

1

2

THÈSE DE DOCTORAT DE

NANTES UNIVERSITÉ

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

Par

Léonard Imbert

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

Thèse présentée et soutenue à Nantes, le 2 Decembre 2024
Unité de recherche : Laboratoire SUBATECH, UMR 6457

Rapporteurs avant soutenance :

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

Composition du Jury :

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

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 Neutrinos physics in JUNO	12
¹⁶	2.1.1 Reactor neutrino oscillation for NMO and precise measurements	12
¹⁷	2.1.2 Other physics	15
¹⁸	2.2 The JUNO detector	17
¹⁹	2.2.1 Detection principle	17
²⁰	2.2.2 Central Detector (CD)	19
²¹	2.2.3 Veto detector	23
²²	2.3 Calibration strategy	23
²³	2.3.1 Energy scale calibration	24
²⁴	2.3.2 Calibration system	25
²⁵	2.3.3 Instrumental non-linearity calibration	25
²⁶	2.4 Satellite detectors	26
²⁷	2.4.1 TAO	26
²⁸	2.4.2 OSIRIS	26
²⁹	2.5 Software	27
³⁰	2.6 State of the art of the Offline IBD reconstruction in JUNO	28
³¹	2.6.1 Interaction vertex reconstruction	28
³²	2.6.2 Energy reconstruction	33
³³	2.6.3 SPMT reconstruction	36
³⁴	2.6.4 Machine learning for reconstruction	36

35	2.7 JUNO sensitivity to NMO and precise measurements	38
36	2.7.1 Theoretical spectrum	38
37	2.7.2 Fitting procedure	39
38	2.7.3 Physics results	39
39	2.8 Summary	40
40	3 Machine learning: Introduction to the methods and algorithms used in this thesis	41
41	3.1 Core concepts in machine learning and neural networks	42
42	3.1.1 Boosted Decision Tree (BDT)	42
43	3.1.2 Artificial Neural Network (NN)	42
44	3.1.3 Training procedure	44
45	3.1.4 Potential pitfalls	47
46	3.2 Neural networks architectures	50
47	3.2.1 Fully Connected Deep Neural Network (FCDNN)	50
48	3.2.2 Convolutional Neural Network (CNN)	50
49	3.2.3 Graph Neural Network (GNN)	52
50	3.2.4 Adversarial Neural Network (ANN)	54
51	4 Image recognition for IBD reconstruction with the SPMT system	55
52	4.1 Method and model	56
53	4.1.1 Model	57
54	4.1.2 Data representation	58
55	4.1.3 Dataset	60
56	4.1.4 Data characteristics	61
57	4.2 Training	63
58	4.3 Results	63
59	4.3.1 J21 results	64
60	4.3.2 J21 Combination of classic and ML estimator	66
61	4.3.3 J23 results	68
62	4.4 Conclusion and prospect	70
63	5 Graph representation of JUNO for IBD reconstruction	73
64	5.1 Data representation	74
65	5.2 Message passing algorithm	76
66	5.3 Data	78
67	5.4 Model	80
68	5.5 Training	80
69	5.6 Optimization	82
70	5.6.1 Software optimization	82
71	5.6.2 Hyperparameters optimization	83
72	5.7 performance of the final version	83
73	5.8 Conclusion	87
74	6 Reliability of machine learning methods	91

75	6.1 Motivations	92
76	6.2 Method	92
77	6.3 Architecture	92
78	6.3.1 Adversarial Neural Network	92
79	6.3.2 Reconstruction Network	93
80	6.3.3 Training	93
81	6.4 Results	93
82	6.4.1 Back to identity	94
83	6.4.2 Breaking of the reconstruction	94
84	6.5 Conclusion and prospect	94
85	7 Joint fit between the SPMT and LPMT spectra	95
86	7.1 Motivations	96
87	7.1.1 Discrepancies between the SPMT and LPMT results	96
88	7.1.2 Charge Non-Linearity (QNL)	97
89	7.2 Approach	98
90	7.2.1 Data production	98
91	7.2.2 Individual fits	99
92	7.2.3 Joint fit	100
93	7.2.4 Data and theoretical spectrum generation	102
94	7.2.5 Limitations	102
95	7.3 Fit software	103
96	7.3.1 IBD generator	103
97	7.3.2 Fit	105
98	7.4 Technical challenges and development	105
99	7.5 Results	106
100	7.5.1 Validation	106
101	7.5.2 Covariance matrix	110
102	7.5.3 Statistical tests	114
103	7.6 Conclusion and perspectives	116
104	8 Conclusion	121
105	A Calculation of optimal α for estimator combination	123
106	A.1 Unbiased estimator	123
107	A.2 Optimal variance estimator	123
108	B Charge spherical harmonics analysis	125
109	C Additional spectrum smearing	133
110	D Correction of E_{vis} bias	135
111	List of Tables	137

112	List of Figures	145
113	List of Abbreviations	147
114	Bibliography	149

¹¹⁵ **Remerciements**

¹¹⁶ **Introduction**

¹¹⁷ **Chapter 1**

¹¹⁸ **Neutrino physics**

¹¹⁹ *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.*

¹²⁰ **Contents**

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

¹³¹ **1.1 Standard model**

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

¹³² **1.1.1 Limits of the standard model**

¹³³ **1.2 Historic of the neutrino**

¹³⁴ **First theories**

¹³⁵ **Discovery**

¹³⁶ **Milestones and anomalies**

¹³⁷ **1.3 Oscillation**

¹³⁸ **1.3.1 Phenomologies**

¹³⁹ **1.4 Open questions**

¹⁴⁰ **Chapter 2**

¹⁴¹ **The JUNO experiment**

¹⁴² “*Ave Juno, rosae rosam, et spiritus rex*”. It means nothing but I found it in tone.

¹⁴³ **Contents**

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

¹⁷⁰ ¹⁷² 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.

¹⁷³ ¹⁷⁵ JUNO is a neutrino detection experiment under construction located in China, in Guangdong province, near the city of Kaiping. Its main objectives are the determination of the mass ordering at the

180 3-4 σ level in 6 years of data taking and the measurement at the sub-percent precision of the oscillation
 181 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

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

188 2.1 Neutrinos physics in JUNO

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

194 2.1.1 Reactor neutrino oscillation for NMO and precise measurements

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

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

195 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
 196 to the NMO in the dependency to Δm_{32}^2 and Δm_{31}^2 causing a phase shift of the spectrum as we can
 197 see in the figure 2.2. By carefully adjusting a theoretical spectrum to the data, one can extract the
 198 NMO and the oscillation parameters. The statistic procedure used to adjust the theoretical spectrum
 199 is reviewed in more details in the section 2.7. To reach the desired sensitivity, JUNO must meet
 200 multiple requirements but most notably:

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

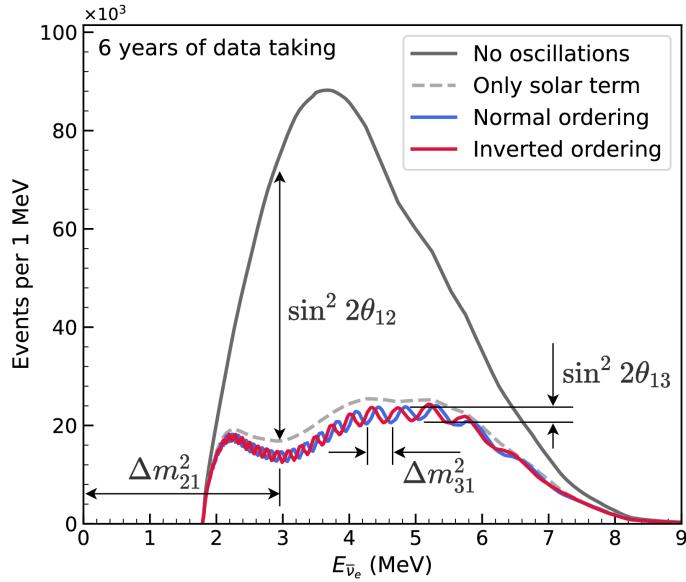


FIGURE 2.2 – Expected number of neutrinos event per MeV in JUNO after 6 years of data taking. The black curve shows the flux if there was no oscillation. The light gray curve shows the oscillation if only the solar terms are taken in account (θ_{12} , Δm_{21}^2). The blue and red curve shows the spectrum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

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

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

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

212 The $\bar{\nu}_e$ coming from reactors are emitted from β -decay of unstable fission fragments. The Taishan
213 and Yangjiang reactors are Pressurised Water Reactor (PWR), the same type as Daya Bay. In those
214 type of reactor more than 99.7 % and $\bar{\nu}_e$ are produced by the fissions of four fuel isotopes ^{235}U , ^{238}U ,
215 ^{239}Pu and ^{241}Pu . The neutrino flux per fission of each isotope is determined by the inversion of the
216 measured β spectra of fission product [4–8] or by calculation using the nuclear databases [9, 10].

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

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

218 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
219 thermal energy released in each fission and $S_i(e_\nu)$ the neutrino flux per fission for this isotope. Using
220 this method, the flux uncertainty is expected to be of an order of 2-3 % [11].

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.

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

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

229 Background in the neutrinos reactor spectrum

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

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

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

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

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

241 Identification of the mass ordering

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

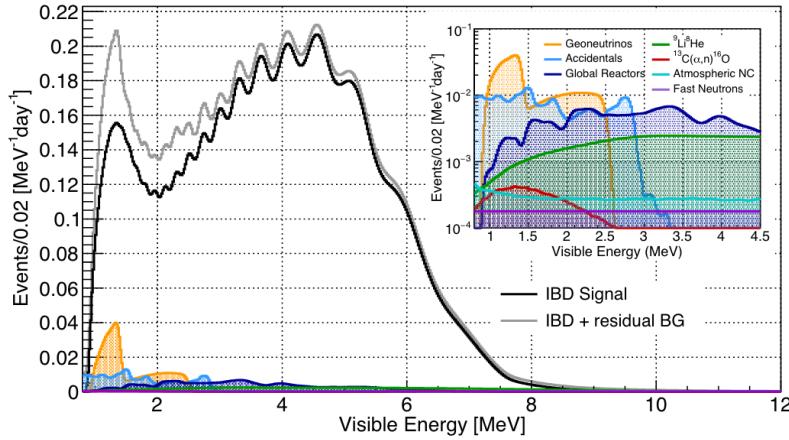


FIGURE 2.3 – Expected visible energy spectrum measured with the LPMT system with (grey) and without (black) backgrounds. The background amount for about 7% of the IBD candidate and are mostly localized below 3 MeV [11]

248 Precise measurement of the oscillations parameters

249 The oscillations parameters θ_{12} , θ_{13} , Δm_{21}^2 , Δm_{31}^2 are free parameters in the fit of the oscillation
 250 spectrum. The precision on those parameters have been estimated and are shown in table 2.2. Wee
 251 see that for θ_{12} , Δm_{21}^2 , Δm_{31}^2 , precision at 6 years is better than the reference precision by an order of
 252 magnitude [11]

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

253 2.1.2 Other physics

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

257 Geoneutrinos

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

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

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

263 Atmospheric neutrinos

264 Atmospheric neutrinos are neutrinos originating from the decay of π and K particles that are pro-
 265 duced in extensive air showers initiated by the interactions of cosmic rays with the Earth atmosphere.
 266 Earth is mostly transparent to neutrinos below the PeV energy, thus JUNO will be able to see neu-
 267 trinos coming from all directions. Their baseline range is large (15km \sim 13000km), they can have
 268 energy between 0.1 GeV and 10 TeV and will contain all neutrino and antineutrinos flavour. Their
 269 studies is complementary to the reactor antineutrinos and can help refine the constraints on the NMO
 270 [2].

271 Supernovae burst neutrinos

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

279 Diffuse supernovae neutrinos background

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

288 Beyond standard model neutrinos interactions

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

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

297 **Proton decay**

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

306 **2.2 The JUNO detector**

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

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

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

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

331 **2.2.1 Detection principle**

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

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

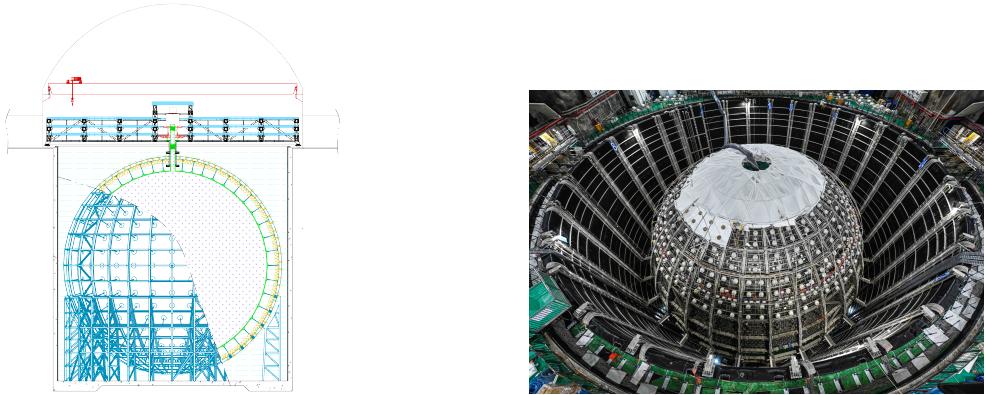


FIGURE 2.4

332 Kinematics calculation shows that this interaction has an energy threshold for the $\bar{\nu}_e$ of $(m_n + m_e -$
 333 $m_p) \approx 1.806$ MeV [19]. This threshold make the experiment blind to very low energy neutrinos.
 334 The residual energy $E_\nu - 1.806$ MeV is be distributed as kinetic energy between the positron and the
 335 neutron. The energy of the emitted positron E_e is given by [19]

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

336 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
 337 equation that the positron energy is strongly correlated to the neutrino energy.

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

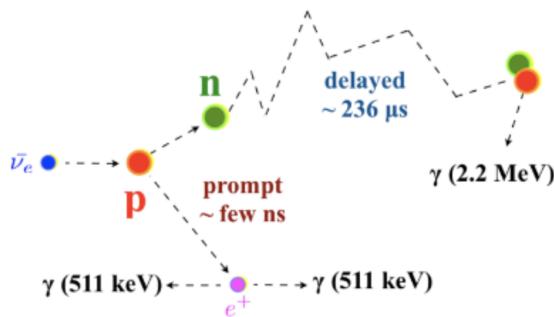


FIGURE 2.5 – Schematics of an IBD interaction in the central detector of JUNO

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

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

350 2.2.2 Central Detector (CD)

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

356 Acrylic vessel

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

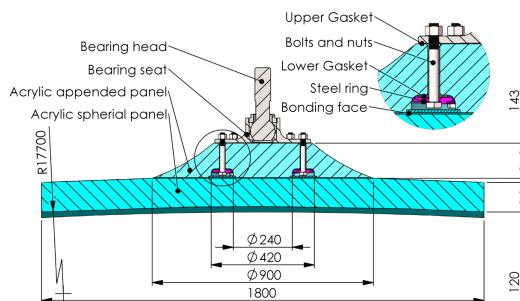


FIGURE 2.6 – Schematics of the supporting node for the acrylic vessel

367 Liquid scintillator

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

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

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

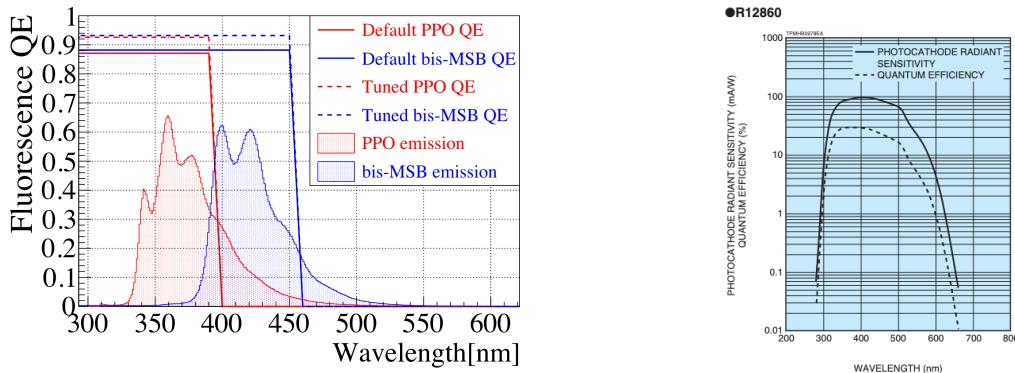


FIGURE 2.7 – On the left: Quantum efficiency (QE) and emission spectrum of the LAB and the bis-MSB [20]. On the right: Sensitivity of the Hamamatsu LPMT depending on the wavelength of the incident photons [22].

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

- 387 — A light yield / MeV of the amount of 10^4 photons to maximize the statistic in the energy
 388 measurement.
- 389 — An attenuation length comparable to the size of the detector to prevent losing photons during
 390 their propagation in the LS. The final attenuation length is 25.8m [24] to compare with the CD
 391 diameter of 35.4m.
- 392 — Uranium/Thorium radiopurity to prevent background signal. The reactor neutrino program
 393 require a contamination fraction $F < 10^{-15}$ while the solar neutrino program require $F <$
 394 10^{-17} .

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

398 Large Photo-Multipliers Tubes (LPMTs)

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

406 The Large Photo-Multipliers Tubes (LPMT), used in the central detector and in the water pool, are
 407 20-inch (50.8 cm) radius PMTs. ~ 5000 dynode-PMTs [22] were produced by the Hamamatsu[®]
 408 company and ~ 15000 Micro-Channel Plate (MCP) [26] by the NNVT[®] company. This system is
 409 the one responsible for the energy measurement with a energy resolution of $3\%/\sqrt{E}$, resolution

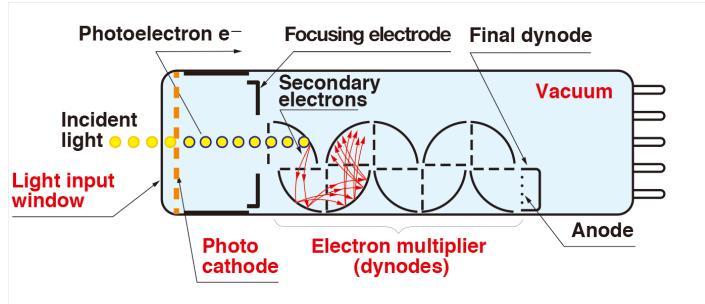


FIGURE 2.8 – Schematic of a PMT

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

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

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

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

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

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

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

436 Small Photo-Multipliers Tubes (SPMTs)

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

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

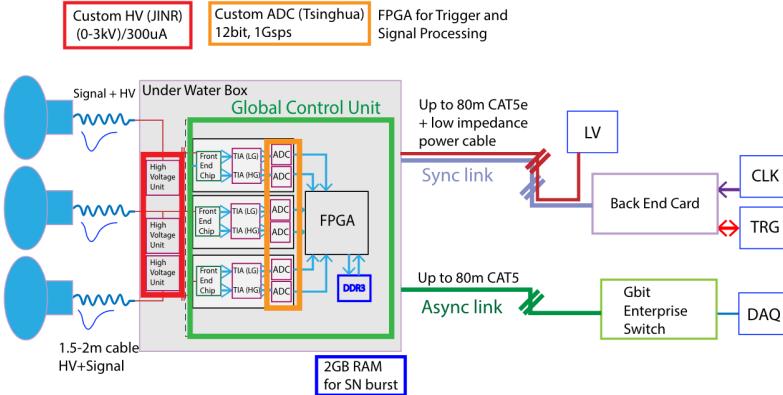


FIGURE 2.9 – The LPMT electronics scheme. It is composed of two part, the *wet* electronics on the left, located underwater and the *dry* electronics on the right. They are connected by Ethernet cable for data transmission and a dedicated low impedance cable for power distribution

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

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

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

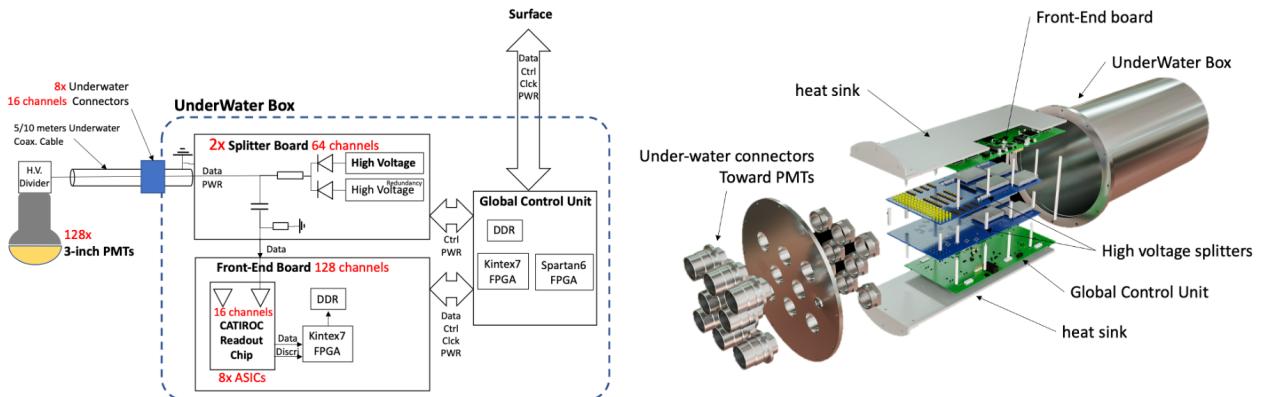


FIGURE 2.10 – Schematic of the JUNO SPMT electronic system (left), and exploded view of the main component of the UWB (right)

458 2.2.3 Veto detector

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

463 **Cherenkov in water pool**

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

472 **Top tracker**

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

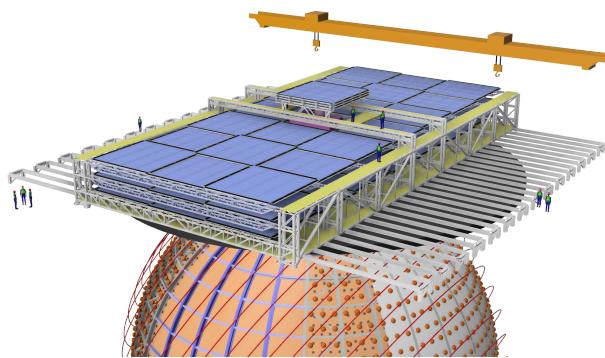


FIGURE 2.11 – The JUNO top tracker

479

480 2.3 Calibration strategy

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

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

2.3.1 Energy scale calibration

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

Sources / Processes	Type	Radiation
^{137}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 2.4 – List of sources and their process considered for the energy scale calibration

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

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

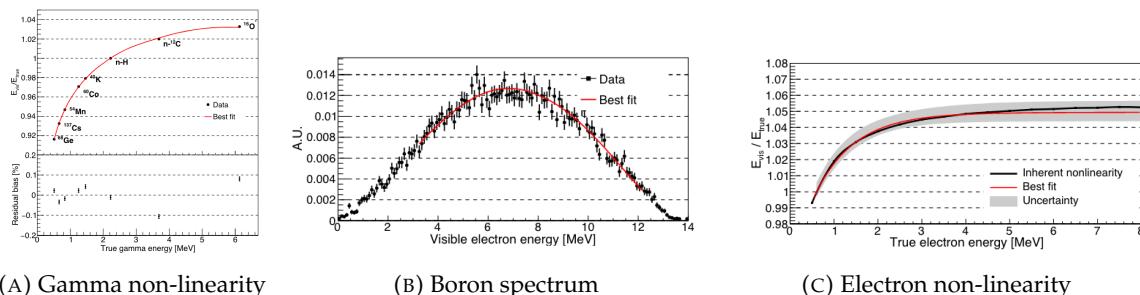


FIGURE 2.12 – Fitted and simulated non linearity of gamma, electron sources and from the ^{12}B spectrum. Black points are simulated data. Red curves are the best fits. Figures taken from [29].

504 **2.3.2 Calibration system**

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

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

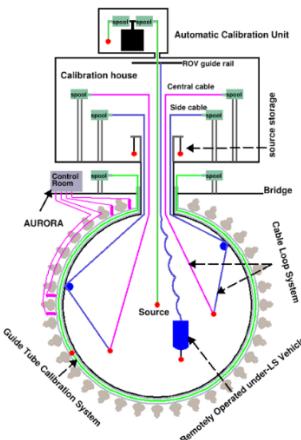


FIGURE 2.13 – Overview of the calibration system

519 The preliminary calibration program is depicted in table 2.5.

520 **2.3.3 Instrumental non-linearity calibration**

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

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

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

TABLE 2.5 – Calibration program of the JUNO experiment

531 2.4 Satellite detectors

532 As introduced in section 2.1.1 and section 2.2.2, 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 2.4.1 TAO

537 The Taishan Antineutrino Observatory (TAO) [12, 30] 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 [13]. 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 2.4.2 OSIRIS

552 The Online Scintillator Internal Radioactivity Investigation System (OSIRIS) [25] is an ultralow back-
 553 ground, 20 m³ 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

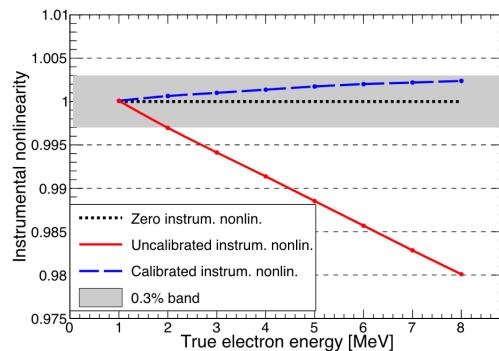


FIGURE 2.14 – Event-level instrumental non-linearity, defined as the ratio of the total measured LPMT charge to the true charge for events uniformly distributed in the detector. The solid red line represents event-level non-linearity without the channel-level correction, with position non-uniformity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

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 contain 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.

