "Determination of the antineutrino spectrum from ^{235}U thermal neutron fission products up to 9.5 MeV". *Physics Letters B* 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: <https://www.sciencedirect.com/science/article/pii/0370269385913371> (visited on 01/16/2024).
- [8] Patrick Huber. "On the determination of anti-neutrino spectra from nuclear reactors". *Physical 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). URL: [http://arxiv.org/abs/1106.0687](https://arxiv.org/abs/1106.0687) (visited on 01/16/2024).
- [9] P. Vogel, G. K. Schenter, F. M. Mann, and R. E. Schenter. "Reactor antineutrino spectra and their application to antineutrino-induced reactions. II". *Physical Review C* 24.4 (Oct. 1, 1981). Publisher: American Physical Society, 1543–1553. DOI: [10.1103/PhysRevC.24.1543](https://doi.org/10.1103/PhysRevC.24.1543). URL: <https://link.aps.org/doi/10.1103/PhysRevC.24.1543> (visited on 01/16/2024).
- [10] D. A. Dwyer and T. J. Langford. "Spectral Structure of Electron Antineutrinos from Nuclear Reactors". *Physical Review Letters* 114.1 (Jan. 7, 2015), 012502. ISSN: 0031-9007, 1079-7114. DOI: [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: [http://arxiv.org/abs/1407.1281](https://arxiv.org/abs/1407.1281) (visited on 01/16/2024).

- [11] Daya Bay Collaboration et al. "Measurement of the Reactor Antineutrino Flux and Spectrum at 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).
- [12] 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).
- [13] 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).
- [14] 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](https://arxiv.org/abs/2109.11174) (visited on 02/28/2024).
- [15] 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](https://arxiv.org/abs/2212.08502) (visited on 08/09/2024).
- [16] 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](https://arxiv.org/abs/astro-ph/0302055) (visited on 01/16/2024).
- [17] 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](https://arxiv.org/abs/2007.00314) (visited on 07/26/2023).
- [18] 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).
- [19] Photomultiplier tube R12860 | Hamamatsu Photonics. URL: https://www.hamamatsu.com/eu/en/product/optical-sensors/pmt/pmt_tube-alone/head-on-type/R12860.html (visited on 02/08/2024).
- [20] 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](https://arxiv.org/abs/2003.12212) (visited on 02/07/2024).
- [21] 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 \[hep-ex, physics:physics\]](https://arxiv.org/abs/1703.01867). URL: [http://arxiv.org/abs/1703.01867](https://arxiv.org/abs/1703.01867) (visited on 07/28/2023).
- [22] 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 \[physics\]](https://arxiv.org/abs/2103.16900). URL: [http://arxiv.org/abs/2103.16900](https://arxiv.org/abs/2103.16900) (visited on 02/07/2024).
- [23] 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 \[hep-ex, physics:physics\]](https://arxiv.org/abs/2205.08629). URL: [http://arxiv.org/abs/2205.08629](https://arxiv.org/abs/2205.08629) (visited on 02/08/2024).

- [24] Yang Han. "Dual Calorimetry for High Precision Neutrino Oscillation Measurement at JUNO Experiment". *AstroParticule et Cosmologie*, France, Paris U. VII, APC, June 2021.
- [25] R. Acquafredda 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) (visited on 02/29/2024).
- [26] 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 \[hep-ex, physics:physics\]](https://arxiv.org/abs/2011.06405). URL: [http://arxiv.org/abs/2011.06405](https://arxiv.org/abs/2011.06405) (visited on 08/10/2023).
- [27] 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 \[physics\]](https://arxiv.org/abs/2209.10387). URL: [http://arxiv.org/abs/2209.10387](https://arxiv.org/abs/2209.10387) (visited on 01/16/2024).
- [28] 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](https://doi.org/10.1088/1742-6596/898/4/042029) (visited on 02/27/2024).
- [29] 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).
- [30] 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).
- [31] 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).
- [32] Angel Abusleme et al. "Potential to Identify the Neutrino Mass Ordering with Reactor Antineutrinos in JUNO" (May 2024). eprint: [2405.18008](https://arxiv.org/abs/2405.18008).
- [33] Xiangpan Ji, Wenqiang Gu, Xin Qian, Hanyu Wei, and Chao Zhang. *Combined Neyman-Pearson Chi-square: An Improved Approximation to the Poisson-likelihood Chi-square*. arXiv.org. Mar. 17, 2019. URL: <https://arxiv.org/abs/1903.07185v3> (visited on 10/03/2024).
- [34] 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).
- [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 \[hep-ex, physics:physics\]](https://arxiv.org/abs/1812.01799). URL: [http://arxiv.org/abs/1812.01799](https://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: [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](https://arxiv.org/abs/2101.08901) (visited on 07/28/2023).
- [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*