565 **2.5 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 [31] framework:

- 569 — Various primary particles simulators for the different kind of events, background and calibra-
 570 tion sources.
- 571 — A Geant4 [32–34] 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 Tran-

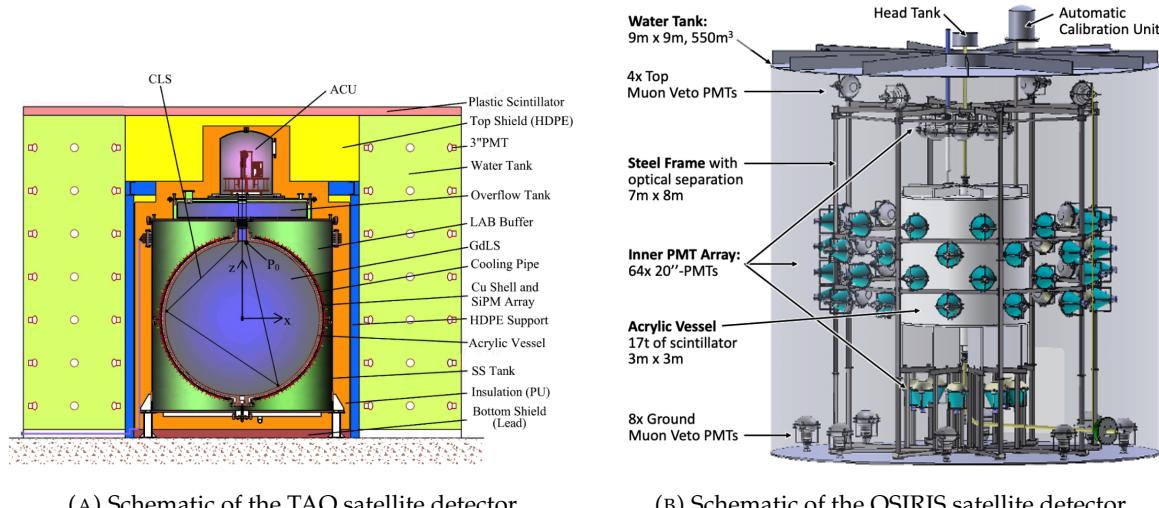


FIGURE 2.15

579 sit Spread (TTS), pre-pulsing, after-pulsing and ringing if 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 2.6 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.

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

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

601 2.6.1 Interaction vertex reconstruction

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

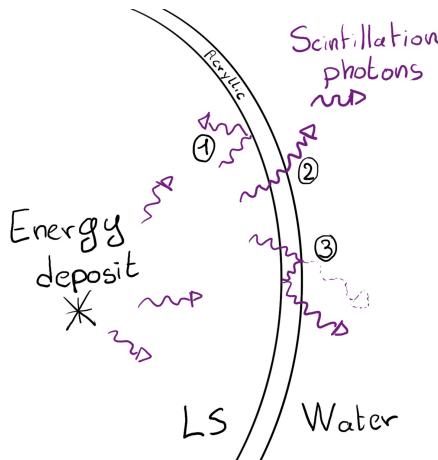
604 **Charge based algorithm**

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

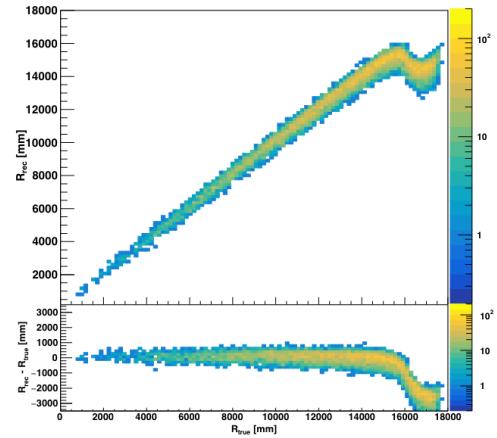
$$\vec{r}_{cb} = a \cdot \frac{\sum_i q_i \cdot \vec{r}_i}{\sum_i q_i} \quad (2.3)$$

606 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
 607 reconstructed interaction position. a is a scale factor introduced because a weighted average over
 608 a 3D sphere is inherently biased. Using calibration we can estimate $a \approx 1.3$ [37]. The results in
 609 figure 2.16b shows that the reconstruction is biased from around 15m and further. This is due to the
 610 phenomena called “total reflection area” or TR Area.

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



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



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

FIGURE 2.16

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

619 **Time based algorithm**

620 The time based algorithm use the distribution of the time of flight corrections Δt (Eq 2.4) of an event
 621 to reconstruct its vertex and t_0 . It follow the following iterations:

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

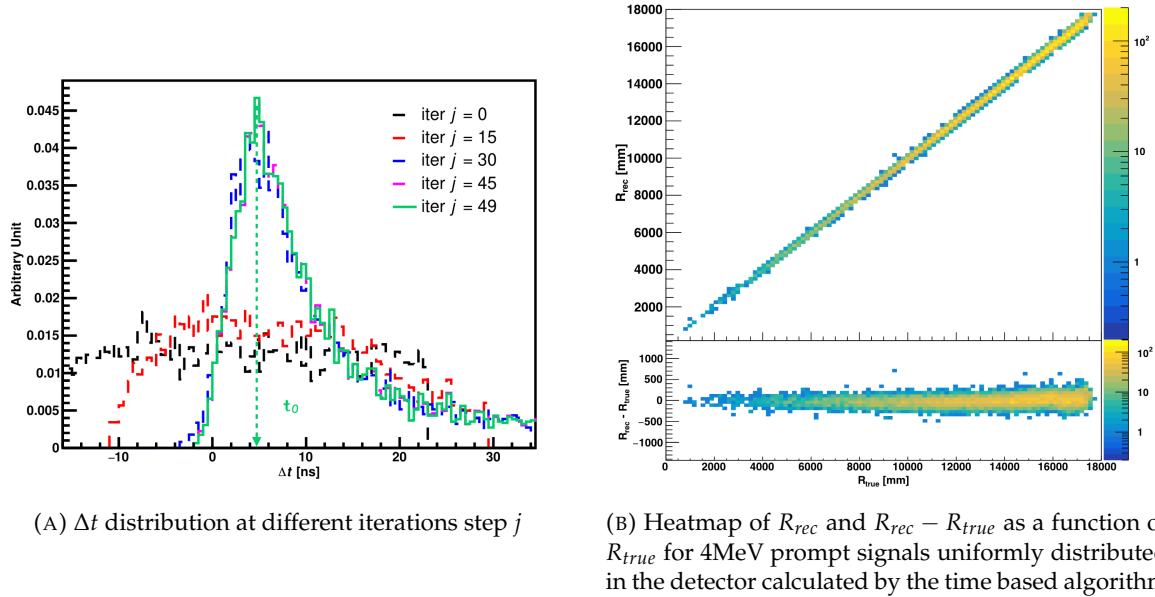


FIGURE 2.17

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

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

624 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
625 photon considering an rectilinear trajectory and an effective velocity in the LS and water (see
626 [37] for detailed description of this effective velocity). Plot the Δt distribution and label the
627 peak position as Δt^{peak} (see fig 2.17a).

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

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

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

633 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
634 step 2.

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

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

638 The parameters (p_0, p_1, p_2) were optimized to (9.42, 0.74, -4.60) for Hamamatsu PMTs and (41.31,
639 -12.04, -20.02) for NNVT PMTs [37]. The results presented in figure 2.17b shows that the time based
640 algorithm provide a more accurate vertex and is unbiased even in the TR area. This results (\vec{r}_0, t_0) is
641 used as initial value for the likelihood algorithm.

642 **Time likelihood algorithm**

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

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

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

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

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

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

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

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

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

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

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

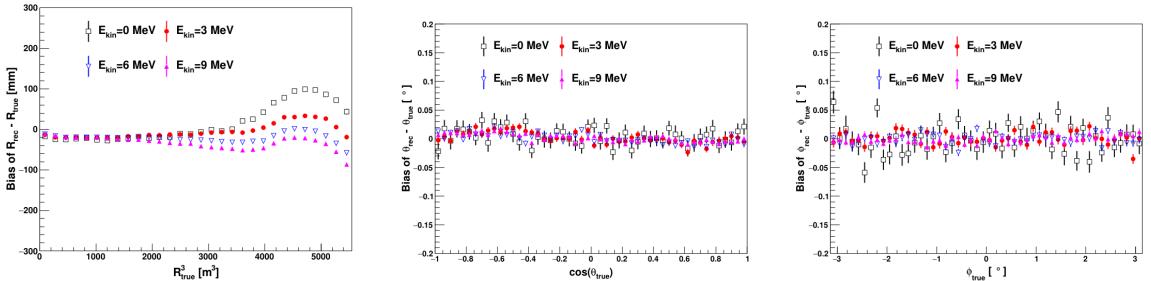


FIGURE 2.18 – Bias of the reconstructed radius R (left), θ (middle) and ϕ (right) for multiple energies by the time likelihood algorithm

658 **Charge likelihood algorithm**

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

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

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

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

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

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

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

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

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

668 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
 669 efficiency, $f(\theta_i)$ its angular response, ζ_m is the attenuation length in the materials and δ_i the expected
 670 number of dark noise.

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

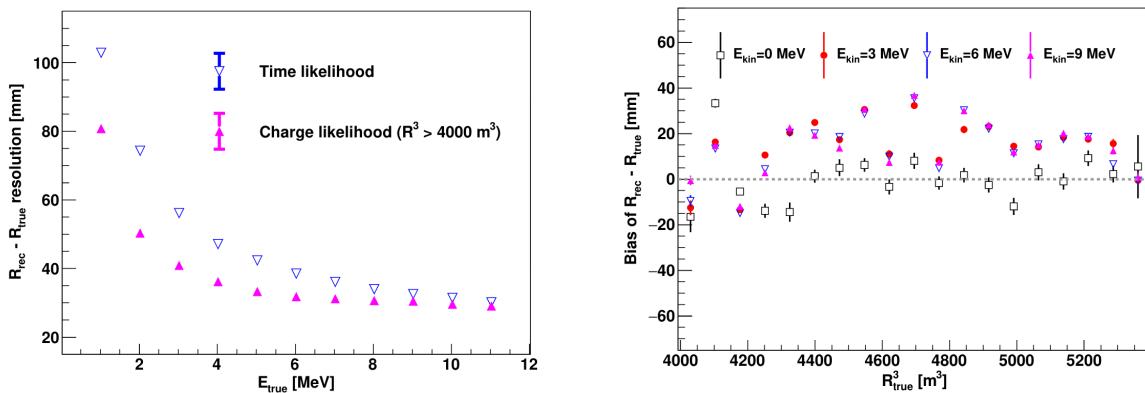


FIGURE 2.19 – On the left: Resolution of the reconstructed R as a function of the energy in the TR area ($R^3 > 4000 \text{ m}^3 \equiv R > 16 \text{ m}$) by the charge and time likelihood algorithms. On the right: Bias of the reconstructed R in the TR area for different energies by the charge likelihood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

676 resolution than the time likelihood. The figure 2.20 shows the radial resolution of the different
 677

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

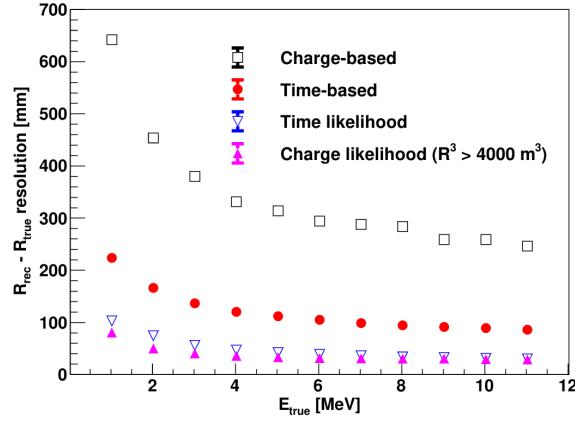
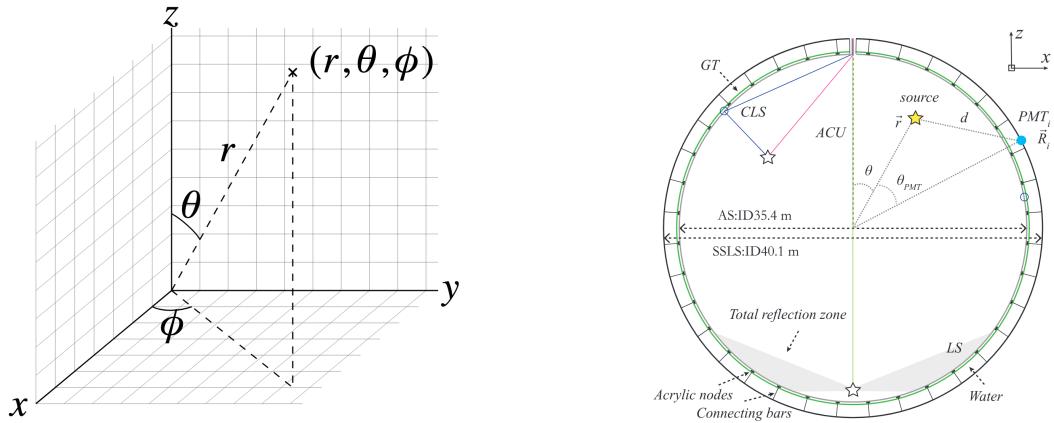


FIGURE 2.20 – Radial resolution of the different vertex reconstruction algorithms as a function of the energy

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

682 2.6.2 Energy reconstruction

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



(A) Spherical coordinate system used in JUNO for reconstruction

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

FIGURE 2.21

687 **Charge estimation**

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

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

694 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
 695 is the expected number of dark noise photoelectrons in the time window L . The time window have
 696 been optimized to $L = 280$ ns [41]. \bar{q}_i is the average recorded photoelectrons in the time window
 697 and \hat{Q}_i is the expected average charge for 1 photoelectron. The N_{pe} map is constructed following the
 698 procedure described in [36].

699 **Time estimation**

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

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

702 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
 703 by a gaussian

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

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

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

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

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

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

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

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

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

⁷¹⁷ and so for the case where only one photon is detected by the PMT ($k = 1$)

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

⁷¹⁸ where $P(k_\alpha, \mu_\alpha)$ is the Poisson probability to detect k_α PE from $\alpha \in \{l, d\}$ with the condition $k_l + k_d =$
⁷¹⁹ k .

⁷²⁰ Now that we have the individual timing and charge probability we can construct the charge likelihood referred as QMLE:
⁷²¹

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

⁷²² 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
⁷²³ charge pdf for k PE. And we can also construct the time likelihood referred as TMLE:

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

⁷²⁴ where K is cut to 20 PE and hit is the set of hits satisfying $-100 < t_{i,r} < 500$ ns.

⁷²⁵ Merging those two likelihood give the charge-time likelihood QTMLLE

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

⁷²⁶ The radial and energy resolutions of the different likelihood are presented in figure 2.22 (from [41]).
⁷²⁷ We can see the improvement of adding the time information to the vertex reconstruction and that
⁷²⁸ an increase in vertex precision can bring improvement in the energy resolution, especially at low
⁷²⁹ energies.

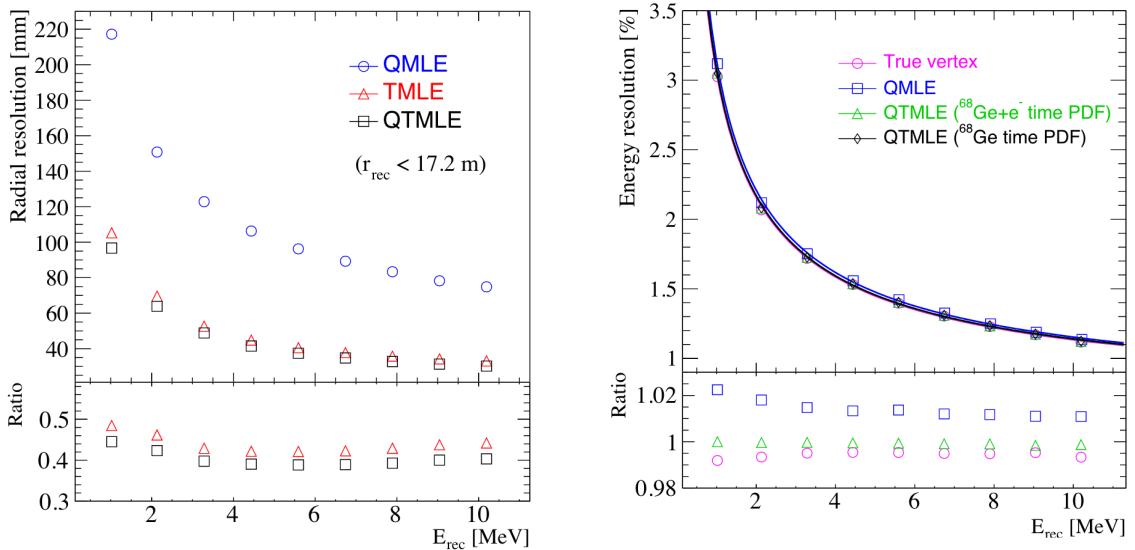


FIGURE 2.22

730 **2.6.3 SPMT reconstruction**

731 **TODO**

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

735 **2.6.4 Machine learning for reconstruction**

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

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

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

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

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

754 **Boosted Decision Tree (BDT)**

755 On of the most classic ML method used in physics in last years is the Boosted Decision Tree (see
 756 Chapter 3.1.1). They have been explored for vertex reconstruction [42] et for energy reconstruction
 757 [42, 43].

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

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 2.6 – Features used by the BDT for vertex reconstruction

760 Its reconstruction performances are presented in figure 2.24.

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

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

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

763 Neural Network (NN)

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

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

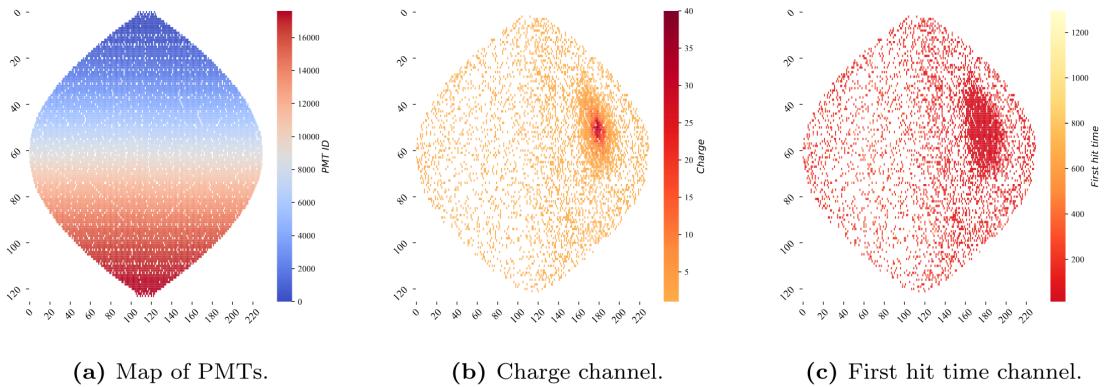


FIGURE 2.23 – Projection of the LPMTs in JUNO on a 2D plane. (a) Show the distribution of all PMTs and (b) and (c) are example of what the charge and time channel looks like respectively

774 Using 2D have the upside of encoding a large part of the informations structurally but loose the rotat-
 775ional invariance of the detector. It also give undefined information to the neural network (what is a
 776 pixel without PMT ? What should be its charge and time ?), cause deformation in the representation
 777 of the detector (sides of projection) and loose topological informations.

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

782 aggregated informations from the PMTs it contains. The edges contains geographic informations of
 783 the nodes relative positions.

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

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

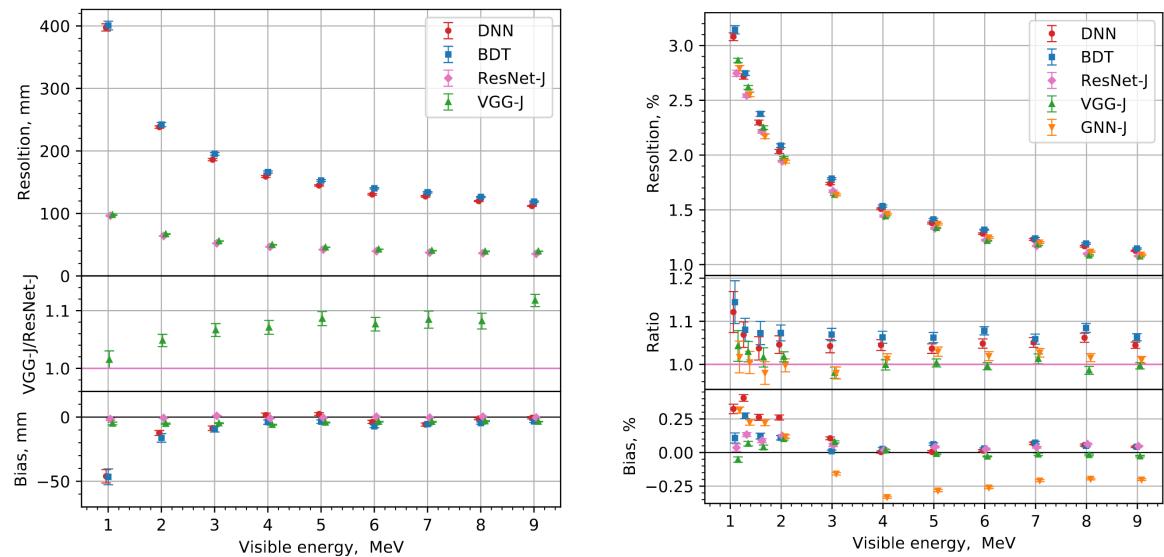


FIGURE 2.24 – Radial (left) and energy (right) resolutions of different ML algorithms.
 The results presented here are from [42]. DNN is a deep neural network, BDT is a BDT,
 ResNet-J and VGG-J are CNN and GNN-J is a GNN.

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

793 2.7 JUNO sensitivity to NMO and precise measurements

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

799 The following details about JUNO measurement is common to the three analysis. The details and
 800 specific of the joint analysis are detailed in Chapter 7.

801 2.7.1 Theoretical spectrum

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

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

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

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

2.7.2 Fitting procedure

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

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

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

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

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

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

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

The χ^2 is minimized by exploring the parameter phase space via gradient descent.

2.7.3 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 [11]. The expected sensitivity to mass ordering is 3σ after 6.5 years [50].

2.8 Summary

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

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

⁸⁵⁵ **Chapter 3**

⁸⁵⁶ **Machine learning: Introduction to the
methods and algorithms used in this
thesis**

⁸⁵⁷

⁸⁵⁸

⁸⁵⁹ “I have the shape of a human being and organs equivalent to those of a
human being. My organs, in fact, are identical to some of those in a
prostheticized human being. I have contributed artistically, literally, and
scientifically to human culture as much as any human being now
alive. What more can one ask?”

Isaac Asimov, The Complete Robot

⁸⁶⁰ **Contents**

⁸⁶¹ 3.1 Core concepts in machine learning and neural networks	⁸⁶² 42
⁸⁶³ 3.1.1 Boosted Decision Tree (BDT)	⁸⁶⁴ 42
⁸⁶⁵ 3.1.2 Artificial Neural Network (NN)	⁸⁶⁶ 42
⁸⁶⁷ 3.1.3 Training procedure	⁸⁶⁸ 44
⁸⁶⁹ 3.1.4 Potential pitfalls	⁸⁷⁰ 47
⁸⁷¹ 3.2 Neural networks architectures	⁸⁷² 50
⁸⁷³ 3.2.1 Fully Connected Deep Neural Network (FCDNN)	⁸⁷⁴ 50
⁸⁷⁵ 3.2.2 Convolutional Neural Network (CNN)	⁸⁷⁶ 50
⁸⁷⁷ 3.2.3 Graph Neural Network (GNN)	⁸⁷⁸ 52
⁸⁷⁹ 3.2.4 Adversarial Neural Network (ANN)	⁸⁸⁰ 54

⁸⁷³ Machine Learning (ML) and more specifically Neural Network (NN) are families of data-driven
⁸⁷⁴ algorithms. They are used in a wide variety of domains including natural language processing,
⁸⁷⁵ computer vision, speech recognition and, the subject of this thesis, scientific studies.

⁸⁷⁶ They are used to model complex distributions from a finite dataset to extract a generalist behavior.
⁸⁷⁷ For example, in our case, it could be an algorithm that would differentiate the nature of a particle
⁸⁷⁸ interacting in the liquid scintillator, between a positron and an electron, based on the readout charge
⁸⁷⁹ and time (Q, t) of the 17612 LPMT of the JUNO experiment. During a first training phase, it would
⁸⁸⁰ learn the discriminative features between the two in the 35224-dimensional charge and time distri-
⁸⁸¹ bution, built from samples of e^+ and e^- events.

⁸⁸² It would learn to derive from a complex, highly dimensional set of data the essential few informations
⁸⁸³ characterizing the interactions: a three body energy deposition (the positron and two annihilation
⁸⁸⁴ gammas) and the single deposit from an electron.

⁸⁸⁵ Ideally, the algorithm would learn to recognize those informations on its own, regardless of the input
⁸⁸⁶ size and complexity. In practice, however, these algorithms are guided by human design through

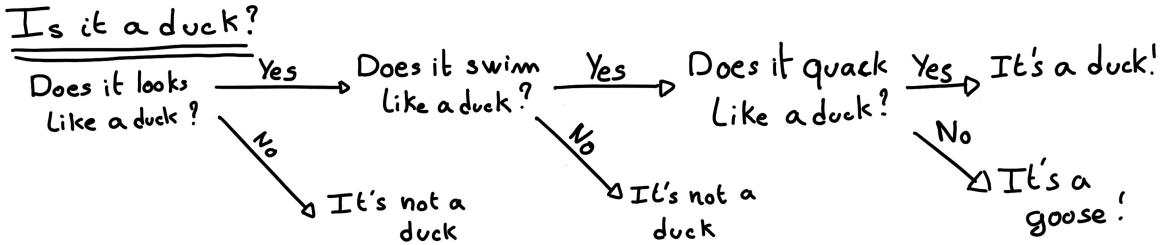


FIGURE 3.1 – Example of a BDT that determine if the given object is a duck

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

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

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

3.1 Core concepts in machine learning and neural networks

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

3.1.1 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

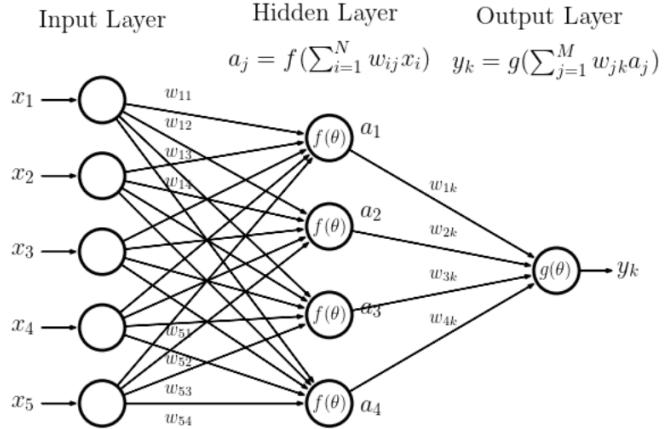


FIGURE 3.2 – Schema of a simple neural network

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

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

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

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

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

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

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

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

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

945

3.1.3 Training procedure

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

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

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

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

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

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

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

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

959 This induce \mathcal{L} needs to be differentiable with respect to θ , thus the layers and their activation functions
 960 also need to be differentiable. This simple gradient descent, designated as Stochastic Gradient
 961 Descent (SGD), can be extended with first and second order momentums like in the Adam optimizer
 962 [53]. More details about the optimizers can be found in section 3.1.3.