- 3543 *A: Accelerators, Spectrometers, Detectors and Associated Equipment* 354.2 (Jan. 30, 1995), 389–399.
 3544 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).
- 3545 [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).
- 3546 [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).
- 3547 [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).
- 3548 [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](https://arxiv.org/abs/2101.04839) [hep-ex, physics:physics]. URL: [http://arxiv.org/abs/2101.04839](https://arxiv.org/abs/2101.04839) (visited on 07/24/2023).
- 3549 [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](https://arxiv.org/abs/2206.09040) [physics]. URL: [http://arxiv.org/abs/2206.09040](https://arxiv.org/abs/2206.09040) (visited on 07/24/2023).
- 3550 [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).
- 3551 [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).
- 3552 [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).
- 3553 [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).
- 3554 [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).
- 3555 [49] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. *Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering*. Feb. 5, 2017. DOI: [10.48550/arXiv.1606.09375](https://doi.org/10.48550/arXiv.1606.09375). eprint: [1606.09375](https://arxiv.org/abs/1606.09375) [cs, stat]. URL: [http://arxiv.org/abs/1606.09375](https://arxiv.org/abs/1606.09375) (visited on 04/04/2024).

- [50] JUNO Collaboration et al. "JUNO Physics and Detector". *Progress in Particle and Nuclear Physics* 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).
- [51] Leo Breiman, Jerome Friedman, R. A. Olshen, and Charles J. Stone. *Classification and Regression Trees*. New York: Chapman and Hall/CRC, Oct. 25, 2017. 368 pp. ISBN: 978-1-315-13947-0. DOI: [10.1201/9781315139470](https://doi.org/10.1201/9781315139470).
- [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: 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).
- [53] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. Jan. 29, 2017. 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).
- [54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep Residual Learning for Image Recognition". *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). ISSN: 1063-6919. June 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).
- [55] Olga Russakovsky et al. *ImageNet Large Scale Visual Recognition Challenge*. Jan. 29, 2015. DOI: [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> (visited on 05/17/2024).
- [56] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image 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: <http://arxiv.org/abs/1409.1556> (visited on 05/17/2024).
- [57] Anna Allen. *generic-github-user/Image-Convolution-Playground*. original-date: 2018-09-28T22:42:55Z. July 15, 2024. URL: <https://github.com/generic-github-user/Image-Convolution-Playground> (visited on 07/16/2024).
- [58] Jason Ansel et al. *PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation*. Publication Title: 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24) original-date: 2016-08-13T05:26:41Z. Apr. 2024. DOI: [10.1145/3620665.3640366](https://doi.org/10.1145/3620665.3640366). URL: <https://pytorch.org/assets/pytorch2-2.pdf> (visited on 07/16/2024).
- [59] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition". *Proceedings of the IEEE* 86.11 (Nov. 1998). Conference Name: Proceedings of the 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).
- [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).
- [61] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. *Neural Message Passing for Quantum Chemistry*. June 12, 2017. DOI: [10.48550/arXiv.1704.01212](https://doi.org/10.48550/arXiv.1704.01212). eprint: [1704.01212 \[cs\]](https://arxiv.org/abs/1704.01212). URL: [http://arxiv.org/abs/1704.01212](https://arxiv.org/abs/1704.01212) (visited on 05/22/2024).
- [62] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. *Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting*. Feb. 22, 2018. DOI: [10.48550/arXiv.1707.01926](https://doi.org/10.48550/arXiv.1707.01926). eprint: [1707.01926 \[cs, stat\]](https://arxiv.org/abs/1707.01926). URL: [http://arxiv.org/abs/1707.01926](https://arxiv.org/abs/1707.01926) (visited on 05/22/2024).
- [63] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. *Generative Adversarial Networks*. June 10, 2014. DOI: [10.48550/arXiv.1406.2661](https://doi.org/10.48550/arXiv.1406.2661). eprint: [1406.2661 \[cs, stat\]](https://arxiv.org/abs/1406.2661). URL: [http://arxiv.org/abs/1406.2661](https://arxiv.org/abs/1406.2661) (visited on 05/29/2024).
- [64] Anatael Cabrera et al. *Multi-Calorimetry in Light-based Neutrino Detectors*. Dec. 20, 2023. DOI: [10.48550/arXiv.2312.12991](https://doi.org/10.48550/arXiv.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).