963

Training lifecycle

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

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

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

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

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

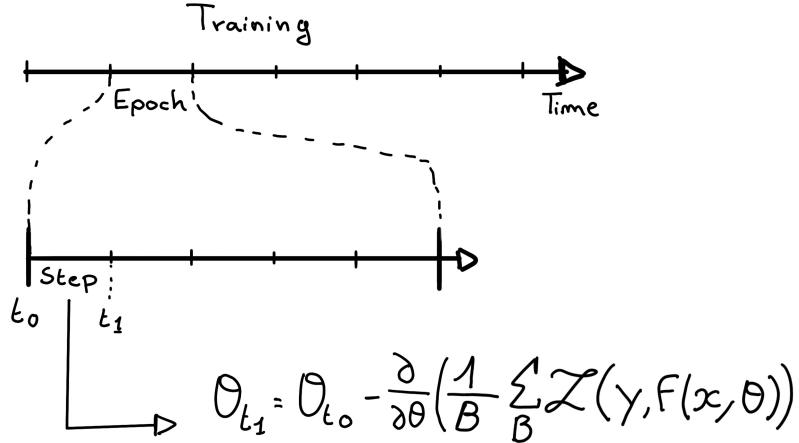


FIGURE 3.3 – Illustration of the training lifecycle

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

986 The optimizer

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

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

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

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

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

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

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

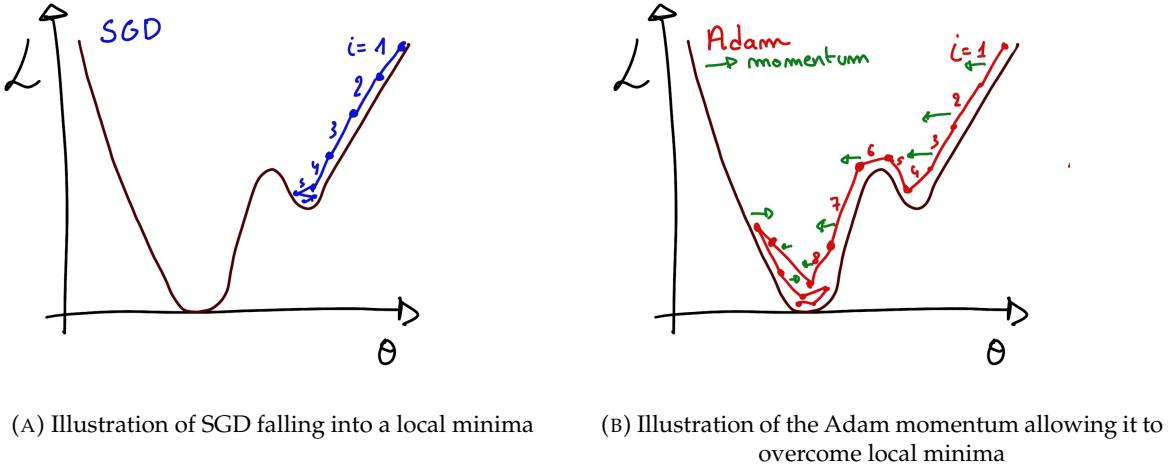


FIGURE 3.4

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

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

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

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

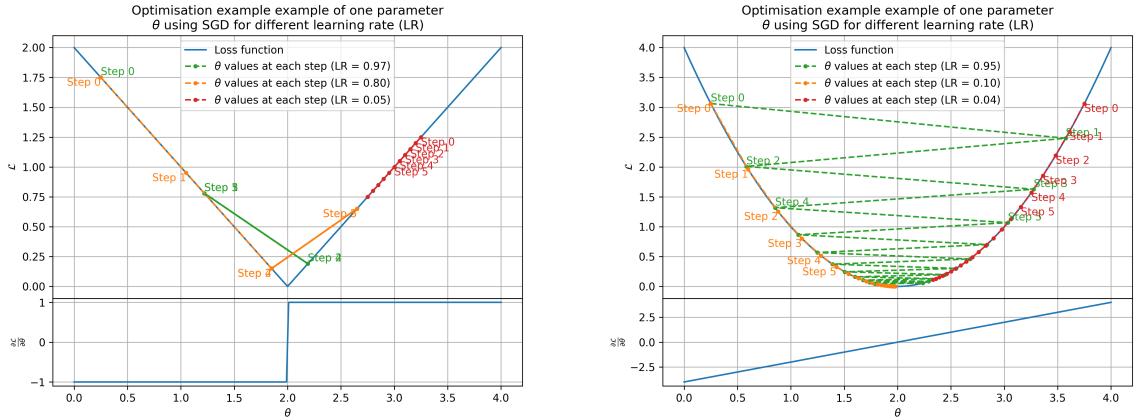
where ϵ is a small number to prevent divergence when v is close to 0. These momentums allow to overcome small local minima in the parameters phase. Imagine ball going down a slope as illustrated in 3.4a, if you ignore the stored momentum you get SGD and get stuck as on the left plot. Now if you consider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

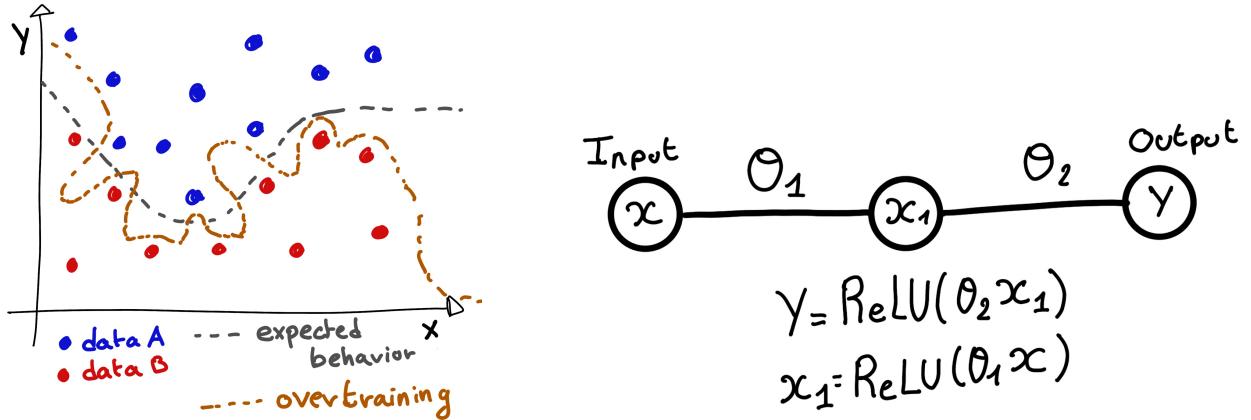
3.1.4 Potential pitfalls

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

1040 Gradient vanishing

1041 Gradient vanishing is the effect of the gradient being so small for the early layers that the parameters
 1042 are barely updated after each step. This cause the network to be unable to converge to the minima.

1043 This comes from the way the gradient descent is calculated. Imagine a simple network composed of
 1044 three fully connected layers: the input layer, a intermediate layer and the output layer. Let L be the
 1045 loss, θ_1 the parameter between the input and the intermediate layer and θ_2 the parameter between
 1046 the intermediate and output layer. This network is schematized in figure 3.6b.

1047 The gradient for θ_1 will be computed using the chain rule presented in equation 3.6. Because θ_1
 1048 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
 1049 much more layer, we can see how the subsequent multiplication of small gradients would lead to
 1050 very small update of the parameters thus "vanishing gradient".

1051 Multiple actions can be taken to prevent this effect such as:

- 1052 — **Batch normalization:** In this case we apply a normalization layer that will normalize the data.
 1053 It means that we transform the input variable X into a variable D which distribution follow
 1054 $\langle D \rangle = 0$ and $\sigma_D = 1$. This helps the parameters of the network to maintain an appropriate
 1055 scale.
- 1056 — **Residual Network (ResNet)** [54]: Residual network is a technique for neural network in
 1057 which, instead of just sequentially feeding the results of each layer to the next one, you
 1058 compute a residual over the input data. This technique is illustrated in figure 3.7. The
 1059 reference [54] show empirical evidence of its relevance.

1060 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

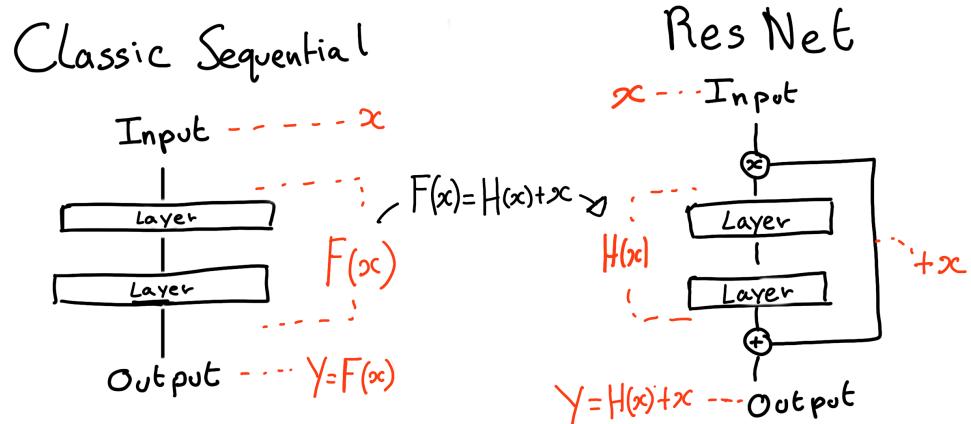


FIGURE 3.7 – Illustration of the ResNet framework

follow

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

1061 The explosion is illustrated in figure 3.8 where we can see that the loss degrades with each step of
 1062 optimization. In this illustration it is clear that reducing the learning rate suffice but this behaviour
 1063 can happens in the middle of the training where the learning rate schedule does not permit reactivity.

1064 There exist solutions to prevent this explosions:

- 1065 — **Gradient clipping:** Is this case we work on the gradient so that the norm of gradient vector
 1066 does not exceed a certain threshold. In our illustration in figure 3.8 the gradient for $\theta > 0$
 1067 could be clipped at 3 for example.
 1068 — **Batch normalization:** For the same reasons as for gradient vanishing, normalizing the input
 1069 data help reduce erratic behaviour.

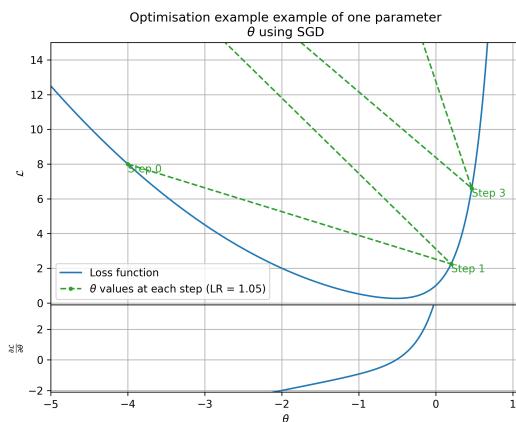


FIGURE 3.8 – Illustration of the gradient explosion. Here it can be solved with a lower learning rate but its not always the case.

1070 **3.2 Neural networks architectures**

1071 **3.2.1 Fully Connected Deep Neural Network (FCDNN)**

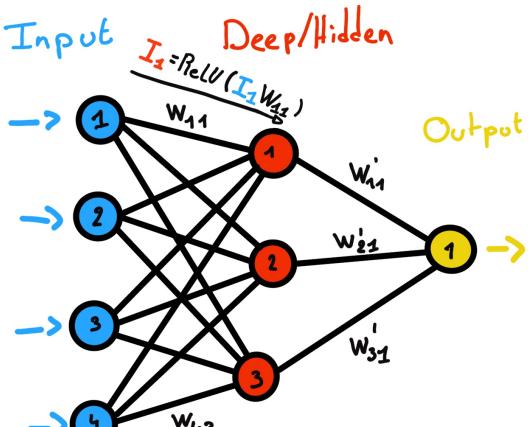
1072 The Fully Connected Deep Neural Network (FCDNN) architecture is the stack of multiple fully
 1073 connected layers as presented in the figure 3.9a. Most of the time, the classic ReLU function

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

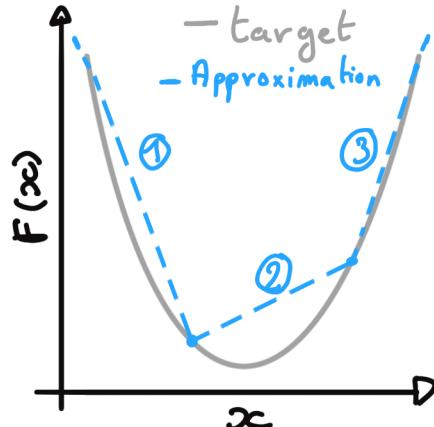
1074 is used as activation function. PReLU and Sigmoid are also popular choices:

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

1075 The reasoning behind ReLU and PReLU is that with enough of them, you can mimic any continuous
 1076 function as illustrated in figure 3.9b. Sigmoid is more used in case of classification, its behavior going
 1077 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

1079 Due to its simplicity, FCDNN are also used as basic pieces for more complex architectures such as
 1080 the CNN and GNN that will be presented in the next sections.

1081 **3.2.2 Convolutional Neural Network (CNN)**

1082 It's not trivial to describe in text the principles of Convolutional Neural Network (CNN) and how
 1083 they works. We try a general description below followed by a step by step description of a concrete
 1084 example.

1085 Convolutional Neural Networks are a family of neural networks that use discrete convolution filters,
 1086 as illustrated in an example in figure 3.10, to process the input data, often images. They are com-
 1087 monly used in image recognition [55] for classification or regression problematics. Concretely, you
 1088 multiply element-wise a portion of the input data, in the case of an image, a small part of the image,

1089 with a kernel of same dimension. In figure 3.10, we multiply the 3×3 pixels sub-image with the
 1090 3×3 kernel.

1091 Their filters scan the input data, highlighting patterns of interest, this scanning procedure making
 1092 them translation-invariant. In the concrete case of figure 3.10, for each pixel of the input image, we
 1093 group it with the 8 neighbours pixel and produce a new pixel that correspond to the output image.
 1094 For the pixel on the edges that do not have neighbours, we either create “imaginary” pixel with the
 1095 value 0 or we just ignore them. If we ignore them, the output image will posses fewer pixels than the
 1096 input image. We see that the operation do not care where is the pattern of interest in the images, the
 1097 filter output will be *invariant* whatever *translation* is applied to the image.

1098 This invariance mean that they are capable of detecting oriented features independently of their
 1099 location on the image. Again taking 3.10 as an example, with only the 9 parameters composing the
 1100 kernel, we can highlight the contour of the duck by looking at the “yellowness” of the pixels.

1101 The learning parameters of CNNs are the kernels components, the network thus learn the optimal
 1102 filters to extract the desired features.

1103 The convolution layers are commonly chained [56], reducing the input dimension while increasing
 1104 the number of filters. The idea behind is that the first layers will process local informations and
 1105 the latest layers will process more global informations, as the latest convolution filters will process
 1106 the results of the preceding that themself have processed local information. To try to preserve the
 1107 amount of information, we tend to grow the numbers of filters for each division of the input data.
 1108 The results of the convolution filters is commonly then flattened and feed to a smaller FCDNN which
 1109 will process the filters results to yield the desired output.

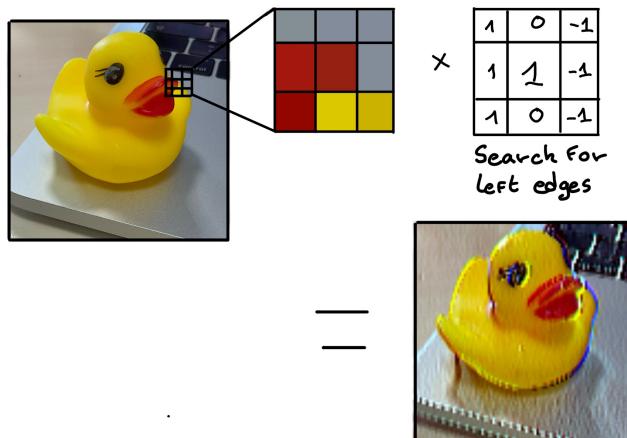


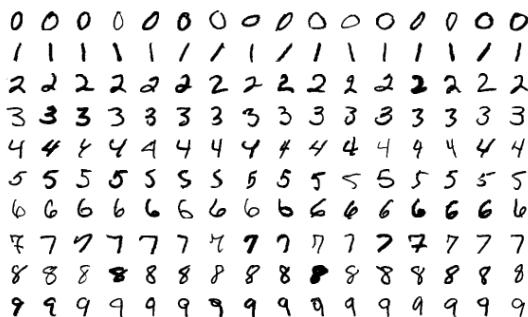
FIGURE 3.10 – Illustration of the effect of a convolution filter. Here we apply a filter with the aim do detect left edges. We see in the resulting image that the left edges of the duck are bright yellow where the right edges are dark blue indicating the contour of the object. The convolution was calculated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

1110 As an example, let’s take the Pytorch [58] example for the MNIST [59], a dataset of black and white
 1111 images of handwritten digits. Those images are 28×28 pixels with only one channel corresponding
 1112 to the grey level of the pixel. Example of images from this dataset are presented in figure 3.11a

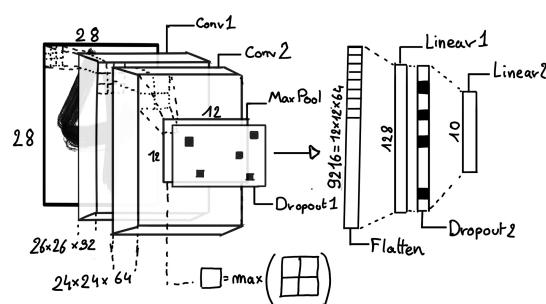
1113 A schema of the CNN used in the Pytorch example is presented in figure 3.11b. Using this schema
 1114 as a reference, the trained network is made of:

- 1115 1. A convolutional layer of (3×3) filters yielding 32 channels. A bias parameter is applied
 1116 to each channel for a total of $(32 \cdot (3 \times 3) + 32) = 320$ parameters. The resulting image is
 1117 $(26 \times 26 \times 32)$ (26 per 26 pixels with 32 channels). The ReLU activation function is applied to
 1118 each pixel.

- 1119 2. A second convolutional layer of (3×3) filters yielding 64 channels. This channel also posses
 1120 a bias parameter for a total of $(64 \cdot (3 \times 3) + 64) = 640$ parameters. Resulting image is $(24 \times$
 1121 $24 \times 64)$. This channel also apply a ReLU activation function.
- 1122 3. Then comes a (2×2) max pool layer with a stride of 1 meaning that for each channel the max
 1123 value of pixels in a (2×2) block is condensed in a single resulting pixel. The resulting image
 1124 is $(12 \times 12 \times 64)$.
- 1125 4. This image goes through a dropout layer which will set the pixel to 0 with a probability of
 1126 0.25. This help prevent overtraining the neural network (see section 3.1.4 for more details).
- 1127 5. The data is the flattened i.e. condensed into a vector of $(12 \times 12 \times 64) = 9216$ values.
- 1128 6. Then comes a fully connected linear layer (Eq. 3.2) with a ReLU activation that output 128
 1129 feature. It needs $(9216 \cdot 128) + 128 = 1'179'776$ parameters.
- 1130 7. This 128 item vector goes through another dropout layer with a probability of 0.5
- 1131 8. The vector is then transformed through a linear layer with ReLU activation. It output 10
 1132 values, one for each digit class $(0, 1, 2, \dots, 9)$. It need $(128 \cdot 10) + 128 = 1408$ parameters.
- 1133 9. Finally the 10 values are normalized using a log softmax function $\text{LogSoftmax}(x_i) = \log\left(\frac{\exp(x_i)}{\sum_j \exp(x_j)}\right)$.
- 1134 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

1135 The final network needs 1'182'144 parameters or, if we consider each parameters to be a double
 1136 precision floating point, 9.45 MB of data. To gives a order of magnitude, such neural network is
 1137 considered "simple", train in a matter of minutes on T4 GPU [60] (14 epochs) and reach an accuracy
 1138 in its prediction of 99%.

1139 3.2.3 Graph Neural Network (GNN)

1140 As seen in the previous section, the CNNs are powerful for image processing, and more generally
 1141 any data that can be expressed as a regular, discrete space and from which the information reside
 1142 in the dispersion in this space. For an image, the edges of an object and how they assemble. A red
 1143 square, straight edges with a sharp angle between them, is much less representative of a duck than
 1144 an yellow sphere, round edges without sharp angles.

1145 This "image" projection is not fitted for every problematics. The signals produced by a detector does
 1146 not always have the properties of images. In the case of JUNO for example, we can create an image
 1147 of two channels, one for the charge Q and one for the timing t but this image should be spheric.
 1148 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

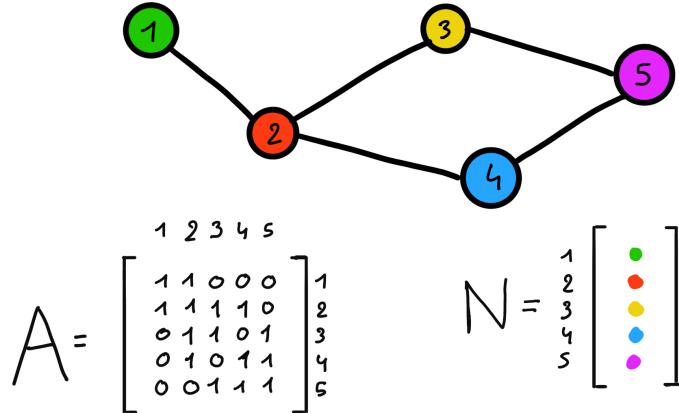


FIGURE 3.12 – Illustration of a graph and its tensor representation.

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_e^{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.

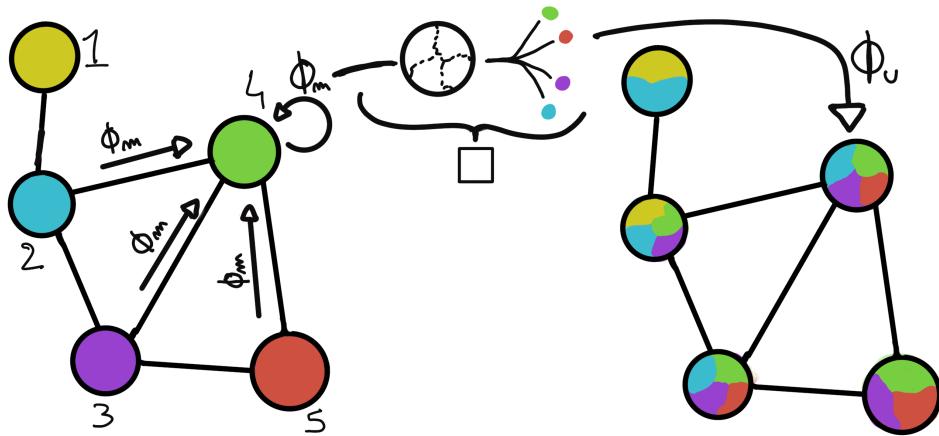


FIGURE 3.13 – Illustration of the message passing algorithm. The detailed explanation can be found in section 3.2.3

1165 To process such object we need specific machine learning algorithms we call Graph neural network.
 1166 To efficiently manipulate graph we need to structurally encode their property in the neural network
 1167 computing architecture: each node is equivalent (as opposite to ordered data in a vector), each node
 1168 has a set of neighbours, ... One of this method is the message passing algorithm presented historically
 1169 in "Neural Message Passing for Quantum Chemistry" [61]. In this algorithm, with each layer of
 1170 message passing a new set of features is computed for each node following

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

1171 where ϕ_u is a differentiable *update* function, \square_j is a differentiable *aggregation* function and ϕ_m is a
 1172 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
 1173 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
 1174 algorithm is applied. The update function need also a few other property if we want to keep the
 1175 graph property, most notably the permutational invariance of its parameters (example: mean, std,
 1176 sum, ...). The differents message, update and aggregation functions can really be any kind of function
 1177 if they follow the constraint presented before, even small Neural Network.

1178 The edges features can also be updated, either by directly taking the results of ϕ_m or by using another
 1179 message function ϕ_e .

1180 To explain this process, let's take the situation presented in figure 3.13. We start with an input graph
 1181 on left, in this case the message passing algorithm is mixing the color on each nodes and produce
 1182 nodes of mixed color. For simplicity, the ϕ_m and ϕ_u function are the identity, they take a color and
 1183 output the same color.

1184 Let's look at what's happening in the node 4. It has 3 neighbours and is a neighbour of itself. The four
 1185 resulting ϕ_m extract the color of each nodes and then feed them to the \square function. The \square function
 1186 just equally distribute the color in the node. Finally the ϕ_u function just update the node with the
 1187 output of \square .

1188 Interestingly we see that the new node 4 does not have any yellow, the color of node 1. But if we were
 1189 to run the message passing algorithm again, it would get some as node 2 is now partially yellow. If
 1190 color here represent information, we see that multiple step are needed so that each node is "aware"
 1191 of the informations the other nodes possess.

1192 Message passing is a very generic way of describing the process of GNN and it can be specialized
 1193 for convolutional filtering [49], diffusion [62] and many other specific operation. GNN are used in a
 1194 wide variety of application such as regression problematics, node classification, edge classification,
 1195 node and edge prediction, ...

1196 It is a very versatile but complex tool.

1197 3.2.4 Adversarial Neural Network (ANN)

1198 The adversarial machine learning, Adversarial Neural Networks (ANN) in the case of neural net-
 1199 work, is a family of unsupervised machine learning algorithms where the learning algorithm (gen-
 1200 erator) is competing against another algorithm (discriminator). Taking the example of Generative
 1201 Adversarial Networks, concept initially developed by Goodfellow et al. [63], the discriminator goal
 1202 is to discriminate between data coming from a reference dataset and data produced by the generator.
 1203 The generator goal, on the other hand, is to produce data that the discriminator would not be able to
 1204 differentiate from data from the reference dataset. The expression of duality between the two models
 1205 is represented in the loss where, at least a part of it, is driven by the results of the discriminator.

¹²⁰⁶ **Chapter 4**

¹²⁰⁷ **Image recognition for IBD
reconstruction with the SPMT system**

¹²⁰⁹

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

¹²¹⁰

Contents

¹²¹¹ 4.1 Method and model	¹²¹² 56
¹²¹³ 4.1.1 Model	¹²¹⁴ 57
¹²¹⁴ 4.1.2 Data representation	¹²¹⁵ 58
¹²¹⁵ 4.1.3 Dataset	¹²¹⁶ 60
¹²¹⁶ 4.1.4 Data characteristics	¹²¹⁷ 61
¹²¹⁷ 4.2 Training	¹²¹⁸ 63
¹²¹⁸ 4.3 Results	¹²¹⁹ 63
¹²¹⁹ 4.3.1 J21 results	¹²²⁰ 64
¹²²⁰ 4.3.2 J21 Combination of classic and ML estimator	¹²²¹ 66
¹²²¹ 4.3.3 J23 results	¹²²² 68
¹²²² 4.4 Conclusion and prospect	¹²²³ 70

¹²²⁴

¹²²⁵

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

¹²³¹

¹²³²

¹²³³

¹²³⁴

¹²³⁵

¹²³⁶

¹²³⁷

¹²³⁸

¹²³⁹

The SPMT system is essential to the deployment of the Dual Calorimetry techniques, already mentioned in Section 2.6 and described in [27, 29, 64]. It is indeed less subject than the LPMTs to charge non linearity effects (QNL). This topic will be studied in more detail in Chapter 7, where the potential of one of the Dual Calorimetry techniques is explored. It consists on combined oscillation analyses based on two antineutrino energy spectra : one reconstructed with the LPMT system, the other one with the SPMT system. For that purpose, it is therefore necessary to have reconstruction tools available. Well maintained tools using the LPMT are available in the collaboration's official software. This is not the case concerning the SPMT system, where algorithms were developed more sporadically. This is one of the reasons why we developed the CNN described in this chapter.

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

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

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

4.1 Method and model

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

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

Some CNN are developed to process spherical images [69] but for the sake of simplicity and as a first approach we decided to go with a planar projection of the detector, approach that has proven its efficiency using the LPMT system (see section 2.6.4). The details about this planar projection will be discussed in section 4.1.2.

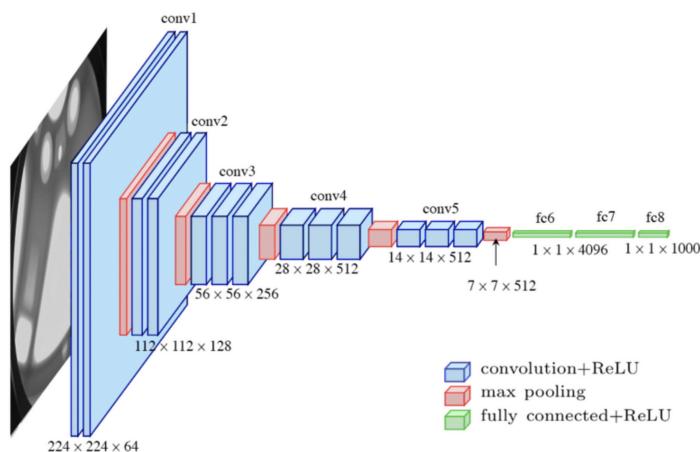


FIGURE 4.1 – Graphic representation of the VGG-16 architecture, presenting the different kind of layer composing the architecture.

1268 4.1.1 Model

1269 The architecture we use is derived from the VGG-16 architecture [56] illustrated in figure 4.1. We
 1270 define a set of hyperparameters that will define the size, complexity and computational power of the
 1271 NN. The chose hyperparameters are detailed below and their values are presented in table 4.1.

- 1272 — **N_{blocks}** : the number of convolution blocks, a block being composed of two convolutional
 1273 layers with 3×3 filters using ReLU activation function, a 3×3 kernel max-pooling layer
 1274 (except for the last block).
- 1275 — **$N_{channels}$** : The number of channels in the first block. The number of channels in the subsequent
 1276 blocks is computed using $N_{channels}^i = i * N_{channels}$, $i \in [1..N_{blocks}]$.
- 1277 — **FCDNN configuration**: The result of the last convolution layer is flattened then fed to a
 1278 FCDNN. Its configuration is expressed as the ouputs of sequenced fully connected linear layer
 1279 using the PReLU activation function. For example $2 * 1024 + 2 * 512$ is the sequence of 2 layers
 1280 which output is 1024 followed by 2 other layers with an output of 512. Finally the last layer
 1281 is a linear layer outputting 4 features without activation function. Each feature of the last layer
 1282 represent a component of the interaction vertex: Energy, X, Y, Z.
- 1283 — **Loss**: The loss function. In this work we study two different loss function $(E + V)$ and $(E_r + V_r)$ detailed below.

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

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

1285 where E_{dep} is the deposited energy and R is the radius of JUNO's CD. With the energy in MeV and
 1286 the distance in meters, we use the factor 0.85 and 10 to balance the two term of the loss function so
 1287 they have the same magnitude.

1288 The loss function $(E + V)$ is close to a simple Mean Squared Error (MSE). MSE is one of the most
 1289 basic loss function, the derivative is simple and continuous in every point. It is a strong starting
 1290 point to explore the possibility of CNNs. The loss $(E_r + V_r)$ can be seen as a relative MSE.

1291 The idea is that: due to the inherent statistic uncertainty over the number of collected Number of
 1292 Photo Electrons (NPE), the absolute resolution $\sigma(E - E_{true})$ will be larger at higher energy than at
 1293 low energy. But we expect the *relative* energy resolution $\frac{\sigma(E - E_{true})}{E_{true}}$ to be smaller at high energy than
 1294 lower energy as illustrated in figure 2.22. Because of this, by using simple MSE the most important
 1295 part in the loss come from the high energy part of the dataset whereas with a relative MSE, the
 1296 most important part become the low energy events in the dataset. We hope that by using a relative
 1297 MSE, the neural network will focus on low energy events where the reconstruction is considered the
 1298 hardest.

1299 The above losses and their parameters values results from fine-tuning after multiples runs and
 1300 adjustments of the full random search.

1301 Each combinations of those hyperparameters (for example ($N_{blocks} = 2, N_{channels} = 32$, FCDNN =
 1302 $(2 * 1024)$, Loss = $(E + V)$)) produce models, hereinafter referred as configurations, are then tested
 1303 and compared to each other over an analysis sample.

1304 On top those generated models, we define 4 hand tailored models:

- 1305 — Gen₀: $N_{blocks} = 4, N_{channels} = 64$, FCDNN configuration: $1024 * 2 + 512 * 2$, Loss $\equiv E + V$
- 1306 — Gen₁: $N_{blocks} = 4, N_{channels} = 64$, FCDNN configuration: $1024 * 2 + 512 * 2$, Loss $\equiv E_r + V_r$
- 1307 — Gen₂: $N_{blocks} = 5, N_{channels} = 64$, FCDNN configuration: $4096 * 2 + 1024 * 2$, Loss $\equiv E + V$
- 1308 — Gen₃: $N_{blocks} = 5, N_{channels} = 64$, FCDNN configuration: $4096 * 2 + 1024 * 2$, Loss $\equiv E_r + V_r$

1310 The resulting models possess between 2'041'034, for Gen₅₂ and Gen₅₃, and 5'759'839'242 parameters,
 1311 for Gen₂₆ and Gen₂₇. The models of interest in this thesis, from which the results are discussed
 1312 in section 4.3, possess 86'197'196 parameters for Gen₃₀ and 332'187'530 parameters for Gen₄₂. For
 1313 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

1314 To rank the various configuration we cannot used directly the mean loss over the validation dataset
 1315 as ($E + V$) and ($E_r + V_r$) are not numerically comparable. We thus use the following quantities,
 1316 directly related to the reconstruction performances:

- 1317 — The mean absolute energy error $\langle E \rangle = \langle |E - E_{true}| \rangle$. It is an indicator of the energy bias of our
 1318 reconstruction.
- 1319 — The standard deviation of the energy error $\sigma E = \sigma(E - E_{true})$. This the indicator on our
 1320 precision in energy reconstruction.
- 1321 — The mean distance between the reconstructed vertex and the true vertex $\langle V \rangle = \langle |\vec{V} - \vec{V}_{true}| \rangle$.
 1322 This an indicator of the bias and precision of our vertex reconstruction.
- 1323 — The standard deviation of the distance between the true and reconstructed vertex $\sigma V = \sigma|\vec{V} -$
 1324 $\vec{V}_{true}|$. This is an indicator if the precision in our vertex reconstruction.

1326 The models were developped in Python using the Pytorch framework [58] using NVIDIA A100 [70]
 1327 and NVIDIA V100 [71] gpus. The A100 was split in two, thus the accessible gpu memory was
 1328 the same as V100, 20 Gb, making it impossible to train some of the architectures due to memory
 1329 consumption.

1330 The training was monitored in realtime by a custom tooling that was developed during this thesis,
 1331 DataMo [72].

1332 The training of one model takes between 4h and 15h depending of its size, overall training the full
 1333 72 models takes around 500 GPU hours. Even with parallel training, this random search hyper-
 1334 optimisation was time consuming.

1335 4.1.2 Data representation

1336 This data is represented as 240×240 images with a charge Q channel and a time t channel. The
 1337 SPMTs are then projected on the plane as illustrated in figure 4.2b using the coordinate system
 1338 presented in 4.2a. The P_y coordinate, the row corresponding to the SPMT in the projection, is
 1339 proportional to θ . The P_x coordinate, the column corresponding to the SPMT in the projection, is
 1340 defined by $\phi \sin \theta$ in spherical coordinates. $\theta = 0$ is defined as being the top of the detector and $\phi = 0$
 1341 is defined as an arbitrary direction in the detector. In practice, $\phi = 0$ is given by the MC simulation.

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

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

1342 where H is the height of the image, W the width of the image and $(0, 0)$ the top left corner of the
 1343 image.

1344 This projection keep the SPMT position in the image proportional to their spherical coordinates while
 1345 keeping the neighbouring information. This proportionality allow us to keep the specificities of the
 1346 detector structure, the vertical bands visible in 4.2b.

1347 When two SPMTs in the same pixel are hit in the event time window, the charges are summed and
 1348 the lowest of the hit-time is chosen. The time window depends on the datasets and are detailed in
 1349 section 4.1.2. The SPMTs being located close to each other, we expect the time difference between
 1350 two successive physics signals, two photons being collected, to be small. The first hit time is chosen
 1351 because it can be considered as the relative propagation time of the photons that went the "straightest", i.e.
 1352 that went under the less perturbation of the two. The timing is thus more representative of
 1353 the event location.

1354 The only potential problem in using this first time come from the Dark Noise (DN). Its time distribution
 1355 is uniform over the signal and could come before a physics signal on the other SPMT in the pixel.
 1356 In that case, the time information in the pixel become irrelevant and we lose the timing information
 1357 for this part of the detector. As illustrated in figure 4.2b the image dimension have been optimized
 1358 so that at most two SPMTs are in the same pixel while keeping the number of empty pixels relatively
 1359 low to prevent this kind of issue.

1360 While it could be possible to use larger images (more pixel) to prevent overlapping, keeping image
 1361 small images gives multiple advantages:

- 1362 — As presented in section 4.1.1, the convolution filter we use are 3×3 convolution filter, meaning
 1363 that if SPMTs would be separated by more than one pixel, the first filter would only see one
 1364 SPMT per filter. This behavior would be kind of counterproductive as the first convolution
 1365 block would basically be a transmission layer and would just induce noise in the data.
- 1366 — It keep the network relatively small, while this do not impact the convolution layers, the
 1367 flatten operation just before the FCDNN make the number parameters in the first layer of
 1368 it dependent on the size of the image.
- 1369 — It reduce the number of empty pixel in the image.

1370
 1371 The question of empty pixel is an important question in this data representation. There is two kind
 1372 of empty pixels in the data.

1373 The first kind is pixel that contain a SPMT but the SPMT did not get hit nor registered any dark noise
 1374 during the event. In this case, the charge channel is zero, which have a physical meaning but then
 1375 come the question of the time layer. One could argue that the correct time would be infinity (or the
 1376 largest number our memory allows us) because the hit "never" happened, so extremely far from the
 1377 time of the event. This cause numerical problem as large number, in the linear operation that are
 1378 happening in the convolution layers, are more significant than smaller value. We could try to encode
 1379 this feature in another way but no number have any significance due to our time being relative to
 1380 the trigger of the experiment so -1 for example is out of question. Float and Double gives us access
 1381 to special value such as NaN (Not a Number) [73] but the behavior is to propagate the NaN which
 1382 leaves us with NaN for energy and position. We choose to keep the value 0 because it's the absorbing
 1383 element of multiplication, absorbing the "information" of the parameter it would be multiplied by.
 1384 It also can be though as no activation in the ReLU activation function. It's important to keep in mind

1385 the fact that a part of the detector that has not been hit is also an information: There is no signal in
 1386 this part of the detector. This problematic will be explored in more details in Chapter 5.

1387 The second kind of pixels are the one that do not represent parts of the detector such as the corners
 1388 of the image. The question is basically the same, what to put in the charge and the time channel. The
 1389 decision is to set the charge and time to 0 following the above reasoning.

1390 Another problematic that happens with this representation, and this is not dependent of the chosen
 1391 projection, is the deformation in the edges of the image and the loss of the neighbouring information
 1392 in the for the SPMTs at the edge of the image $\phi \sim 180^\circ$. This deformation and neighbouring loss
 1393 could be partially circumvented as explained in section 4.4

1394 4.1.3 Dataset

1395 In this study we will discuss two datasets of one millions prompt signal of IBD events.

1396 J21

1397 The first one comes from the JUNO official MC simulation J21v1r0-Pre2 (released the 18th August
 1398 2021). This historical version is the one on which the classical SPMT reconstruction algorithm,
 1399 presented in section 2.6.3 [65], was developed. This dataset is used as a reference for comparison
 1400 to classical algorithm performances. The data in this dataset is *detsim* level (see section 2.5) which
 1401 includes no digitization, no DAQ and therefore no reconstruction of PMT signals. Only the number
 1402 of PEs that hit a PMT and the hit times are provided. A fast simulation based on gaussian drawings
 1403 produces charges, with bias and variability, and the equivalent for times. The drawings parameters
 1404 were adjusted based on [26, 74]. Because there is no charge reconstruction, the timing on the event
 1405 is based on the Geant4 simulation, and so $t = 0$ is the moment the positron is created in the CD. To
 1406 prevent correlation between the numerical value of the time of the first hit t_0 and the radius of the
 1407 event, we offset all time by this first hit time. Without simulation of the charge reconstruction, we
 1408 cannot simulate the event trigger, we thus add an arbitrary time cut at a $t_0 + 1000$ ns.

1409 J23

1410 The second comes from the JUNO official monte-carlo simulations J23.0.1-rc8.dc1 (released the 7th
 1411 January 2024). The data is *calib* level (see section 2.5). Here the charge comes from the waveform
 1412 integration, the time window resolution and trigger decision are all simulated inside the software.

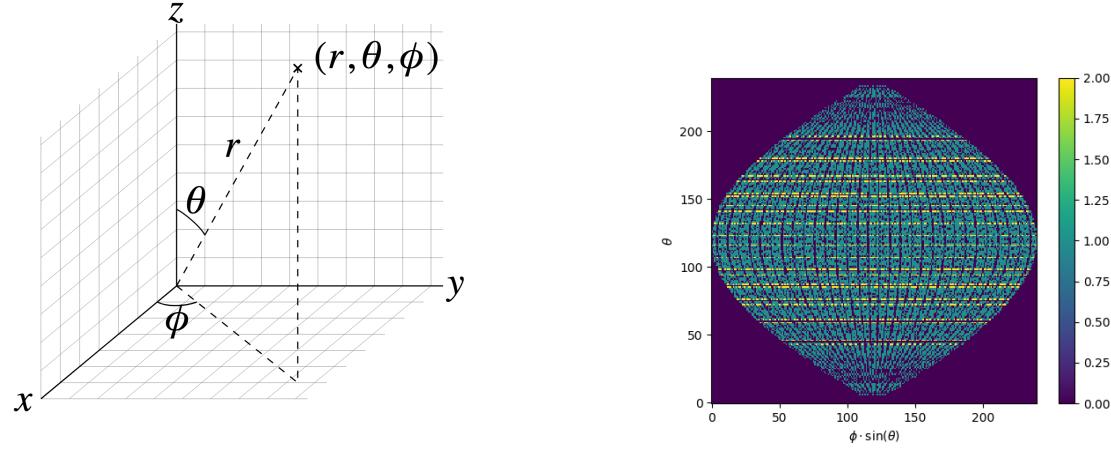
1413 To put in perspective this amount of data, the expected IBD rate in JUNO is 47 / days. Taking into
 1414 account the calibration time, and the source reactor shutdown, it amount to $\sim 94'000$ IBD events
 1415 in 6 years. With this million of event, we are training the equivalent of ~ 10 years of data. With
 1416 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
 1417 event scenarios that could be happening in the detector.

1418 While we expect and hope the MC simulation to give use a realistic representation of the detector,
 1419 there could be effect, even after the fine-tuning on calibration data, that the simulation cannot handle.
 1420 Thus, once the calibration will be available, we will need to evaluate, and if needed retrain, the
 1421 network on calibration data to establish definitive performances.

1422 The simulated data is composed of positron events, uniformly distributed in the CD volume and in
 1423 kinetic energy over $E_k \in [0; 9]$ MeV producing a deposited energy $E_{dep} \in [1.022; 10.022]$ MeV. This is
 1424 done to mimic the signal produced by the IBD prompt signal. Uniform distributions are used so that
 1425 the CNN does not learn a potential energy distribution, favoring some part of the energy spectrum
 1426 instead of other.

4.1.4 Data characteristics

To delve a bit into the kind of data we will use, you can find in figure 4.2b the repartition of the 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

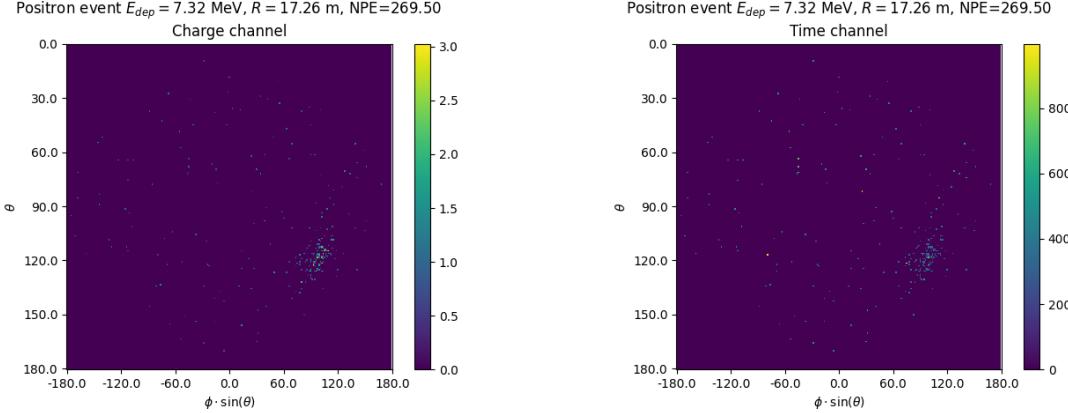


FIGURE 4.3 – Example of a high energy, radial event. We see a concentration of the charge on the bottom right of the image, clear indication of a high radius event. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

See also figures 4.3 to 4.6 - and the explanation in their captions - which present events from J23 for different positions and energies. We see some characteristics and we can instinctively understand how the CNN could discriminate different situations.

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

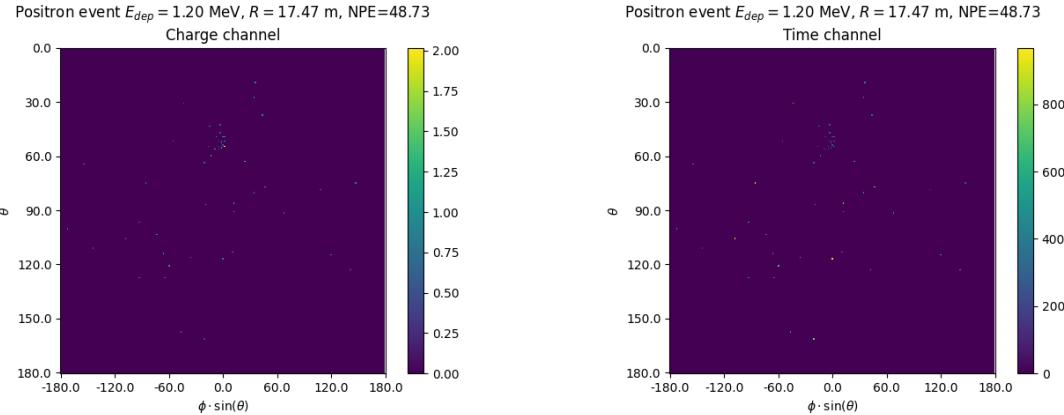


FIGURE 4.4 – Example of a low energy, radial event. The signal here is way less explicit, we can kind of guess that the event is located in the top middle of the image. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in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

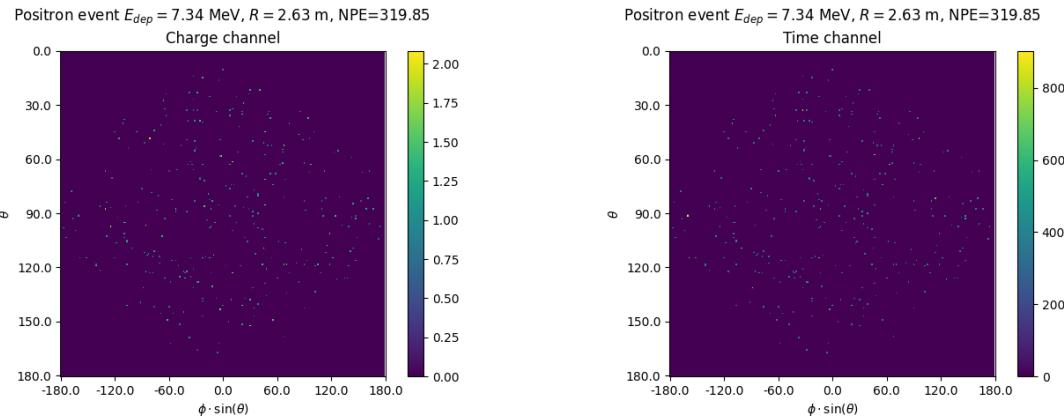


FIGURE 4.5 – Example of a high energy, central event. In this image we can see a lot of signal but uniformly spread, this is indicative of a central event. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

LS we can model:

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

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

¹⁴³³ where NPE_{tot} is the total number of PE detected by the event, P_{mev} is the mean number of PE detected
¹⁴³⁴ per MeV and D_N is the dark noise contribution that is considered energy independent. In the case
¹⁴³⁵ where the readout time window is dependent of the energy the dark noise contribution become
¹⁴³⁶ energy dependant, also the LS response is realistically energy dependant but figure 4.7a shows that
¹⁴³⁷ we are heavily dominated by the stochastic behavior of light emission and detection.

¹⁴³⁸ The fit shows a light yield of 40.78 PE/MeV and a dark noise contribution of 4.29 NPE. As shown in
¹⁴³⁹ figure 4.7b, the physics makes for 90% of the signal at low energy.

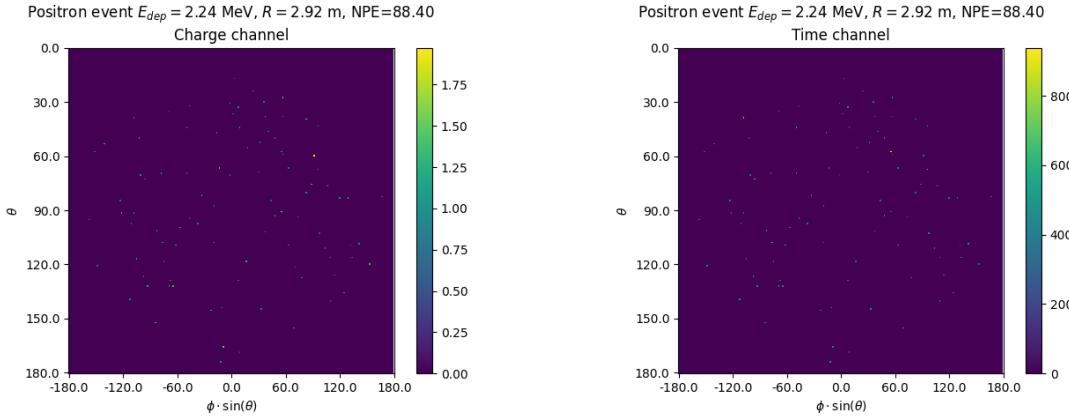


FIGURE 4.6 – Example of a low energy, central event. Here there is no clear signal, the uniformity of the distribution should make it central. **On the left:** the charge channel. The color is the charge in each pixel in NPE equivalent. **On the right:** The time channel in nanoseconds.