- [65] Victor Lebrin. "Towards the Detection of Core-Collapse Supernovae Burst Neutrinos with the 3-inch PMT System of the JUNO Detector". These de doctorat. Nantes Université, Sept. 5, 2022. URL: <https://theses.fr/2022NANU4080> (visited on 05/22/2024).
- [66] Dan Cireşan, Ueli Meier, and Juergen Schmidhuber. *Multi-column Deep Neural Networks for Image Classification*. version: 1. Feb. 13, 2012. DOI: [10.48550/arXiv.1202.2745](https://doi.org/10.48550/arXiv.1202.2745). eprint: [1202.2745\[cs\]](https://arxiv.org/abs/1202.2745). URL: [http://arxiv.org/abs/1202.2745](https://arxiv.org/abs/1202.2745) (visited on 06/27/2024).
- [67] R. Abbasi et al. "A Convolutional Neural Network based Cascade Reconstruction for the Ice-Cube Neutrino Observatory". *Journal of Instrumentation* 16.7 (July 1, 2021), P07041. ISSN: 1748-0221. DOI: [10.1088/1748-0221/16/07/P07041](https://doi.org/10.1088/1748-0221/16/07/P07041). eprint: [2101.11589\[hep-ex\]](https://arxiv.org/abs/2101.11589). URL: [http://arxiv.org/abs/2101.11589](https://arxiv.org/abs/2101.11589) (visited on 06/27/2024).
- [68] D. Maksimović, M. Nieslony, and M. Wurm. "CNNs for enhanced background discrimination in DSNB searches in large-scale water-Gd detectors". *Journal of Cosmology and Astroparticle Physics* 2021.11 (Nov. 2021). Publisher: IOP Publishing, 051. ISSN: 1475-7516. DOI: [10.1088/1475-7516/2021/11/051](https://doi.org/10.1088/1475-7516/2021/11/051). URL: <https://dx.doi.org/10.1088/1475-7516/2021/11/051> (visited on 06/27/2024).
- [69] Taco S. Cohen, Mario Geiger, Jonas Koehler, and Max Welling. *Spherical CNNs*. Feb. 25, 2018. DOI: [10.48550/arXiv.1801.10130](https://doi.org/10.48550/arXiv.1801.10130). eprint: [1801.10130\[cs,stat\]](https://arxiv.org/abs/1801.10130). URL: [http://arxiv.org/abs/1801.10130](https://arxiv.org/abs/1801.10130) (visited on 07/13/2024).
- [70] NVIDIA A100 GPUs Power the Modern Data Center. NVIDIA. URL: <https://www.nvidia.com/en-gb/data-center/a100/> (visited on 08/06/2024).
- [71] NVIDIA V100. NVIDIA. URL: <https://www.nvidia.com/en-gb/data-center/v100/> (visited on 08/06/2024).
- [72] Leonard Imbert. *leonard-IMBERT/datamo*. original-date: 2023-10-17T12:37:38Z. Aug. 9, 2024. URL: <https://github.com/leonard-IMBERT/datamo> (visited on 08/09/2024).