4.2 Training

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

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

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

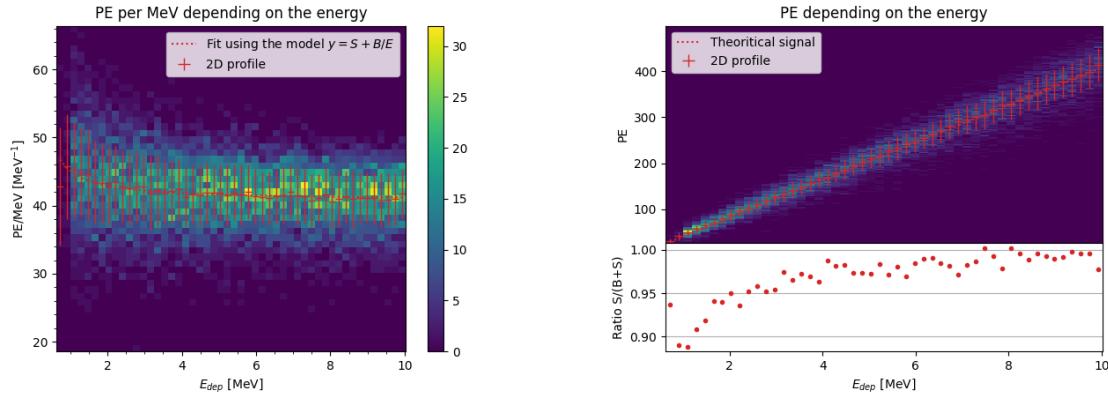
4.3 Results

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

The events are considered point-like in this study. The target truth position, or vertex, is the mean position of the energy deposits of the positron and the two annihilation gammas. This approximation for point-like interaction is also used for the likelihood study presented in section 2.6 and in previous ML studies presented in section 2.6.4 [42].

Due to the symmetries of the detector, we mainly consider and discuss the bias and precision evolution depending on the radius R but we will still monitor the performances depending on the spherical angle θ and ϕ . From the detector construction and effect we expect dependency in radius due to the TR area effect presented in section 2.6 and the possibility for the positron or the gammas to escape from the CD for positrons interacting near the edge. We also expect dependency on θ , the top of the experiment being non-instrumented due to the filling chimney. It is also to be noted that the events in the dataset are uniformly distributed in the CD, and so are uniformly distributed in R^3 and ϕ . The θ distribution is not uniform and we will have more events for $\theta \sim 90^\circ$ than $\theta \sim 0^\circ$ or $\theta \sim 180^\circ$.

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 2.6 reconstruct E_{vis} while the ones presented in section 2.6.4 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).

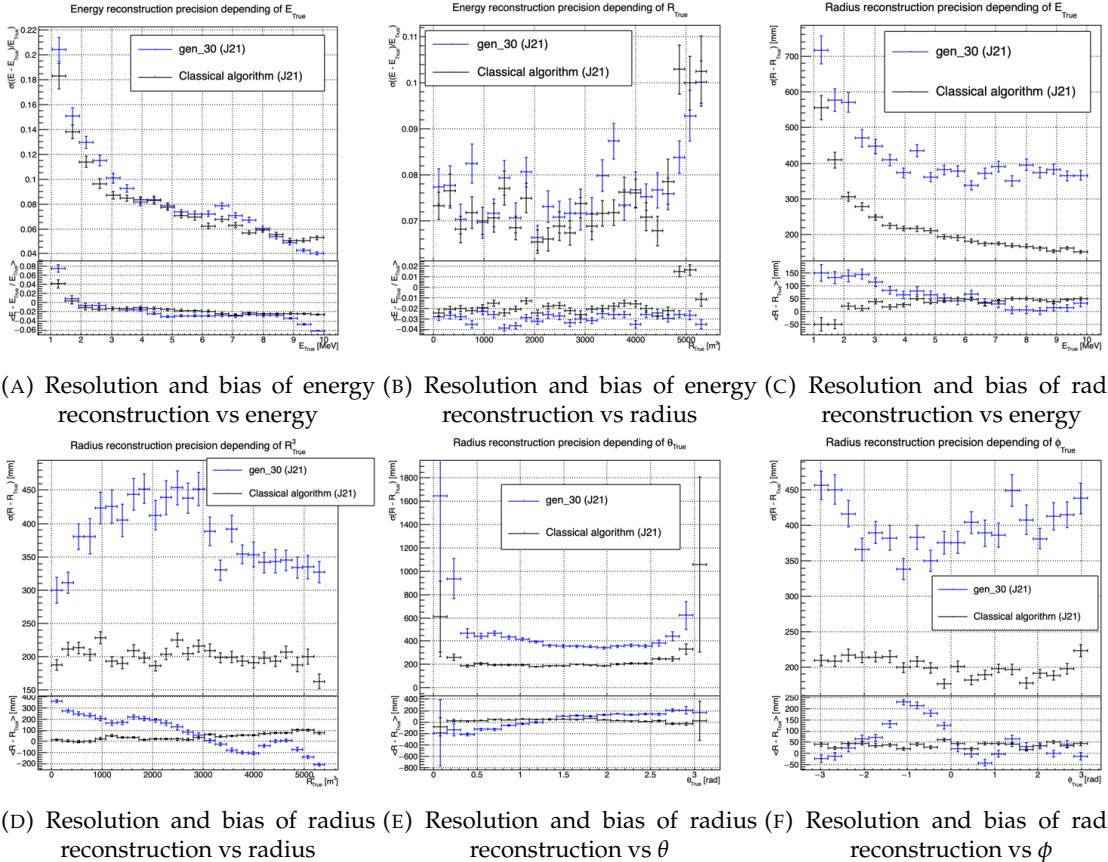


FIGURE 4.8 – Reconstruction performance of the Gen₃₀ model on J21 data and its comparison to the performances of the classic algorithm “Classical algorithm” from [65]. The top part of each plot is the resolution and the bottom part is the bias.

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.

1502 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.

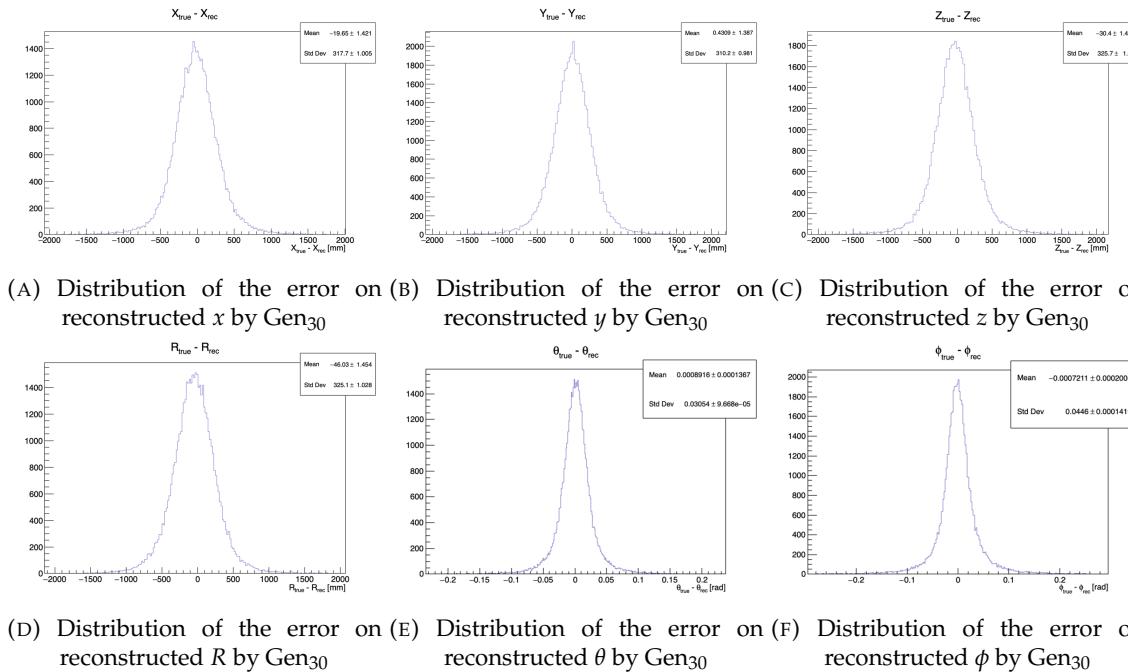
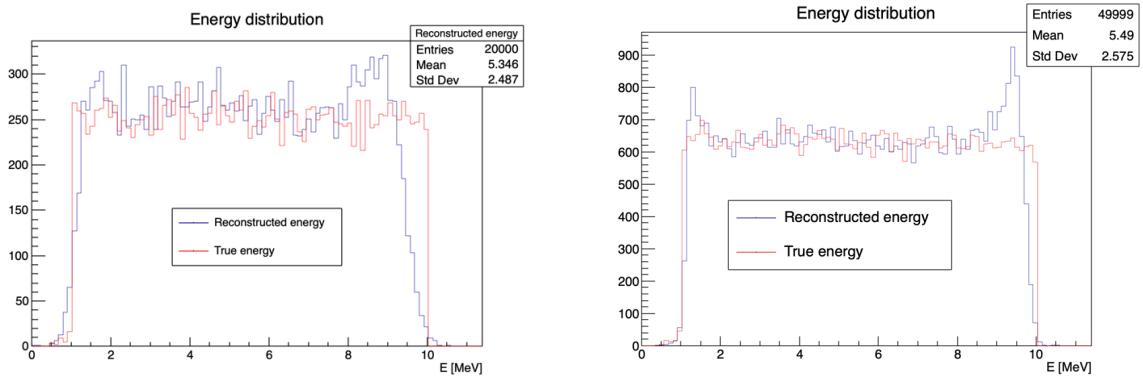


FIGURE 4.9 – Residual distribution of the different component of the vertex by Gen₃₀. The reconstructed component are x , y and z but we see similar behavior in the error of R , θ and ϕ .

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

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

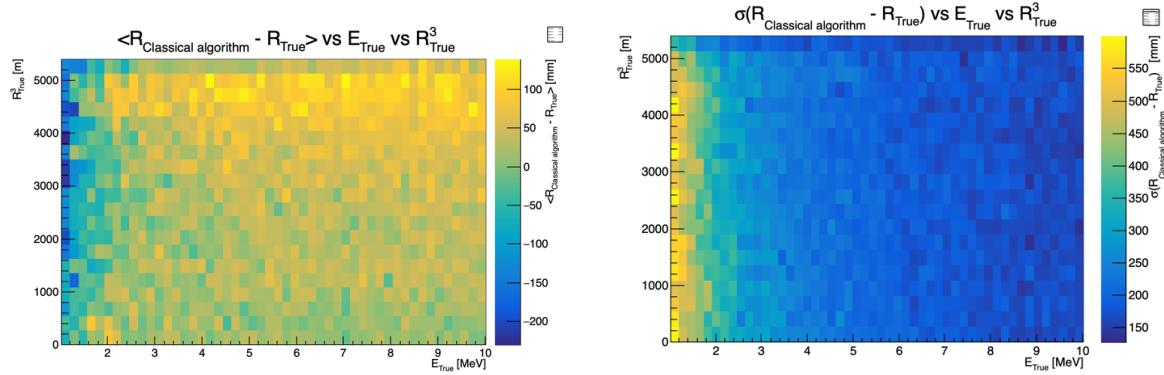


FIGURE 4.11 – Radius bias (on the left) and resolution (on the right) of the classical algorithm in a E, R^3 grid

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 m³ (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

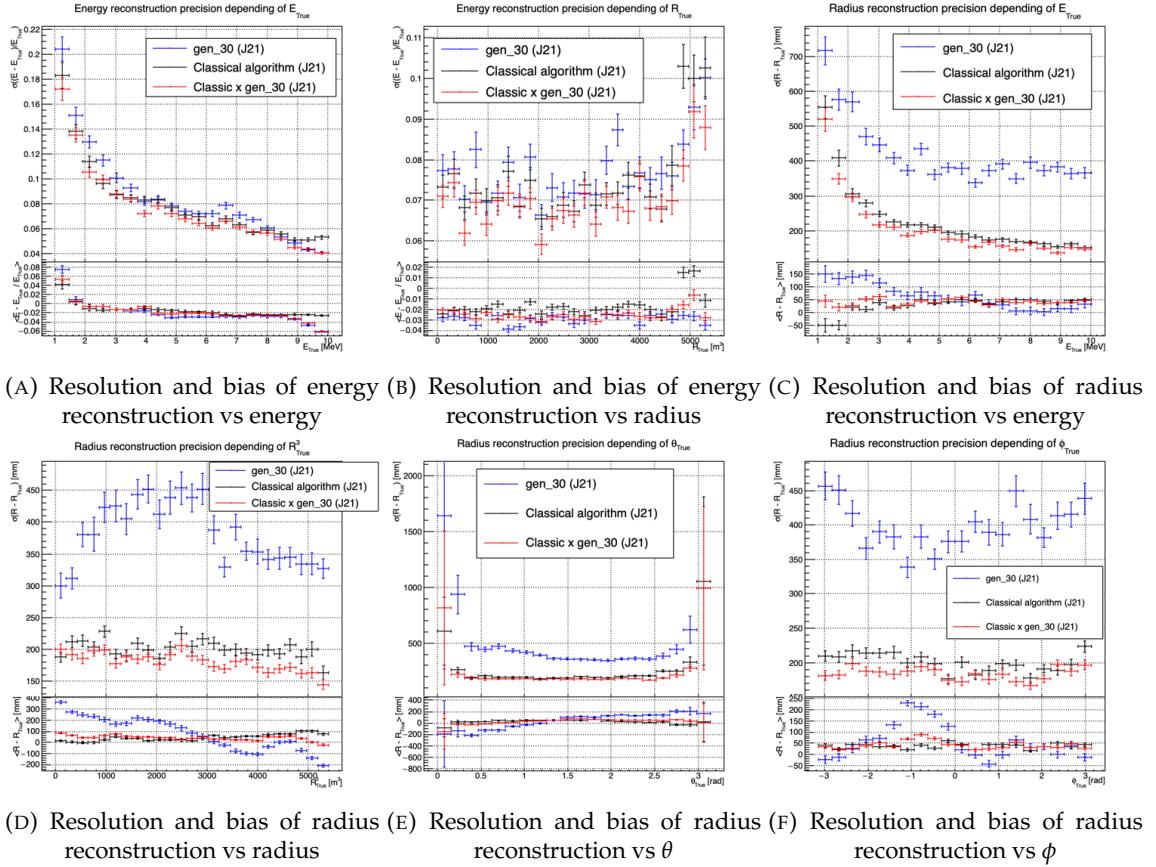


FIGURE 4.12 – Reconstruction performance of the Gen30 model on J21, the classic algorithm “Classical algorithm” from [65] and the combination of both using weighted mean. The top part of each plot is the resolution and the bottom part is the bias.

1559 light yield, reflection indices of materials of the detector, structural elements such as the connecting
1560 structure and more realistic dark noise. Additionally, the trigger, waveform integration and time
1561 window are defined using the algorithms that will ultimately be used by the collaboration to process
1562 real physics events.

1563 We retrained the models defined in 4.1.1 on the J23 data and used the same hyperparameter optimisation
1564 procedure. The results from the best architecture, Gen₄₂, are presented in figure 4.14. Following
1565 the table 4.1, Gen₄₂: $N_{blocks} = 3$, $N_{channels} = 64$, FCDNN configuration: $4096 * 2$, Loss $\equiv E + V$.

1566 Energy reconstruction

1567 The results of the energy reconstruction are presented in figures 4.14a and 4.14b. The resolution is
1568 close to the one of the classical algorithm with the exception of the start and end of the spectrum.
1569 This is the same effect that we saw with Gen₃₀, events are pulled from the edge of the distribution,
1570 resulting in smaller resolution but heavy biases.

1571 Vertex reconstruction

1572 The vertex reconstruction, presented in figures 4.14c, 4.14d, 4.14e and 4.14f is not yet to the level of
1573 the classical reconstruction but the degradation is smaller than for Gen₃₀ being at most a difference

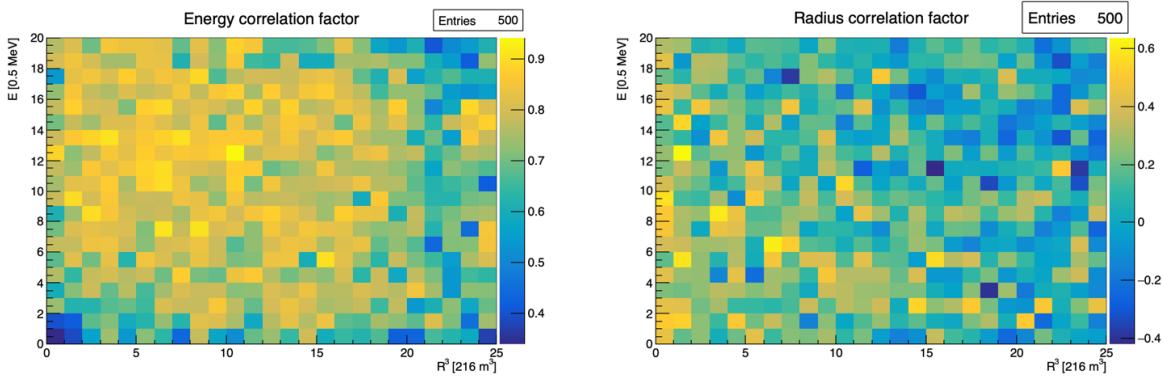


FIGURE 4.13 – Correlation between CNN and classical method reconstruction (on the left) for energy and (on the right) for radius in a E, R^3 grid

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.

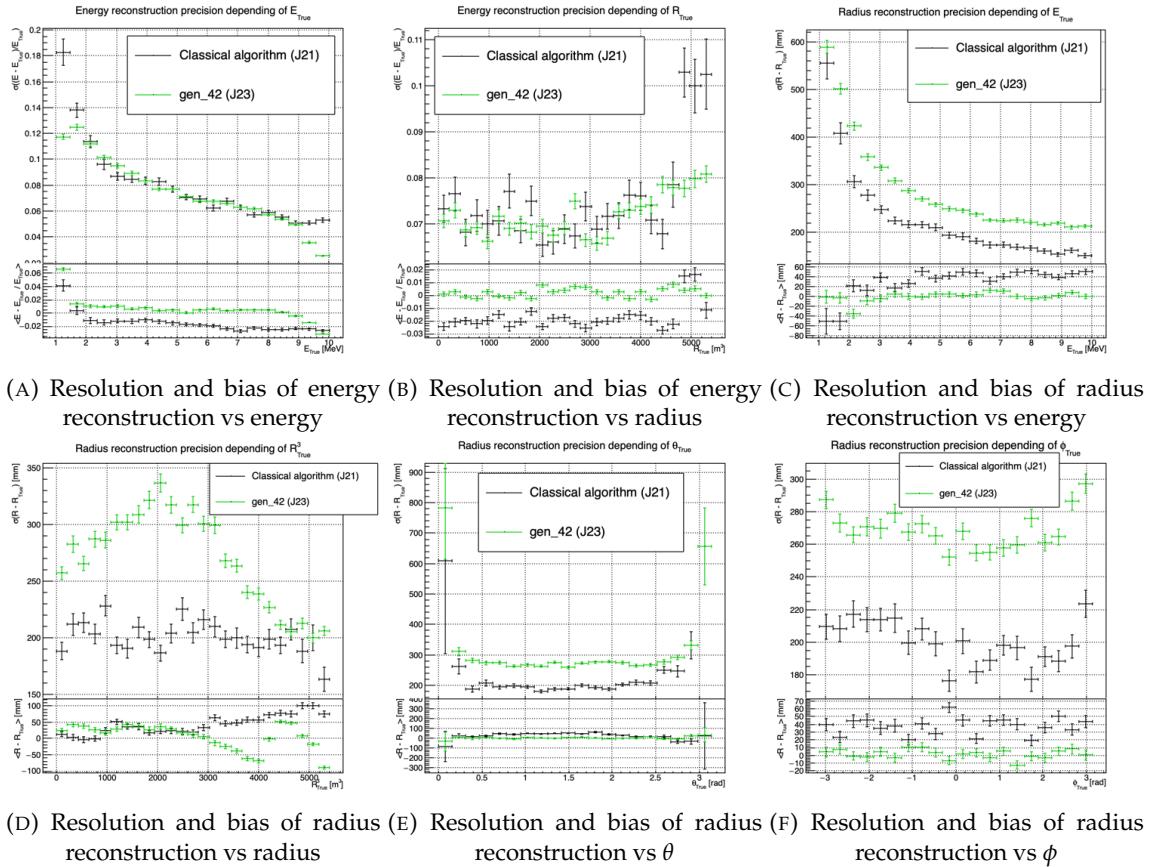


FIGURE 4.14 – Reconstruction performance of the Gen42 model on J23 data and its comparison to the performances of the classic algorithm “Classical algorithm” from [65]. The top part of each plot is the resolution and the bottom part is the bias.

¹⁶⁰³ **Chapter 5**

¹⁶⁰⁴ **Graph representation of JUNO for
IBD reconstruction**

¹⁶⁰⁶

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

Douglas Adams, The Hitchhiker's Guide to the Galaxy

¹⁶⁰⁷

Contents

¹⁶⁰⁸

¹⁶⁰⁹

¹⁶¹⁰

¹⁶¹¹

¹⁶¹²

¹⁶¹³

¹⁶¹⁴

¹⁶¹⁵

¹⁶¹⁶

¹⁶¹⁷

¹⁶¹⁸

¹⁶¹⁹

¹⁶²⁰

¹⁶²¹

5.1	Data representation	74
5.2	Message passing algorithm	76
5.3	Data	78
5.4	Model	80
5.5	Training	80
5.6	Optimization	82
5.6.1	Software optimization	82
5.6.2	Hyperparameters optimization	83
5.7	performance of the final version	83
5.8	Conclusion	87

¹⁶²²

¹⁶²³

¹⁶²⁴

¹⁶²⁵

¹⁶²⁶

¹⁶²⁷

In section [2.6.4](#), 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.

¹⁶²⁸

¹⁶²⁹

¹⁶³⁰

¹⁶³¹

¹⁶³²

¹⁶³³

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.

¹⁶³⁴

¹⁶³⁵

¹⁶³⁶

¹⁶³⁷

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.

¹⁶³⁸

¹⁶³⁹

¹⁶⁴⁰

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,

1641 i.e. the energy and radius of the event do change by rotation our referential. For more details see
1642 section 3.2.3. Graph representation also has the advantage to be able to encode global and higher
1643 order informations.

1644 5.1 Data representation

1645 In section 2.6.4, we mentioned a GNN developed before the beginning of this thesis to reconstruct
1646 IBD energies in JUNO [42]. In their approach: nodes of the graph correspond to 3072 pixels representing
1647 geometric regions of the detector and the information of the ~ 6 LPMTs found in a pixel are then
1648 aggregated on those nodes. The network then process the data using the equivalent of convolution
1649 but on graph [49]. In the first layer, each node is connected only with its direct neighbours.

1650 To determine the energy released by an IBD in the LS, it is helpful to determine the position of
1651 the main energy deposit. Therefore, relative Q and t's of PMTs all around the sphere is a useful
1652 information. If in the first layer only neighbour nodes are linked, several layers are necessary to
1653 access this detector-wide information. In an ideal world, we would develop a Graph NN where each
1654 PMT is a node (even if it has not been hit in the event under consideration, since this is in itself an
1655 information) and where each node is connected to all the other ones. This makes the detector-wide
1656 information available as early as the first layer. This architecture might help the network to better
1657 learn. Such an architecture can also be motivated this way: one of the strength of GNN's is their
1658 capacity to encompass the characteristics of a detector. A node can be the representation of a detector
1659 element, and the edge can represent its relationship with other elements. In the case of JUNO, any
1660 measurement is collective : an interaction is seen by all the PMTs, with no a priori hierarchy in the
1661 role of each. A fully connected GNN, in that respect, seems to make sense.

1662 Another advantage of a GNN is also that it is well adapted to inhomogenous detectors. We therefore
1663 tried to build GNNs including both LPMTs and SPMTs.

1664 With 17612 LPMTs and 25600 SPMTs, the ideal fully connected Graph mentioned above is impossible:
1665 even excluding self relation and considering the relation to be undirected (the edge from a node A
1666 to a node B being the same from as the one from B to A) the amount of necessary edges would be
1667 $n(n - 1)/2$ with $n = 43212$ nodes. This amounts to 933'616'866 edges. If we encode an information
1668 with double precision (64 bits) in what we call an adjacency matrix, illustrated in figure 3.12, each
1669 information we want to encode in the relation would consume 4 GB of data. When adding the
1670 overhead due to gradient computation during training, this would put us over the memory capacity
1671 of a single V100 gpu card (20 GB of memory). We could use parallel training to distribute the training
1672 over multiple GPU but we considered that the technical challenge to deploy this solution was too
1673 high.

1674 We finally decided of a middle ground where we define three *families* of nodes:

- 1675 — The core of the graph is composed of nodes representing geometric regions of the detector.
1676 We call those nodes *mesh* nodes. Those mesh nodes are all connected to each other. We keep
1677 their number low to gain in memory consumption.
- 1678 — PMTs in which Photo-Electrons (PE) are found are represented by *fired* nodes. Fired nodes
1679 are connected to the mesh node they geometrically belong to.
- 1680 — A final node is called the input/output node (*I/O*). It is connected to every mesh node. Its
1681 features are combinations of signals found in the whole detector.

1682

1683 Those nodes and their relations are illustrated in figure 5.1a. From this representation, we end up
1684 with three distinct adjacency matrix

- 1685 — A $N_{\text{fired}} \times N_{\text{mesh}}$ adjacency matrix, representing the relations between fired and mesh. Those
1686 relations are undirected.

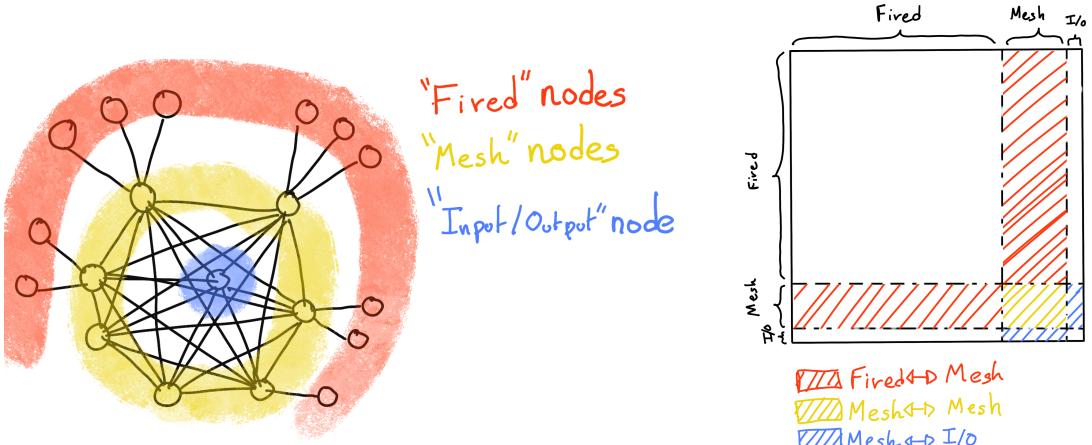


FIGURE 5.1

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

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

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

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

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

1702 harmonics in JUNO, see annex [B](#).

- 1703 2. \mathbb{A} and \mathbb{B} are informations that are related the likeliness of the interaction vertex to be on the
1704 segment between the center of two meshes.

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

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

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

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

1705 where \vec{i} is the position vector of the mesh i , D_{ij} is the distance between the center of the meshes
1706 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
1707 i and n the optical index of the LS. \mathbb{A} is the vertex between center of meshes distance ratio
1708 between i and j based on the time information. For \mathbb{B} , the charge ratio evolve with the square
1709 of the distance, so the mesh couple with the smallest \mathbb{B} should be the one with the interaction
1710 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]
1711 and position in [m].

1712 Q and t are the reconstructed charge and time of the hit PMTs. (x, y, z) is the position
1713 of the PMTs and the last parameter represent the type of the PMT. It's 1 for LPMT and
1714 -1 for SPMT

1715 Q_m and t_m is the set of charges and time of the PMT belonging the mesh m .
1716 (X_m, Y_m, Z_m) i the position of the center of the geometric region represented by the
1717 mesh m

1718 $(\langle X \rangle, \langle Y \rangle, \langle Z \rangle)$ is the position of the charge barycenter, $\sum Q$ the sum of the collected
1719 charge in the detector and P_l^h is the relative power of the l th harmonic. See annex [B](#) for
1720 details.

1721 Since our different nodes do not have the same number of features, they exist in distinct spaces.
1722 Traditional graph neural networks only handle homogeneous graphs, where the nodes and edges
1723 have the same number of features at each layer. Therefore, the libraries and publicly available
1724 algorithms we found were not suited to our needs. As a result, we had to develop and implement a
1725 custom message-passing algorithm capable of handling our heterogeneous graph.

1726 5.2 Message passing algorithm

1727 As introduced in previous section and in the tables [5.1](#) and [5.2](#), our graphs nodes and edges will
1728 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

1717 message passing function. We thus need to define a message passing function for each transition
1718 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

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

1719 to update the nodes after each layers. Following the illustration in figure 5.3, for each transition
1720 between families or inside a family we need an aggregation, a message and an update function. For
1721 the aggregation, we use the sum. We use the same, simple, formalism for every ϕ_u :

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

1722 using the Einstein summation notation. The second order tensor, or matrix, I_i^n is holding the nodes
1723 informations with i the node index and n the feature index. n represent the features of the previous
1724 layer and n' the features of this layer.

1725 $A_{i',e}^i$ is the adjacency tensor, discussed in the previous section, representing the edges between the
1726 node i' and the node i , each edges holding the features indexed by e . If the edge does not exist, the
1727 features are set to 0. This choice is justified by the linearity of the operation in equation 5.6 : whatever
1728 the weights, when multiplied by 0 the results is 0 and the sum result is unchanged.

1729 The learnable parameters are composed of:

- 1730 — The third order tensor $W_n^{e,n'}$ which represent the passage from the previous combined feature
1731 space between the node and the edge features $n \otimes e$, the previous layer, to the current space
1732 n' , this layer.
- 1733 — The first order tensor $B^{n'}$ which is a learnable bias on the new features n' .
- 1734 — The second order tensor $S_n^{n'}$, which can be viewed as a self loop relation where the node update
1735 itself based on the previous layer informations, going from the previous space n to the current
1736 space n' .

1737 If a node have neighbours in different families, the different IAW coming from the different families

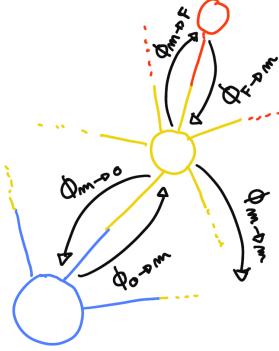


FIGURE 5.3 – Illustration of the different update function needed by our GNN

1738 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

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

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

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

1739 We thus have a S , W and B for each of the ϕ_u function we defined above. The IAW sum can be
 1740 seen as the ϕ_m function and $IS + B$ as the second part of the ϕ_u function. Eq 5.5 gave the generic
 1741 form of message passing : to update a node i , one first combines informations from the surrounding
 1742 nodes and edges and then combine the result ($\square_j \phi_m$) with the current features of node i . Many
 1743 practical ways to combine can be tried. In our implementation of message passing (Eq. 5.6 and 5.7)
 1744 the latter combination is the simple sum of the former (IAW , the equivalent of $\square_j \phi_m$) with a linear
 1745 combination of the current features of node i ($IS + B$).

1746 Interestingly, the number of learnable weight in those layer is independent of the number of nodes
 1747 in each family and depends solely on the number of features on the nodes and the edges.

1748 The expression above only update the node features. We could update the edges, using the results of
 1749 ϕ_m for example, but for technical simplicity we only update the nodes and keep the edges constant.
 1750 Preserving the edges after each layers allow to share the adjacency matrix between all layers, saving
 1751 memory and computing time.

1752 This operation of message passing is the constituent of our message passing layers, designed in this
 1753 work as *JWGLayer*, each of them owning their own set of parameter W , S and B . To those layers, we
 1754 can adjoin an activation function such as *PReLU*

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

1755 5.3 Data

1756 For this study we will be using a 1M positrons event dataset, uniformly distributed in energy with
 1757 $E_k \in [0, 9]$ MeV and uniformly distributed in the detector. Those events come from the JUNO

1758 official simulation version J23.0.1-rc8.dc1. All the event are *calib* level, with simulation of the physics,
 1759 electronics, digitizations and triggers. 900k events will be used for the training, 50k for validation
 1760 and loss monitoring and 50k for the results analysis in section 5.7. Each events is between 2k and
 1761 12k fired PMTS, resulting in fired nodes being the largest family in our graphs in all circumstances
 1762 as illustrated in figure 5.4c.

1763 As expected, by comparing the scale between the figure 5.4a and 5.4b we see that the LPMT system
 1764 is predominant in term of informations in our data. The number of PMT hits grow with energy but
 1765 do not reach 0 for low energy event due to the dark noise contribution which seems to be around
 1766 1000 hits per event for the LPMT system (left limit of figure 5.4a) and around 15 hits per event for the
 1767 SPMT system (left limit of figure 5.4b) which is consistent with the results show in section 4.1.2.

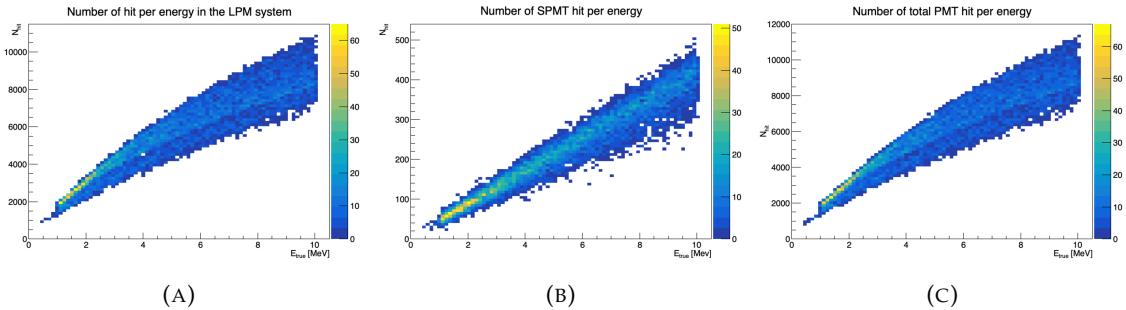


FIGURE 5.4 – Distribution of the number of hits depending on the energy. **On the right:** for the LPMT system. **In the middle :** for the SPMT system. **On the left:** For both system.

1768 The structure seen in the distribution in figure 5.4a comes from the shape of the number of hits
 1769 depending on the radius as shown in figures 5.5a and 5.5b where the number of hit decrease with
 1770 radius. It is important to understand that this is not representative of the number of PE per event
 1771 and the decrease in hits over the radius means that the PE are just more concentrated in a smaller
 1772 number of PMTs.

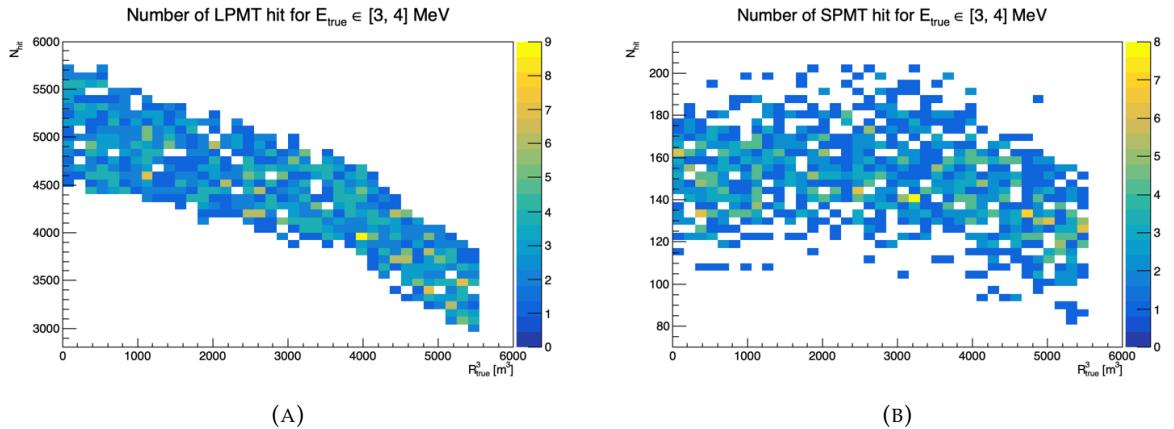


FIGURE 5.5 – Distribution of the number of hits depending on the radius. **On the right:** for the LPMT system. **On the right :** for the SPMT system. To prevent the superposition of structure of different scales we limit ourselves to the energy range $E_{true} \in [0, 9]$.

1773 No quality cut is applied here, we rely only on the trigger system. It means that event that would not
 1774 trigger are not present in the dataset but for events that triggered twice, it happens rarely, the two
 1775 trigger are considered as two separate event.

1776 5.4 Model

1777 In this section, we discuss the different layers that compose the final version of the model. The number
 1778 of layers, their dimensions, and their arrangement were fine-tuned through multiple iterations.
 1779 As mentioned earlier, each JWGLayer is defined by the number of features on the nodes and edges of
 1780 the output graph, assuming it takes as input the graph from the previous layer. For simplicity, when
 1781 discussing a graph configuration, it will be presented as follow: { N_f , N_m , N_{IO} , $N_{f \rightarrow m}$, $N_{m \rightarrow m}$, $N_{m \rightarrow f}$
 1782 } where

- 1783 — N_f is the number of feature on the fired nodes.
- 1784 — N_m is the number of features on the mesh nodes.
- 1785 — N_{IO} is the number of features on the I/O node.
- 1786 — $N_{f \rightarrow m}$ is the number of features on the edges between the fired and mesh nodes.
- 1787 — $N_{m \rightarrow m}$ is the number of features on the edges between two mesh nodes.
- 1788 — $N_{m \rightarrow f}$ is the number of features on the edges between the mesh nodes and the I/O node.

1789 Because we do not change the number of features on the edges, we can simplify the notation to { N_f ,
 1790 N_m , N_{IO} }. As an example, the input graph configuration, following the tables 5.1 and 5.2 is { 6, 8, 13,
 1791 5, 8, 5 } or, without the edge features, { 6, 8, 13 }.

1792 The final version of the model, called JWGV8.4.0 is composed of

- 1793 — An JWGLayer, converting the input graph { 6, 8, 13 } to { 64, 512, 2048 } with a PReLU activation
 function.
- 1794 — 3 resnet layers, each of them composed of
 - 1796 1. 2 JWG layers with a PReLU activation function. They do not change the dimension of the
 graph
 - 1797 2. A sum layer that sums the features in the input graph with the one computed from the
 JWG layers
- 1800 — A flatten layer that flatten the features of the I/O and mesh nodes in a vector.
- 1801 — 2 fully connected layers of 2048 neurons with a PReLU activation function.
- 1802 — 2 fully connected layers of 512 neurons with a PReLU activation function.
- 1803 — A final, fully connected layer of 4 neurons acting as the output of the network.

1804 A schematic of the model is presented in figure 5.6.

1805 We use the Mean Square Error (MSE) for the loss

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

1806 as it was the best resulting loss in Chapter 4.

1807 5.5 Training

1808 The optimizer used for training is the Adam optimizer and default hyperparameters ($\beta_1 = 0.9$,
 1809 $\beta_2 = 0.999$ and $\epsilon = 1e-8$) with a learning rate $\lambda = 1e-8$. The training last 200 epochs of 800
 1810 steps. We use a batch size of 32, the largest we can have with 40GB of GPU ram. The learning rate
 1811 is constant during the first 20 epochs then exponentially decrease with a rate of 0.99. We save two
 1812 set of parameters, the set of parameters the set that yield the lowest validation loss and the set of
 1813 parameters at the end of the training. The validation is computed over a single batch.

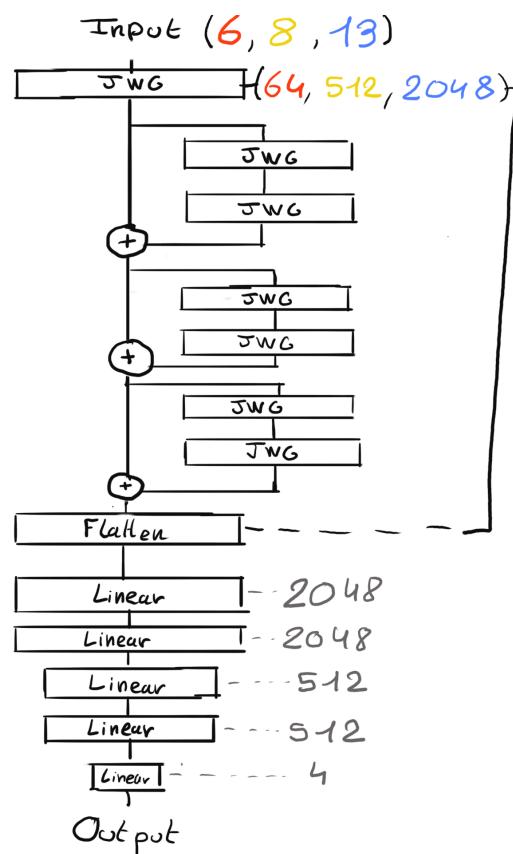


FIGURE 5.6 – Schema of the JWGv8.4.0 architecture, the colored triplet is the graph configuration after each JWG layers

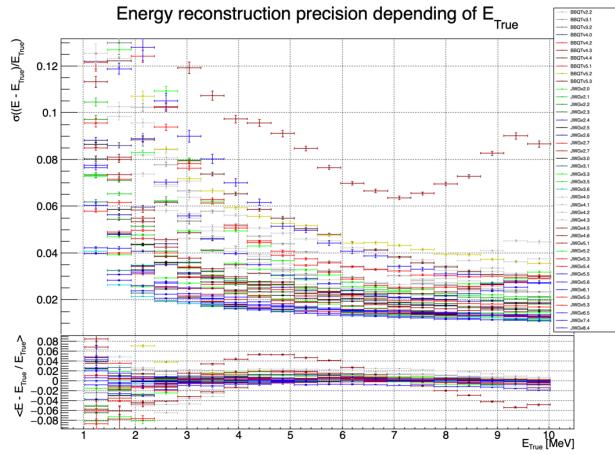


FIGURE 5.7 – Energy reconstruction depending on the true energy for samples of the different versions of the GNN

1814 5.6 Optimization

1815 The GNN model presented in previous sections is the result of a long work of optimization. Indeed,
 1816 the innovative architecture we propose left us with an infinity of possible configurations with no
 1817 guidance from prior works in literature nor in JUNO.

1818 In the end, more than 60 different configurations have been tested. This effort is illustrated on Figure
 1819 5.7¹, where the 40 configurations are compared in their ability to reconstruct the positron energy.
 1820 Although all configurations share the fundamental principles we base our innovative architecture
 1821 on (three different kinds of nodes and edges, usage of raw level features on some of them, usage of
 1822 higher level data on others, division of JUNO’s surface into regional pixels to form mesh nodes, the
 1823 very large number of edges connected to each mesh node, etc.), performances can vary a lot between
 1824 our first attempts (far beyond any acceptable energy resolution, and not even on this figure) and
 1825 recent ones. Therefore: the precise way to choose hyperparameters mattered a lot, regardless of the
 1826 relevance of the global architectural principles.

1827 The spectacular improvement between early and later configurations also explains the length of this
 1828 process : for long we hoped we would finally reach the classical performance, and it was tempting
 1829 to test yet another configuration.

1830 5.6.1 Software optimization

1831 A substantial effort was devoted to the data processing workflow. Transforming JUNO simulation
 1832 outputs into graphs is a computationally expensive task. Furthermore, due to the ever-changing
 1833 nature of the graph dimensions and features during optimization, preprocessing JUNO’s files by
 1834 precalculating the graphs and then reading them from files was not viable, as it would require a
 1835 large amount of disk space to store events for each version of the graph.

1836 Therefore, the software does not rely on preprocessed data and instead computes the observables,
 1837 adjacency matrix, etc., during training. This data processing is performed in parallel on the CPU.
 1838 The raw data comes from ROOT files produced by the collaboration software, and the Event Data
 1839 Model (EDM), used internally by the collaboration [76], had to be interfaced with our software,
 1840 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.

1841 power calculation, we migrated from the Healpix library to Ducc0 [77] for more precise control over
1842 multithreading.

1843 5.6.2 Hyperparameters optimization

1844 The first kind of hyper-parameters that received a lot of effort concern the network's detailed archi-
1845 tecture:

- 1846 — Message passing layers where originally not JWG layers, we started by using small FCDNN
1847 in place of ϕ_u and ϕ_m . Due to low performances and memory consumption issues, we pivoted
1848 to the message passing algorithm presented in section 5.2.
- 1849 — The ResNet architecture was brought after issue with the gradient vanishing.
- 1850 — The number of layers was varied between 5 and 12.
- 1851 — The number of node features after each given message passing layer (64, 512, 2048 in the final
1852 version) was varied.
- 1853 — The Final FCDNN after the message passing layers is not present in all versions.
- 1854 — At some point, the PReLU activation function replaced the ReLU function.

1855

1856 For some of them, software work was necessary. In any case, each configuration required a training
1857 of about 90h. Adding the analysis time necessary to the verification of its performance and the
1858 comparison with other versions, one understands the number of tests had to be limited.

1859 Other hyperparameters were also tested :

- 1860 — The higher level variables described in section 5.1 (powers of various spherical harmonics, \mathbb{A} ,
1861 \mathbb{A} , $(Q_{m1} - Q_{m2})/(Q_{m1} + Q_{m2})$) were added progressively. Notice that our choice to focus
1862 our search on this kind of variables is also due to the fact that JWGLayer involves linear
1863 operations. It is therefore difficult for such a network to propose variables of this kind among
1864 the node features learned layers after layers (i.e. it's difficult for the network to understand
1865 these variables are important, or only after many layers).
- 1866 — Time allocated to training, the Learning Rate, the size of batches, etc.
- 1867 — The number of pixels (ie of mesh nodes) was varied between 192 and 768.
- 1868 — Several definitions loss functions where tried. In particular, we tried some focussed only on
1869 the E resolution, only on the vertex resolution (R) or trying to optimize both.

1870

1871 To make a long story short, each new configuration was the result of our reflections after having
1872 analysed the previous configurations, or after having thought over again about JUNO's detailed
1873 response to energy deposits – seeking for variables that could help the GNN.

1874 Another, quite common, approach was in principle possible : a random search. However, due to the
1875 extensive training time, up to 90h per training, the heavy memory consumption of the models that
1876 would often exceed the 20GB limit of the V100, this approach was not realistic in our case, though we
1877 were able to extend the memory limit to 40GB thanks to a local A100 GPU card available at Subatech.

1878 5.7 performance of the final version

1879 The reconstruction performance of "JWGv8.4" are presented in figures 5.8, 5.9, 5.10 and compared to
1880 the "Omilrec" algorithm, the official IBD reconstruction algorithm in JUNO. Omilrec is based on the
1881 QTMLE reconstruction method that was presented in section 2.6.

1882 This comparison required to use a consistent definition of E_{true} . This is not trivial since at JUNO,
1883 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.

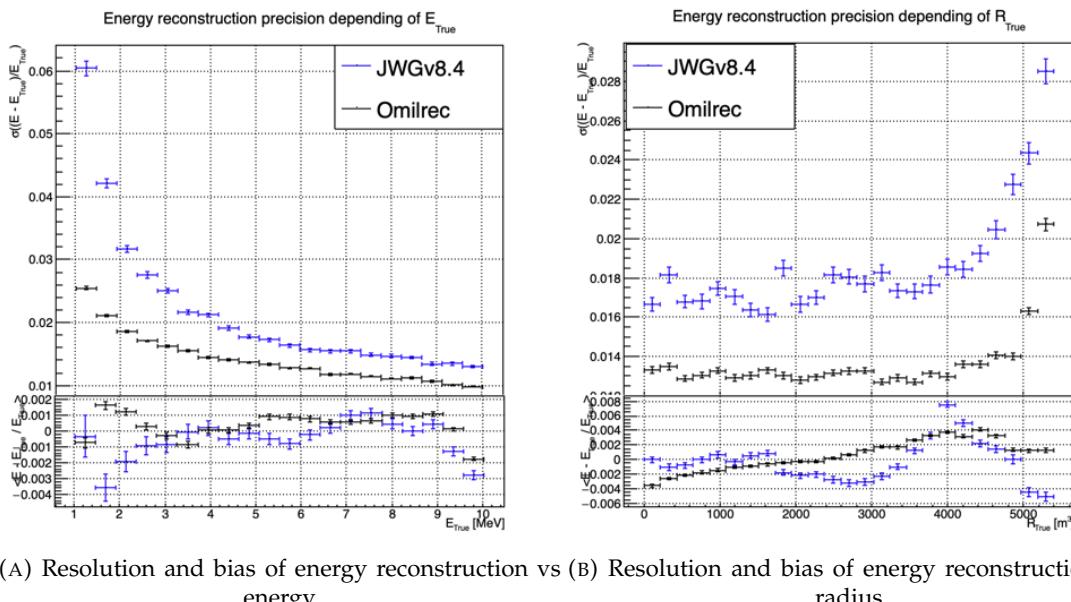


FIGURE 5.8 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.6, JWGV8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.

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.

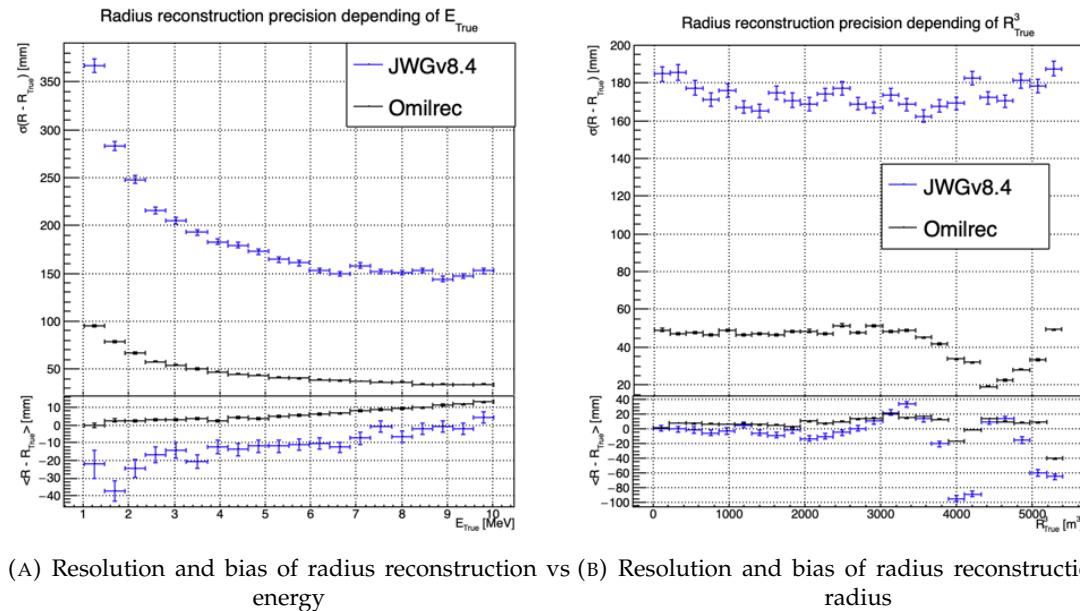


FIGURE 5.9 – Reconstruction performance of the Omilrec algorithm based on QTMLR presented in section 2.6, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.

1906 Finer granularity

1907 We tried to recover some redundancy by increasing the number of mesh nodes from 198 to 768. The
 1908 improvement we observed was small, and did not allow to get close to OMILREC's performance.

1909 To explore further in this direction, we would ideally try 3072 pixels (the next HEALPIX rank).
 1910 However, this is not possible for our GNN due to hardware limitations, mainly the available GPU
 1911 memory. Instead, we discussed the problem with Gilles Grasseau, calculus research engineer with
 1912 whom we collaborate on the subject of ML reliability (see Chapter 6). In the framework of this ac-
 1913 tivity, Gilles needs to develop reconstruction algorithms to be "attacked" by a prototype Adversarial
 1914 NN. One of them is a pseudo-spherical CNN using oriented filters, called HCNN.