- [73] "IEEE Standard for Floating-Point Arithmetic". *IEEE Std 754-2019 (Revision of IEEE 754-2008)* (July 2019). Conference Name: IEEE Std 754-2019 (Revision of IEEE 754-2008), 1–84. DOI: [10.1109/IEEESTD.2019.8766229](https://doi.org/10.1109/IEEESTD.2019.8766229). URL: <https://ieeexplore.ieee.org/document/8766229> (visited on 07/03/2024).
- [74] Chuanya Cao et al. "Mass production and characterization of 3-inch PMTs for the JUNO experiment". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1005 (July 2021), 165347. ISSN: 01689002. DOI: [10.1016/j.nima.2021.165347](https://doi.org/10.1016/j.nima.2021.165347). eprint: [2102.11538\[hep-ex,physics:physics\]](https://arxiv.org/abs/2102.11538). URL: [http://arxiv.org/abs/2102.11538](https://arxiv.org/abs/2102.11538) (visited on 02/08/2024).
- [75] K. M. Gorski, E. Hivon, A. J. Banday, B. D. Wandelt, F. K. Hansen, M. Reinecke, and M. Bartelman. "HEALPix – a Framework for High Resolution Discretization, and Fast Analysis of Data Distributed on the Sphere". *The Astrophysical Journal* 622.2 (Apr. 2005), 759–771. ISSN: 0004-637X, 1538-4357. DOI: [10.1086/427976](https://doi.org/10.1086/427976). eprint: [astro-ph/0409513](https://arxiv.org/abs/astro-ph/0409513). URL: [http://arxiv.org/abs/astro-ph/0409513](https://arxiv.org/abs/astro-ph/0409513) (visited on 08/10/2023).
- [76] Teng Li, Xin Xia, Xing-Tao Huang, Jia-Heng Zou, Wei-Dong Li, Tao Lin, Kun Zhang, and Zi-Yan Deng. "Design and development of JUNO event data model*". *Chinese Physics C* 41.6 (June 2017). Publisher: IOP Publishing, 066201. ISSN: 1674-1137. DOI: [10.1088/1674-1137/41/6/066201](https://doi.org/10.1088/1674-1137/41/6/066201). URL: <https://dx.doi.org/10.1088/1674-1137/41/6/066201> (visited on 08/16/2024).
- [77] 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).
- [78] Mario Schwarz, Sabrina M. Franke, Lothar Oberauer, Miriam D. Plein, Hans Th J. Steiger, and Marc Tippmann. *Measurements of the Lifetime of Orthopositronium in the LAB-Based Liquid Scintillator of JUNO*. Apr. 25, 2018. DOI: [10.1016/j.nima.2018.12.068](https://doi.org/10.1016/j.nima.2018.12.068). eprint: [1804.09456\[physics\]](https://arxiv.org/abs/1804.09456). URL: [http://arxiv.org/abs/1804.09456](https://arxiv.org/abs/1804.09456) (visited on 09/17/2024).
- [79] Narongkiat Rodphai, Zhimin Wang, Narumon Suwonjandee, and Burin Asavapibhop. "20-inch photomultiplier tube timing study for JUNO". *Journal of Physics: Conference Series* 2145.1 (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).