1915 To produce its input image, this algorithms split the Sphere into 3072 pixels. Each channel of this
 1916 image is an aggregation of the (Q, t) values found in all the PMTs. The charge are summed and the
 1917 lowest time is kept. The performance of this algorithm can be seen on Figures ?? and ??, compared
 1918 to OMILREC. With 3072 pixels, the performance of HCNN does not match that of OMILREC, but is
 1919 closer to it than our GNN. The granularity of the pixels, and the way to summarize the individual
 1920 PMTs information when going from 17000 LPMTs to only 3072 pixels indeed seems to play a role.

1921 This is consistent with the results obtained by the first GNN tried at JUNO on reactor neutrinos
 1922 (already described in section 2.6.4). It used 3072 pixels, and also obtained an uncompetitive R
 1923 reconstruction.

1924 Information reduction, from fired to Meshes

1925 The problem described above is somehow classical. ML algorithms, ideally, would start from the full
 1926 information present in the detector, and learn to reduce it optimally.

1927 In cases where only 3072 pixels can be used instead of the complete information from 17000 PMTs,
 1928 one needs to understand how to combine the individual from the 5 or 6 PMT found in each pixel into

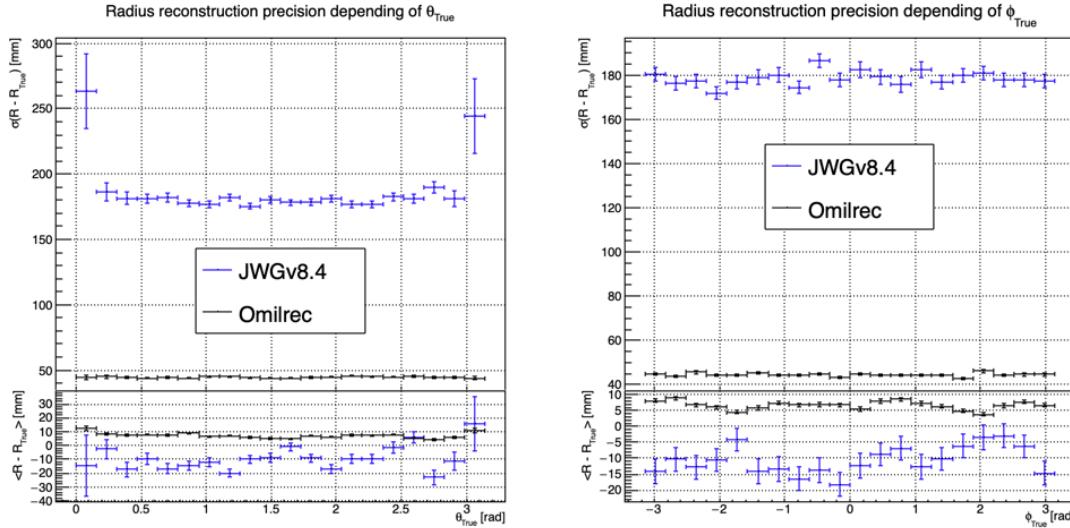
(A) Resolution and bias of radius reconstruction vs θ (B) Resolution and bias of radius reconstruction vs ϕ

FIGURE 5.10 – Reconstruction performance of the Omilrec algorithm based on QTMLE presented in section 2.6, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.

1929 pixel-level features, without loosing important information.

1930 In the case of our GNN, we hoped that by connecting each mesh node to its corresponding 5 or 6
1931 fired nodes, we could keep the full information. In reality, it seems that the message passing between
1932 fired and mesh does not work efficiently. When nodes are updated by the first (may be also by the
1933 subsequent) layer, the new mesh features might be dominated by the original features in the second
1934 column of tables 5.1, themselves a simple version of aggregation. Layer after layer, we might be
1935 limited to that level of time information, lacking time redundancy.

1936 We have verified this by testing version of the GNN in which the link between fired and mesh was
1937 cut, or in which no time info was included among the fired nodes features. It had only a small effect
1938 which seems to confirm a problem in the way the full information, from all the individual PMTs, is
1939 used by our GNN.

1940 Possible improvements

1941 It appears that the network is unable to aggregate the timing information correctly. While this could
1942 be addressed by using a finer segmentation, with more mesh nodes, improvements might also arise
1943 from refining the message-passing algorithm. The algorithm presented in this thesis is still quite
1944 basic, relying on a simple linear combination of features. We have seen through examples in CNNs,
1945 GNNs, and other architectures, both in research and industry, that specializing the network — for
1946 instance, by incorporating convolutional filters — can lead to improvements that were previously
1947 unattainable with simpler FCDNNs. Applying this approach to the message-passing algorithm, by
1948 utilizing a GNN with a more advanced message-passing, could yield better results.

1949 Regarding the timing information, we provided high-level features, assuming this would assist the
1950 neural network in converging to the solution. However, by offering such information upfront,
1951 the GNN might be taking the “easy” path, settling for a local and broader minimum, rather than
1952 extracting the features that could lead to better performance.

1953 If there are difficulties in transferring information between the fired and mesh nodes, it may stem

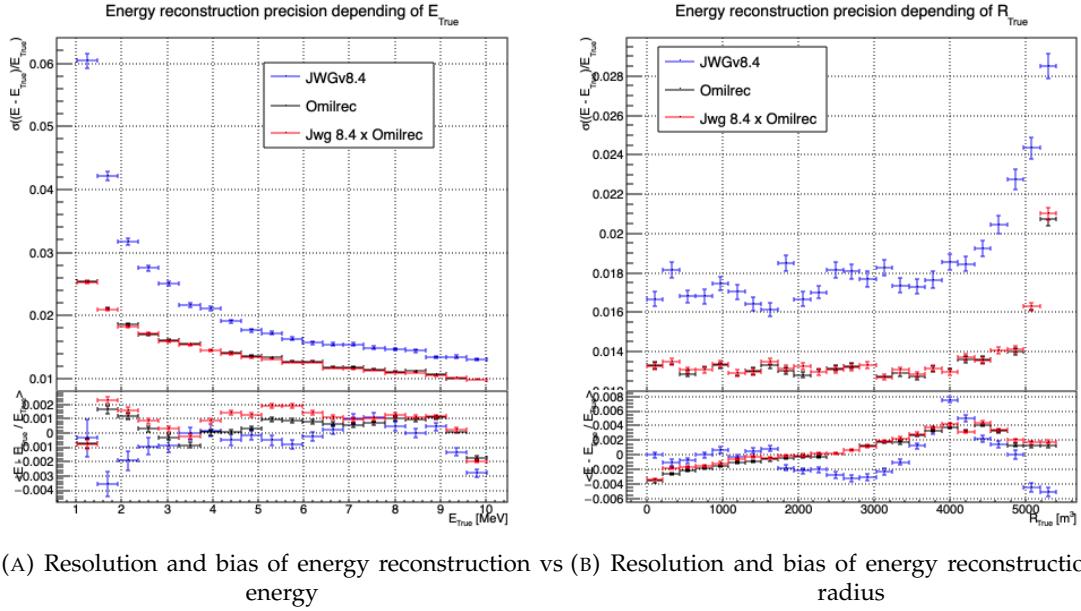


FIGURE 5.11 – Reconstruction performance of the Omilrec algorithm, JWGV8.4 and the combination between the two using the optimal variance estimator presented in annex A.2. The top part of each plot is the resolution and the bottom part is the bias.

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.

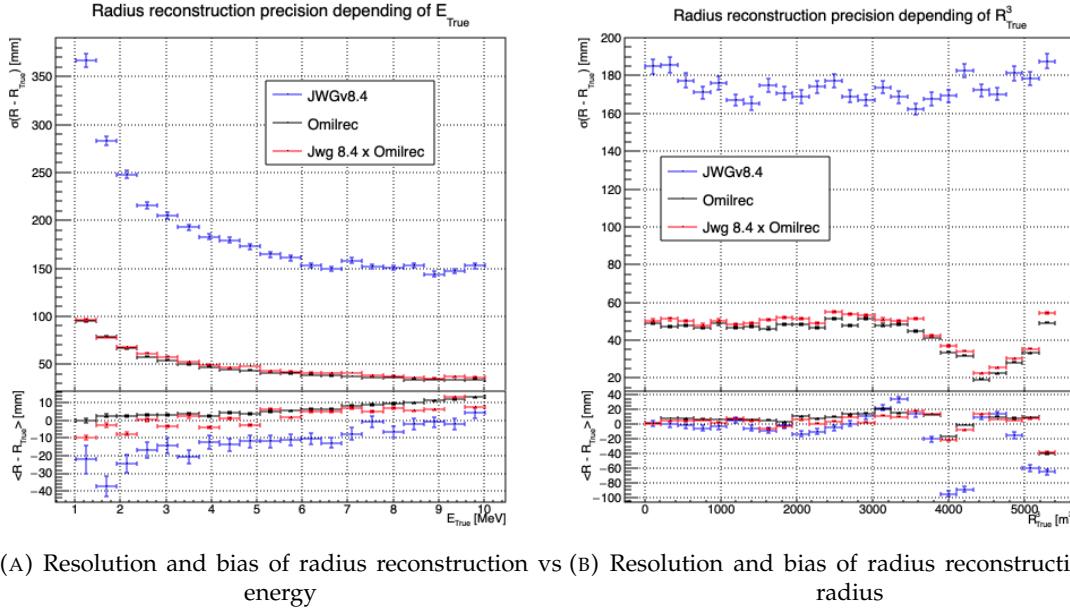
5.8 Conclusion

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 2.6.

Machine learning algorithms were developed to challenge this classical approach, and they are presented in Section 2.6.4. 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.

These three types, or families, of nodes do not have the same number of features, resulting in a het-



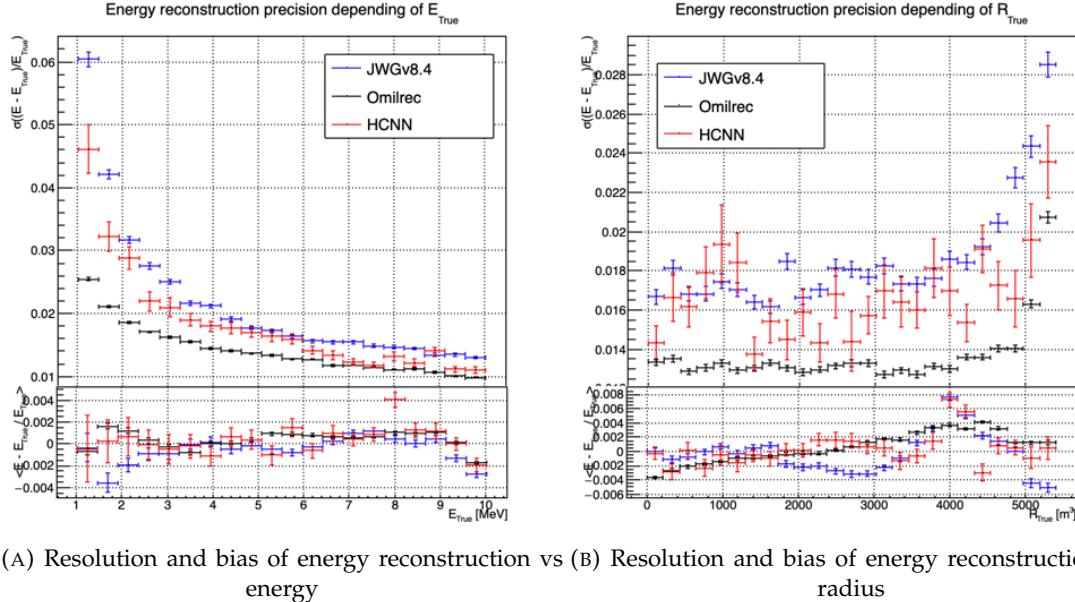
(A) Resolution and bias of radius reconstruction vs (B) Resolution and bias of radius reconstruction vs
energy radius

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.

erogeneous graph. Publicly available algorithms for graph processing are designed for homogeneous graphs, so we had to develop a custom algorithm adapted to heterogeneous graphs.

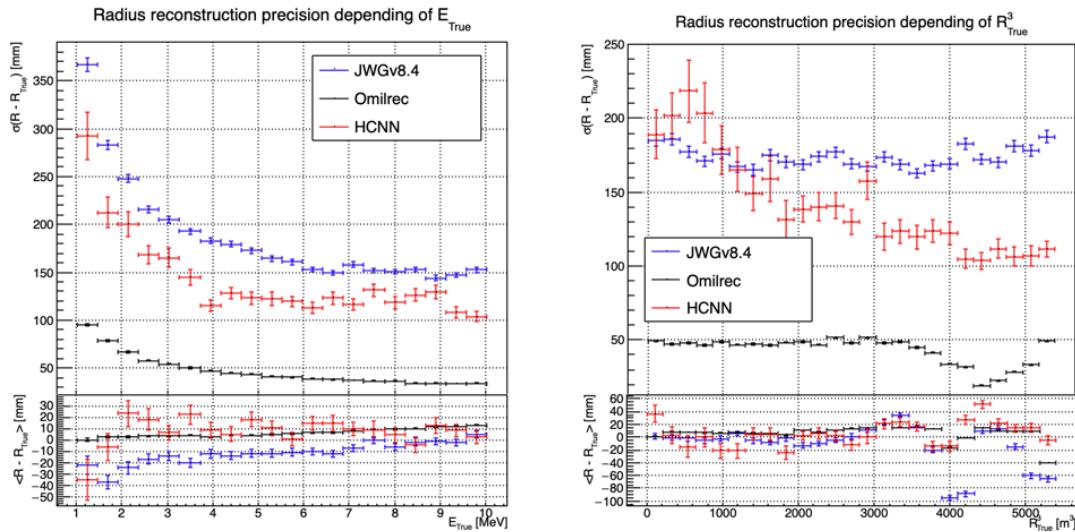
This GNN required significant technical development, but the results are not at the level of the classical algorithm. The tests we conducted suggest that the problem may lie in the aggregation of raw information from the fired nodes onto the mesh nodes, as removing the fired nodes does not degrade the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

Perhaps by incorporating rawer information, such as the waveform, refining the message-passing algorithm, or adjusting the features on the different nodes, we could match the precision of the classical algorithm. However, it is also possible that deeper, more radical changes are needed to 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 2.6, 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 2.6, 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.

1990 **Chapter 6**

1991 **Reliability of machine learning
methods**

1992

1993

"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

1994

Contents

1995

1996

1997

1998

1999

2000

2001

2002

2003

2004

2005

2006

2007

2008

6.1 Motivations	92
6.2 Method	92
6.3 Architecture	92
6.3.1 Adversarial Neural Network	92
6.3.2 Reconstruction Network	93
6.3.3 Training	93
6.4 Results	93
6.4.1 Back to identity	94
6.4.2 Breaking of the reconstruction	94
6.5 Conclusion and prospect	94

2009 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 2010 uncertainties could lead to the impossibility to run the measurements, or even worse, wrong our mass 2011 ordering measurements. While the liquid scintillator technology is well known and straightforward, 2012 this is the first time it is deployed to such scale, and for such precision. This novelty brings its fair 2013 share of elements, effects or assumption, that, if they were to be overlooked, could cause issue. 2014

2015 We already showed a large variety of reconstruction algorithms, OMILREC for LPMT reconstruction 2016 in section 2.6, numerous machine learning algorithms in section 2.6.4 and our own work in chapters 2017 4 and 5. Those algorithms were compared to each other based on their performance as in [42] but 2018 we are the first that looked into the correlation between the reconstruction. The combinations of 2019 algorithms shown in Chapter 4 and Chapter 5 show that some information eludes the algorithms. 2020 We used this fact to try to improve our performance but this could also lead the algorithm to being 2021 vulnerable to some effect that could affect the detector and wrong the measurements.

2022 The search for such effect could be done by hand, but the process would be tedious. We propose in 2023 this thesis a machine learning method to probe for those effects. In section 6.1, I delve further in the 2024 motivations of this work. In section 6.2, I describe the method behind the algorithm. In section 6.3 2025 I detail the architecture of our algorithm and in section 6.4 the results of it. Finally, in section 6.5, I 2026 conclude and discuss about the prospect and possible improvements to bring to this work.

2027 6.1 Motivations

2028 As introduced above, JUNO needs a very good understanding of the biases and effects affecting its
 2029 reconstruction as a small bias could wrong the mass ordering measurement. To calibrate those biases
 2030 and effect, JUNO rely on multiples sources that will be located at various point in the detector. The
 2031 calibration strategy was already discussed in section 2.3 and show calibrations sources of gammas,
 2032 neutrons and positrons, with the catch that the positrons will annihilate inside the encapsulation and
 2033 only the two 511 keV gammas will be seen. All those sources will be located at the center of the
 2034 detector, impervious to non-uniformity.

2035 A second, natural, source will be used for calibration: The ^{12}B spectrum. The ^{12}B is a cosmogenically
 2036 produced isotope through the passage of muons inside the LS. The ^{12}B decays via β^- emissions with
 2037 a Q value of 13.5 MeV with more than 98% of the decay resulting in ground state ^{12}C . The ^{12}B event
 2038 will be cleanly identified by looking for delayed high energy β events after an energetic muon. Due
 2039 to its natural causes, the ^{12}B events will be uniformly distributed in the detector. The calibration
 2040 strategy consist in fitting the energy spectrum of ^{12}B with the results of the simulation to adjust the
 2041 simulation parameters.

2042 We see that, while the calibration strategy is pretty complete, its missing a few points. First, none
 2043 of the calibrations sources considered are positrons. While electrons and positrons events should
 2044 be pretty similar in their interaction with the electronic cloud of the LS atoms, electron events are
 2045 missing the two annihilations γ and the potential of forming a positronium [78]. The topology of the
 2046 event thus differ of the order of magnitude of our reconstruction performance, a few nanoseconds
 2047 for the energy deposit and positronium annihilation against a time transit spread between 3 and 6
 2048 ns depending on the PMT type [79–81] and the γ will travel distances of the order of magnitude
 2049 of the typical LPMT resolution of 8 cm (see section 2.6). Moreover, where for calibration sources
 2050 the localization will be well known, the individual truth of ^{12}B will be unknown. We thus need to
 2051 compare our model to higher order observables such as energy distribution more than individual
 2052 comparison.

2053 If there is potential failure point in those considerations, we need to search for them efficiently.

2054 6.2 Method

2055 All of the considerations could hide potential unknown or undetected effect that could lead to issue
 2056 in the mass ordering analysis. But, while we have idea from where the issue could come, the
 2057 production by hand of event perturbations that would not show in the calibration would be tedious.
 2058 That's why we propose to use a Neural Network to produce those perturbations if they exists. A
 2059 schematic of the concept is presented in figure 6.1.

2060 6.3 Architecture

- 2061 — Expliquer la problematique dans l'architecture
- 2062 — Ambition de pouvoir etre appliqu  a toutes les methodes, pas que NN
- 2063 — Pb technique: descente de gradient
- 2064 — Pr senter la loss

2065 6.3.1 Adversarial Neural Network

- 2066 — Decrire l'architecture de l'ANN

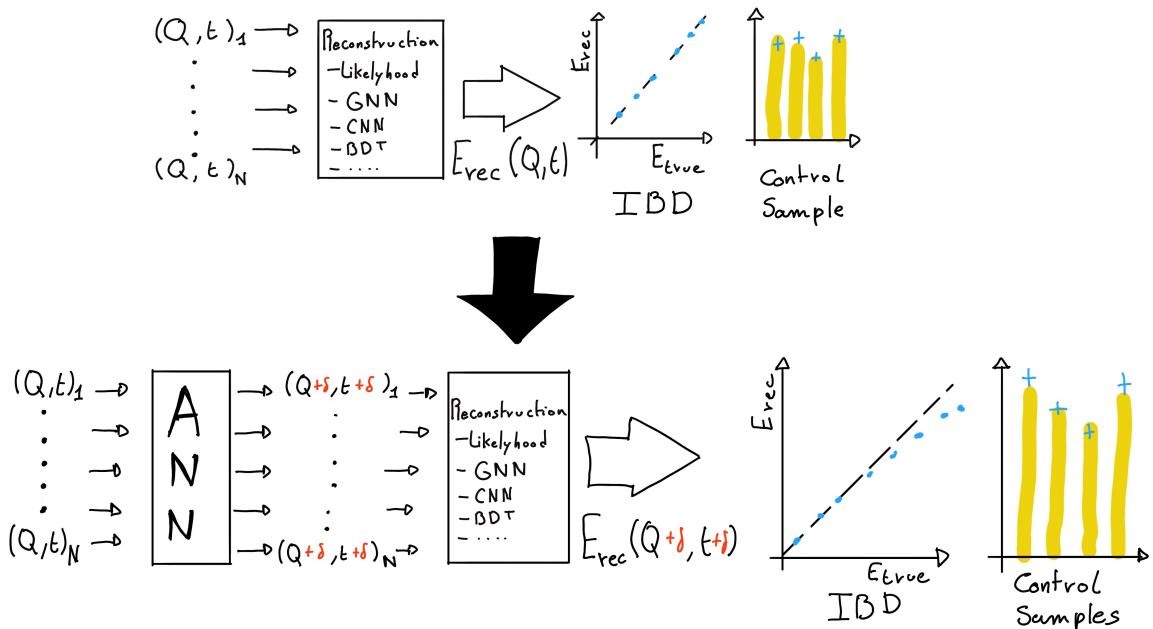


FIGURE 6.1 – Schema of the method to discover vulnerabilities in the reconstruction methods

6.3.2 Reconstruction Network

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

6.3.3 Training

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

Hyperparameter optimization

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

6.4 Results

- Voir slide Gilles

2085 **6.4.1 Back to identity**

2086 **6.4.2 Breaking of the reconstruction**

2087 **6.5 Conclusion and prospect**

- 2088 — Not enough
2089 — Probably guide the ANN

2090 **Chapter 7**

2091 **Joint fit between the SPMT and LPMT
spectra**

2093 “We demand rigidly defined areas of doubt and uncertainty!”

2094 Douglas Adams, *The Hitchhiker’s Guide to the Galaxy*

Contents

<small>2095</small> 7.1 Motivations	<small>2096</small> 96
<small>2097</small> 7.1.1 Discrepancies between the SPMT and LPMT results	<small>2098</small> 96
<small>2098</small> 7.1.2 Charge Non-Linearity (QNL)	<small>2099</small> 97
<small>2099</small> 7.2 Approach	<small>2100</small> 98
<small>2100</small> 7.2.1 Data production	<small>2101</small> 98
<small>2101</small> 7.2.2 Individual fits	<small>2102</small> 99
<small>2102</small> 7.2.3 Joint fit	<small>2103</small> 100
<small>2103</small> 7.2.4 Data and theoretical spectrum generation	<small>2104</small> 102
<small>2104</small> 7.2.5 Limitations	<small>2105</small> 102
<small>2105</small> 7.3 Fit software	<small>2106</small> 103
<small>2106</small> 7.3.1 IBD generator	<small>2107</small> 103
<small>2107</small> 7.3.2 Fit	<small>2108</small> 105
<small>2108</small> 7.4 Technical challenges and development	<small>2109</small> 105
<small>2109</small> 7.5 Results	<small>2110</small> 106
<small>2110</small> 7.5.1 Validation	<small>2111</small> 106
<small>2111</small> 7.5.2 Covariance matrix	<small>2112</small> 110
<small>2112</small> 7.5.3 Statistical tests	<small>2113</small> 114
<small>2113</small> 7.6 Conclusion and perspectives	<small>2114</small> 116

2116 JUNO is an experiment of precise measurements, where we try to observe small fluctuation in the energy spectrum and with the goal to achieve sub-percent precision on the oscillation parameters measurement. A precise and complete understanding of the reconstruction and detector effects is thus crucial. The challenge reside in the technology used in the detector, which, while based on well known technology: scintillator observed by PMT, is being deployed on a scale never seen before, in term of scintillator volume and PMT size. Understanding every effects that goes in the detector can become extremely complicated. The ability to compare the results of the same experiment with two systems is thus extremely precious, this is the origin the dual calorimetry with the LPMT and SPMT system.

2126 The resolution and bias of the reconstruction needs to be extremely well characterized: the target resolution of 3% [50] is unprecedented and is necessary to be able to distinguish between Normal

2128 Ordering (NO) and Inverse Ordering (IO). The non-linearity uncertainty needs to be constrained
 2129 under 1% as exceeding this value, the risk appear to measure the wrong ordering [27].

2130 One of the possible source of non-linearity, which will be used as a reference in this chapter, is the
 2131 charge non-linearity (QNL) that will be discussed in next section. The dual calorimetry can address
 2132 this issue, using calibrations methods and measurements that will be employed to correct it [27].

2133 More generally, comparing the results of the two systems will allow for the detection of potential
 2134 issues on the calibration or reconstruction. This is done in this thesis by comparing directly the
 2135 spectra and oscillation parameters measurements of the two systems.

2136 The study of the independent results of the two system can provide some informations [64] but this
 2137 is missing the important correlation that should be present between the two systems: they see the
 2138 same events, in the same scintillator, they're bound to be correlated. We explore in this chapter a
 2139 preliminary study of the impact of those correlations via multiple methods and the impact of QNL
 2140 at various degrees.

2141 In the next section we will discuss the motivations behind this study. In section 7.2, I present the
 2142 approaches and assumptions in this study. In section 7.3, I present the fit framework used, and then,
 2143 in section 7.4 the technical improvement brought and the difficulties faced during the development.
 2144 To end this chapter I present the results in 7.5 and discuss the conclusions and perspectives in 7.6.

2145 7.1 Motivations

2146 7.1.1 Discrepancies between the SPMT and LPMT results

2147 As discussed in the introduction of this chapter, the SPMT and LPMT systems will observe the same
 2148 events. This mean that, after calibration, if the two system show significant differences in their results
 2149 this is the signal of potential overlook of an effect or problem. Being able to detect such differences
 2150 is thus crucial, as discussed above, even the smallest deviation from our model could lead to the
 2151 impossibility to measure the Mass Ordering (MO) or even worse, wrong our measurement.

2152 The two systems are expected to have the same sensitivity to the oscillation parameters θ_{12} and Δm^2_{21}
 2153 [11]. We will thus rely on the measurement of those two parameters to detect potential discrepancies.