- 2145/1/012017. URL: <https://dx.doi.org/10.1088/1742-6596/2145/1/012017> (visited on 09/17/2024).

[80] Dong-Hao Liao et al. "Study of TTS for a 20-inch dynode PMT*". *Chinese Physics C* 41.7 (July 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).

[81] Nan Li et al. "Characterization of 3-inch photomultiplier tubes for the JUNO central detector". *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 09/17/2024).

[82] B. Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil Yogamani, and Patrick Pérez. *Deep Reinforcement Learning for Autonomous Driving: A Survey*. Jan. 23, 2021. eprint: [2002.00444\[cs\]](https://arxiv.org/abs/2002.00444). URL: [http://arxiv.org/abs/2002.00444](https://arxiv.org/abs/2002.00444) (visited on 10/02/2024).

[83] Oriol Vinyals et al. "Grandmaster level in StarCraft II using multi-agent reinforcement learning". 575.7782 (Nov. 2019). Publisher: Nature Publishing Group, 350–354. ISSN: 1476-4687. DOI: [10.1038/s41586-019-1724-z](https://doi.org/10.1038/s41586-019-1724-z). URL: <https://www.nature.com/articles/s41586-019-1724-z> (visited on 10/02/2024).

[84] Daya Bay Collaboration et al. *A high precision calibration of the nonlinear energy response at Daya Bay*. arXiv.org. Feb. 21, 2019. URL: <https://arxiv.org/abs/1902.08241v2> (visited on 10/01/2024).

[85] 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).

[86] X. B. Ma, W. L. Zhong, L. Z. Wang, Y. X. Chen, and J. Cao. "Improved calculation of the energy release in neutron-induced fission". *Physical Review C* 88.1 (July 12, 2013). Publisher: American 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).

[87] Timo Gnambs. "A Brief Note on the Standard Error of the Pearson Correlation". *Collabra: 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).

[88] "Note Sur Une Méthode de Résolution des équations Normales Provenant de L'Application de la MéThode des Moindres Carrés a un Système D'équations Linéaires en Nombre Inférieur a Celui des Inconnues. — Application de la Méthode a la Résolution D'un Système Defini 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).

[89] Pauli Virtanen et al. "SciPy 1.0: fundamental algorithms for scientific computing in Python". *Nature Methods* 17.3 (Mar. 2020). Publisher: Nature Publishing Group, 261–272. ISSN: 1548-7105. 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).

3739

3740

3741 **Titre :** Méthode Deep Learning and analyse Double Calorimétrique pour la mesure de haute
 3742 précision des paramètres d'oscillation des neutrinos dans JUNO
 3743

3744 **Mot clés :** Neutrinos; expérience JUNO; Deep Learning; reconstruction d'IBD; oscillations des
 3745 neutrinos; double calorimetrie

3746 **Résumé :** JUNO est un observatoire de
 3747 neutrinos à scintillateur liquide, polyvalent et
 3748 medium baseline (environ 52 km), situé en
 3749 Chine. Ses principaux objectifs sont de
 3750 mesurer les paramètres d'oscillation θ_{12} , Δm_{21}^2
 3751 et Δm_{31}^2 avec une précision au pour-mille
 3752 et de déterminer l'ordre des masses des
 3753 neutrinos avec un niveau de confiance de
 3754 3σ . Atteindre ces objectifs nécessite une
 3755 résolution énergétique sans précédent de
 3756 $3\%/\sqrt{E(\text{MeV})}$ avec cette technologie. Cela
 3757 demande une compréhension approfondie
 3758 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.

3772

3773 **Title:** Deep learning methods and Dual Calorimetric analysis for high precision neutrino oscil-
 3774 lation measurements at JUNO
 3775

3776 **Keywords:** Neutrinos; JUNO experiment; Deep learning; IBD reconstruction; neutrinos Oscil-
 3777 lation; dual Calorimetry

3778 **Abstract:** JUNO is a multipurpose, medium
 3779 baseline (~ 52 km) liquid scintillator neutrino
 380 observatory located in China. Its primary
 381 objectives are to measure the oscillation
 382 parameters θ_{12} , Δm_{21}^2 , and Δm_{31}^2 with per mil
 383 precision and to determine the neutrino mass
 384 ordering at a 3σ confidence level. Achiev-
 385 ing these goals requires an unprecedented
 386 energy resolution of $3\%/\sqrt{E(\text{MeV})}$ with this
 387 technology. This demands a comprehensive
 388 understanding of the various effects within the

detector. 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.