2154 We could just look at the value and compare them to the estimated independent error of the two
 2155 system, but we believe and will demonstrate in this chapter that the independent study of the two
 2156 system is missing a lot of informations, and that, by taking into account the statistic and systematic
 2157 correlations between the two systems, we can produce much more powerful statistical tests.

2158 Our work in this chapter is to develop such tools. The first step is, of course, to verify that in the
 2159 case of no discrepancies, the results are coherent with the independent analysis. This will give us the
 2160 distribution of those statistical test in absence of discrepancies. When we will have real data, we will
 2161 be able to compare it to those distributions to compute a p-value characterizing the absence of those
 2162 potential discrepancies.

2163 To evaluate the power of our methods, we need to simulate a concrete difference between the two
 2164 spectra. We have decided to study a plausible effect, the Charge Non-Linearity (QNL) that is detailed
 2165 next section. But the goal of those tools is to be discrepancy agnostic, as those discrepancies could
 2166 come from a variety of source (calibration issue, insufficient simulation tuning, etc...)

2167 7.1.2 Charge Non-Linearity (QNL)

2168 The CD energy response is subject to two kinds of non-linearity, the first one is the LS response
 2169 non-linearity, where the LS photo-production is not linear with the deposited energy as illustrated
 2170 in figure 2.12a. The second one is the LPMT response non-linearity where the charge read from the
 2171 LPMT is not linear with respect to the number of collected Photo-Electrons (PE) (see section 2.3).

2172 The LS non-linearity comes from physic sources. Particle interactions in the LS will produce mainly
 2173 scintillation light, as discussed in section 2.2.2, but will also produce some Cherenkov light (< 10%
 2174 of the collected light). Both mechanisms possess intrinsic non-linearity, for the Cherenkov emission
 2175 it depends on the velocity of charged particle velocity while the scintillation photon-yield follows a
 2176 so-called Birk's law with a "quenching" effect depending on the energy and type of particle [16]. This
 2177 results in am event-wise QNL.

2178 The LPMT response non-linearity can come from sheer saturation when subject to a high photon rate
 2179 inducing a gain non-linearity or come from readout effects such as electronic noise, overshoot, the
 2180 integration time window and even the waveform algorithm. All of these effects result in a channel-
 2181 wise QNL.

2182 Precedent studies [27] suggest a model to emulate the non-linearity response that will be used in this
 2183 work. We define the channel wise non-linearity that would be applied to each LPMT readout

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

2184 where Q_{rec} is the reconstructed number of PE by the PMT, Q_{true} is true number of PE that hit the
 2185 PMT, and γ_{qnl} is a factor representing the amplitude of the non-linearity.

2186 We also define an event-wise non-linearity characterized by

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

2187 where E_{vis} is the visible energy that is collected by the detector and E_{true} is the true deposited energy.
 2188 An example of the effect of such event-wise QNL is presented in figure 7.1.

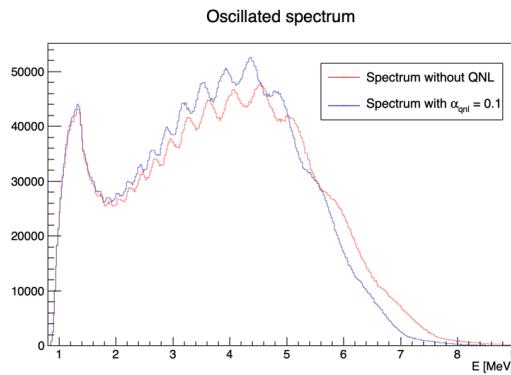
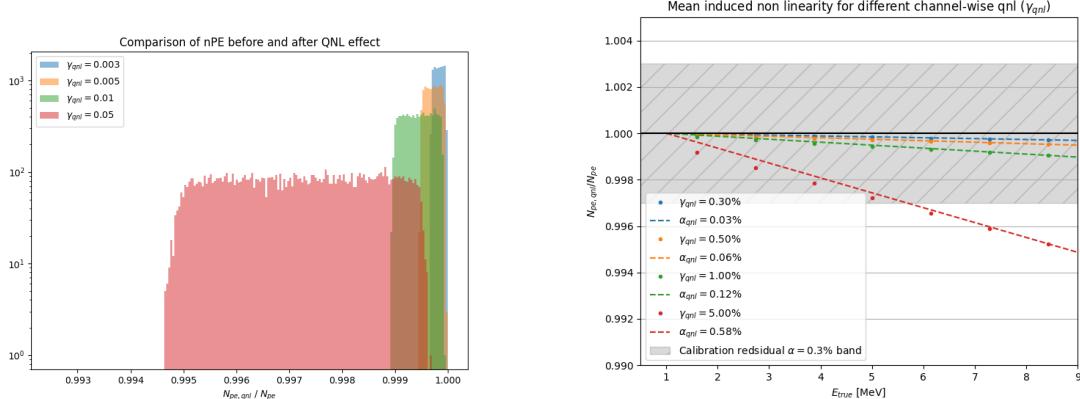


FIGURE 7.1 – Two oscillated spectra of 1e7 event expected in JUNO. In red the spectrum without supplementary QNL. In blue the same spectrum but where an event-wise QNL $\alpha_{qnl} = 10\%$ is introduced.

2189 Using 1M events from the JUNO official simulation J23.0.1-rc8.dcl (released on 7th January 2024), we
 2190 simulated events up to the photon collection in LPMTs and introduced an additional channel-wise
 2191 QNL by using the equation 7.1 to modify the number of collected photons.



(A) Distribution of ratio of collected nPE after the additional QNL over the number of nPE that would be collected for different 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 7.2a and the dashed line are the equivalent event-wise non-linearity from eq 7.2. The hatched zone is the residual non-linearity expected after calibration [29].

FIGURE 7.2

In figure 7.2a we show the distribution of the ratio $\frac{Q_{rec}}{Q_{true}}$ for central events ($R < 4\text{m}$) and different values of γ_{qnl} . In figure 7.2a, 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.

7.2.1 Data production

IBD spectra

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

2210 the reconstruction resolution for each system to their respective spectra, resulting in the final LPMT
 2211 and SPMT spectra.

2212 We will study the effect of exposure on our methods at different threshold: 100 days, 1 year, 2 year
 2213 and finally 6 years which is the nominal data taking period for the NMO analysis.

2214 These spectra are generated for different QNL, $\alpha_{qnl} = 0$ (no spectrum distortion) and for $\alpha_{qnl} \in$
 2215 $\{0.01, 0.005, 0.003, 0.002, 0.001\}$. As a reminder, the calibration guarantees a residual event-wise non-
 2216 linearity of $\alpha_{qnl} \leq 0.003$ [29].

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

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

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

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

2217 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
 2218 the covariance matrix of the LPMT + SPMT spectra. A is a transformation matrix defined as:

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

2219 Thus, $A_{ij} = 1$ if $i = j$, and $A_{ij} = -1$ if j is the SPMT bin corresponding to the i LPMT bin.

2220 This χ^2_{spe} is minimal when the statistic between the bins of the LPMT and SPMT spectra follow the
 2221 covariance matrix V . By looking at the distribution of this χ^2_{spe} when $\alpha_{qnl} = 0$ we can produce
 2222 p-values for the values found when $\alpha_{qnl} \neq 0$.

2223 Background spectra

2224 The JUNO common inputs provide only LPMT background spectra. These background spectra are
 2225 already smeared by the LPMT resolution and thus need to be regenerated to be smeared to account
 2226 for the SPMT resolution. Fortunately the SPMT resolution is greater than that of the LPMT, allowing
 2227 us to apply additional smearing to the spectrum using

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

2228 Where $S(E)$ is the SPMT spectrum, $L(E)$ the LPMT spectrum, σ_L and σ_S the LPMT and SPMT resolution
 2229 respectively. This formula is valid under the assumption that the LPMT and SPMT smearing are
 2230 gaussian and that the LPMT and SPMT have the same bias. Those two assumptions are valid in the
 2231 context of the IBD spectrum production as detailed in section 7.3.1. The demonstration of equation
 2232 7.7 can be found in annex C.

2233 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 $\theta_{12,L}, \Delta m_{21,L}^2, \theta_{12,S}$ and $\Delta m_{21,S}^2$. A is the transformation matrix that allow us to compute the covariance matrix de $\vec{\Delta}$ from V following

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

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

7.2.3 Joint fit

Standard joint fit

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

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

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

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

T is the fitted function and depend on multiple parameters

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

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

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

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

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

2259 7.1.

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

2260 θ_{13} is the parameter on which we are least accurate. It's fixed to nominal value to prevent degeneracy
2261 (table 7.1).

2262 The covariance matrix is produced from a correlation matrix C

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

2263 where σ_i is the uncertainty on the number of event in the i th bin. We consider in this study that the
2264 content of each bin follow a Poisson statistic, thus the uncertainty is $\sigma_i = \sqrt{N_i}$ where N_i is the content
2265 of the i th bin. The bin content used for the uncertainty can come from two sources: the data and the
2266 theoretical spectra $\sigma_i = \sqrt{D_i}$ (Pearson test) and $\sigma_i = \sqrt{T_i}$ (Neyman test). Precedent studies have
2267 show that both Pearson and Neyman tests show bias at low statistic, we thus use the Pearson V test
2268 where

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

2269 and the covariance matrix V is computed using the data spectrum for the uncertainty.

2270 The estimation of the covariance is crucial in this study as the strength of this test rely on the sys-
2271 tematic and statistical correlations between the LPMT and SPMT spectrum. The generation methods
2272 and results of this matrix is detailed in section 7.5.2.

2273 Delta joint fit

2274 Using the same structure we define a second joint fit, the Delta joint fit where, in addition to every-
2275 thing that was discussed above, we add two other parameters $\delta\theta_{12}$ and $\delta\Delta m_{21}^2$ and split the theoretical
2276 $T(\theta_{12}, \Delta m_{21}^2, \dots)$ spectrum in two

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

2277 If the there is no additional distortion between the LPMT and the SPMT spectra, the fit should
2278 converge to $\delta\theta_{12} = \delta\Delta m_{21}^2 = 0$. By observing the dispersion of those parameters we can define
2279 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 $\alpha_{qnl} \neq 0$
2280 to define a p-value.

2281 The last test we explore in this thesis is to fit the same spectrum with the Standard Joint fit, that
2282 we consider as the hypothesis without distortion H_0 , and the Delta Joint fit, designated as the H_1
2283 hypothesis. By looking at the dispersion of $\chi_{joint,H_0}^2 - \chi_{joint,H_1}^2$ we can extract a sensitivity to potential
2284 distortion.

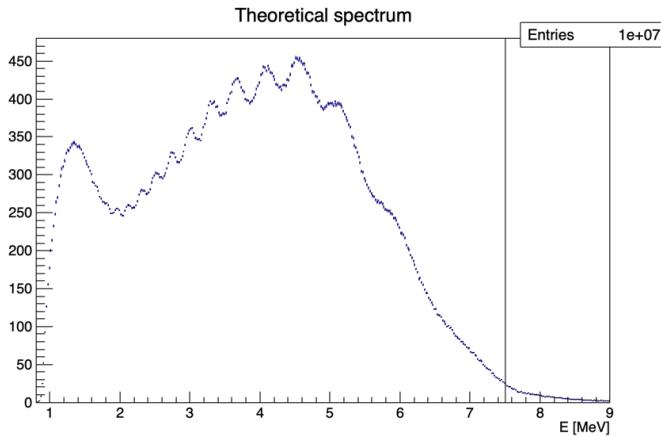


FIGURE 7.3 – Theoretical LPMT spectrum at nominal oscillation values binned using 410 bins from 0.8 to 9 MeV. It is rescaled to 6 years statistic. The black line represent the 335 bin cut

2285 7.2.4 Data and theoretical spectrum generation

2286 To implement the joint fit, we have technically two data spectra and two theoretical spectra. The data
 2287 in this study are produced using an IBD generator *IBD gen*, see section 7.3.1. The theoretical spectrum
 2288 are produced the same way as data spectrum but with much higher statistics, 10^7 events to compare
 2289 with the $\approx 10^5$ events for 6 years statistic. The two spectrum, that we get as a collection of events,
 2290 are binned in two histograms from 0.8 to 9 MeV of reconstructed energy with bins of 0.02 MeV each,
 2291 resulting in 410 bins per spectrum. An illustration of the theoretical spectrum can be found in figure
 2292 7.3. The low number of events in the tail of the spectrum can cause instability due to the low statistic,
 2293 we thus cut the spectrum at 7.5 MeV / 335 bins for the fit.

2294 All the IBD spectra presented and used in this study are produced assuming Normal Ordering using
 2295 the PDG nominal value [16] for the oscillation parameters. Those values are reported in table 7.1.

2296 7.2.5 Limitations

2297 In this work we are only working considering the statistical errors. We can ignore systematic effects,
 2298 such as effects that would affect the neutrino spectrum or the background spectrum, as they are
 2299 entirely correlated between the two systems. The details of those systematic effects can be found in
 2300 [11].

2301 Most of our results assume decorrelated detection effects between the SPMT and LPMT systems.
 2302 Their respective reconstruction effects are simulated using simple gaussian drawing on the resolution,
 2303 independently from the event position. This approach was used in previous sensitivity and
 2304 precision studies [11, 82]. The potential effect of those reconstruction effects and a first attempt to
 2305 take them into account are explored in section 7.5.2.

2306 Even if the goal of this work is to propose deformation agnostic tools, the QNL we use in this study is
 2307 simplistic as we consider event-wise, position uniform deformation. We show in figure 7.2a and 7.2b
 2308 that event-wise QNL is equivalent to the mean behaviour of channel-wise QNL but a more complete
 2309 study would simulate channel-wise deformation for each event.

7.3 Fit software

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

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

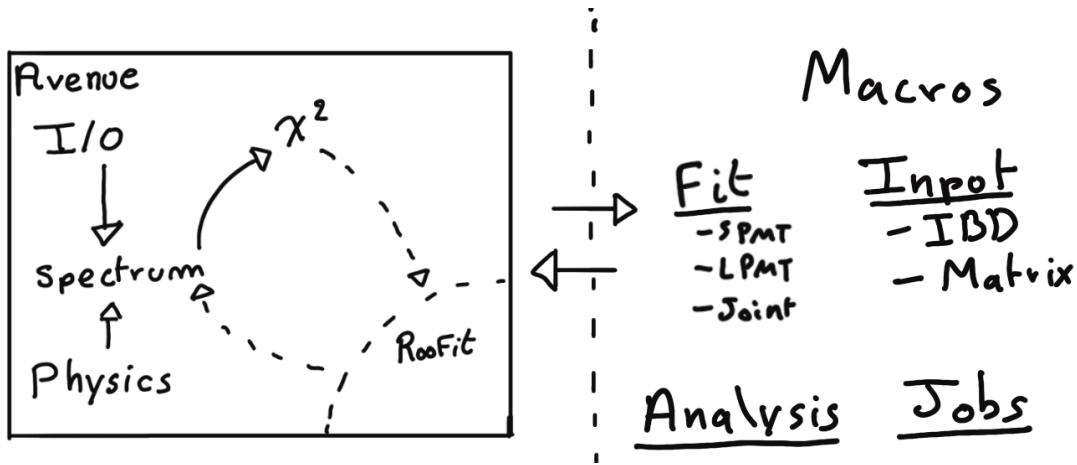


FIGURE 7.4 – Schematic description of the fit framework

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

7.3.1 IBD generator

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

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

Then for each event we

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

- (b) Summing the contribution of each isotopes following the respective fission fraction [0.58, 0.07, 0.30, 0.05] as reported in [84].
- (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).
- (d) The total spectrum is then finally adjusted by taking into account the correction of the Day Bay bump [85], adjustment due to spent nuclear fuel and due to the non-equilibrium.
4. (Optional) Compute the survival probability due to oscillation at nominal oscillation parameters value. If the neutrino does not survive, the event is rejected and the algorithm restart from step (1).
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).
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).
7. Correct the deposited energy with the expected event-wise non-linearity from [29] to obtain the visible energy E_{vis} .
8. (Optional) Add a custom non-linearity as described in section 7.1.2. This non linearity is characterized by α_{qnl} to obtain E_α .
9. Finally, using the expected resolution of the LPMT and SPMT systems, provided in the JUNO common inputs, we draw from a gaussian characterized by those resolution the reconstructed energy E_{rec} or E_{lpmt} and E_{spmt} for each systems. The resolutions are provided as ABC parameters using

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

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

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

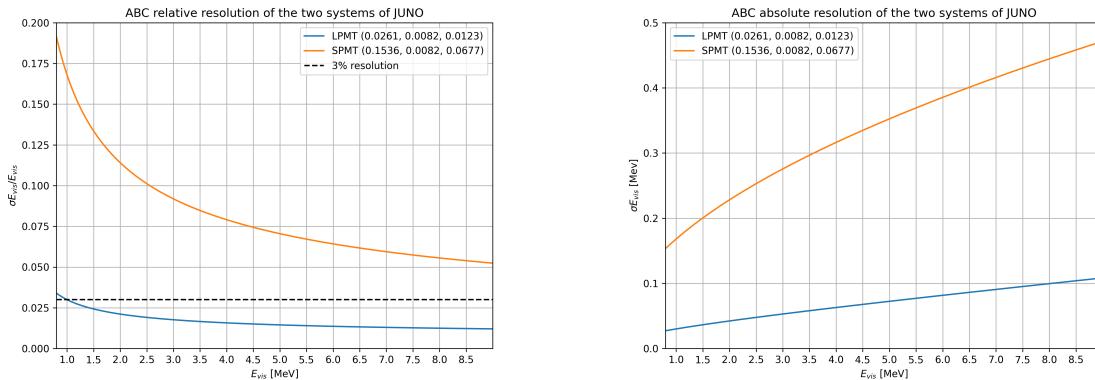


FIGURE 7.5 – Relative (On the left) and absolute (On the right) resolutions of the LPMT and SPMT systems used in this study. The number in parenthesis are the parameter A, B and C respectively for each systems.

2365 **7.3.2 Fit**

2366 The fit macro is the core of this fitting procedure. This macro is responsible for loading the fit
 2367 configuration and setup the Avenue framework. Using Avenue, it will setup the data files, theoretical
 2368 spectrum, choose the binning, χ^2 , etc... It also have the possibility to generate toys on the fly based
 2369 on the theoretical spectrum. Given this theoretical spectrum we can randomize the bin content either
 2370 by:

- 2371 1. Drawing the bin content in a Poisson distribution with the bin content as parameter.
- 2372 2. Drawing the bin content in a Gaussian distribution with the bin content as mean and variance.
 The bin content is then rounded to the nearest integer.
- 2373 3. Drawing the bin difference following a given covariance matrix using the Choleski decomposi-
 tion. This matrix is at least the statistical covariance matrix but can also contain systematic
 uncertainties.

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

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

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

$$(7.23)$$

2374 where V is covariance matrix used to produce the fluctuations, \mathbf{R} is drawn in a multinomial
 2375 distribution of mean 0 and variance 1, \mathbf{h} the bin content of the theoretical spectrum and $\tilde{\mathbf{h}}$ the
 2376 bin content of the generated toy.

2377 The first two methods allow for the fast production of independent toys while the third allow for
 2378 the production of statistical and systematical dependent toys. Unfortunately, none of those methods
 2379 are fitted to produce toy with a QNL different from the theoretical spectrum. The uncertainty on the
 2380 reconstructed energy σE_{rec} being dependent on E_{vis}/E_α makes that we would need to deconvolute
 2381 the reconstruction effect from the theoretical spectrum. It is much easier to just produce those toys
 2382 from the IBD generator.

2383 **7.4 Technical challenges and development**

2384 The fit framework Avenue was already partially developed with multispectra fitting in mind but
 2385 a lot technical development was necessary to allow for a joint fit. The first step was to migrate
 2386 the framework from ROOT5 (last release in March 2018) to ROOT6 (v6.26.06 released in July 2022)
 2387 to ensure compatibility with the data coming from the JUNO collaboration, and benefiting of the
 2388 improvement and corrections that came with ROOT6. This allow us to upgrade the C++ standard
 2389 from C++11 to C++17. A substantial effort has been done to modernize the code, generalizing the
 2390 functions and methods via templating to help readability and using smart pointer to prevent possible
 2391 memory leaks.

2392 The Avenue framework had to be adapted, notably on the chi-square calculation and spectrum gen-
 2393 eration to correctly take into account the correlation between the SPMT and LPMT spectra. The delta
 2394 joint fit requiring two more parameters over a spectrum twice as large as before with LPMT takes
 2395 much more time, around 15h for 6 years exposure, than the single LPMT fit. Thus the framework
 2396 and the fit macro had to be updated for distributed computing. Notably the aggregation of fit results
 2397 can now be done in a single file instead of managing a file per fit. In case of numerous toy, the hard
 2398 drive access time could lead to long analysis time.

2399 While the IBD generator was already able to generate LPMT and SPMT spectrum, it was not designed
 2400 for generating correlated spectrum. As detailed in section 7.3.1, up to the reconstruction effect, the

2401 two spectrum need to share the same generation else the two spectrum would be decorrelated and it
2402 would be like we would run two different experiment.

2403 7.5 Results

2404 7.5.1 Validation

2405 The first step is to confirm that the updated fit framework is able to reproduce existing results and
2406 that the joint fit behave as expected, meaning

- 2407 — Without QNL, the individual (*LPMT* and *SPMT*) fit converge to the parameters nominal
2408 values and their errors are similar to the ones reported in existing analysis such as [11].
- 2409 — The standard joint fit with an independent covariance matrix (*Indep Standard joint*), meaning
2410 that the covariance between the LPMT and SPMT spectra is 0, believe to have twice as much
2411 informations, and thus believe to have a greater precision than the individual fits.
- 2412 — The standard joint (*Standard joint*) fit with a correlated covariance matrix has errors similar to
2413 the LPMT individual fit as the LPMT drive the precision on θ_{13} and Δm_{31}^2 and that the LPMT
2414 as SPMT are expected to have close precision on θ_{12} and Δm_{21}^2 .
- 2415 — The delta joint (*Delta joint*) fit with covariance matrix have the same resolution as the standard
2416 joint fit. The supplementary parameter $\delta\theta_{12}$ and $\delta\Delta m_{21}^2$ should not bring supplementary
2417 precision.

2418 The italicized name are the name used in the results reports to identify each fit. We also look into the
2419 *Indep Delta joint*, which is the Delta Joint fit but the covariance between the LPMT and SPMT spectra
2420 is 0, and the *Weighted* results where

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

2421 We expect the weighted resolution to be similar to the *Indep Standard joint* as, in both of those test, we
2422 do not consider the correlation between the SPMT and LPMT results.

2423 Asimov studies

2424 We ran Asimov studies on the tests presented above on the updated framework, the results are
2425 reported in table 7.2. All those test are ran considering statistics error only, 6 years exposure with
2426 all backgrounds, Pearson χ^2 (covariance is estimated using data spectrum) and θ_{13} fixed to nominal
2427 value. For the *SPMT* fit Δm_{31}^2 is fixed at nominal value as the SPMT system is not expected to be
2428 sensitive to this parameter.

2429 In every cases presented above, the fit converges to the parameters nominal value thus only the
2430 errors are presented.

2431 We observe, as expected, that $\sigma_{\text{Weighted}} \approx \sigma_{\text{Indep Standard joint}}$ with the exception of $\sigma\theta_{12}$. This could
2432 from the slight difference in statistic between the SPMT and LPMT spectra. Indeed, due to a larger
2433 smearing in energy resolution, events that would be inside the spectrum range [0.8, 7.5] MeV are
2434 smeared outside it. This deficit is partially compensated by event outside the spectrum coming back
2435 in it but we expect very few event outside the spectrum in comparison to event at the edges of it.
2436 Thus the event deficit is not totally compensated. θ_{12} being mainly driven by the amplitude of the
2437 spectrum (see illustration 2.2), that's why we think this the origin of the difference.

2438 The second observation is that $\sigma_{\text{Standard joint}} \approx \sigma_{\text{LPMT}}$. Once the covariance matrix between the
2439 LPMT and SPMT is correctly introduced, the fit “understand” that it does not have supplementary
2440 information and the LPMT system, which have the best precision, dominate the resolution.

	Δm_{21}^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_{21}^2 and Δm_{31}^2						
Indep Standard joint			9.33082e-04			4.82955e-26
LPMT			1.27032e-03			2.58849e-26
SPMT			1.31070e-03			2.24106e-26
Weighted			9.12193e-04			
Fixed 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 7.2 – 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

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

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

Toy studies

Once we validated that the asimov study is yielding coherent results, we study the behaviour of toy studies. The above asimov study was using the Pearson χ^2 (Eq. 7.13) without pull parameter. We show in figure 7.6 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_{21}^2 biased of about 0.1σ . When introducing the PearsonV χ^2 (Eq. 7.17) the bias disappear as reported in figure 7.7.

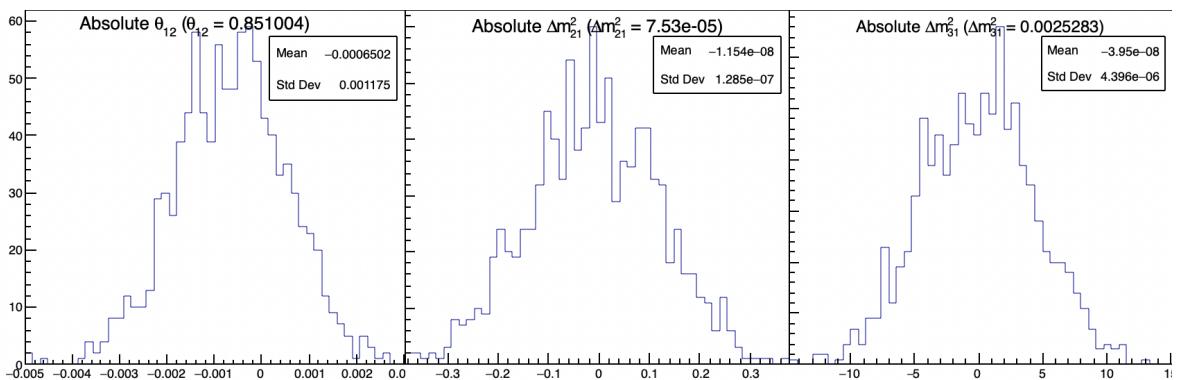


FIGURE 7.6 – Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, Pearson χ^2 , θ_{13} fixed.

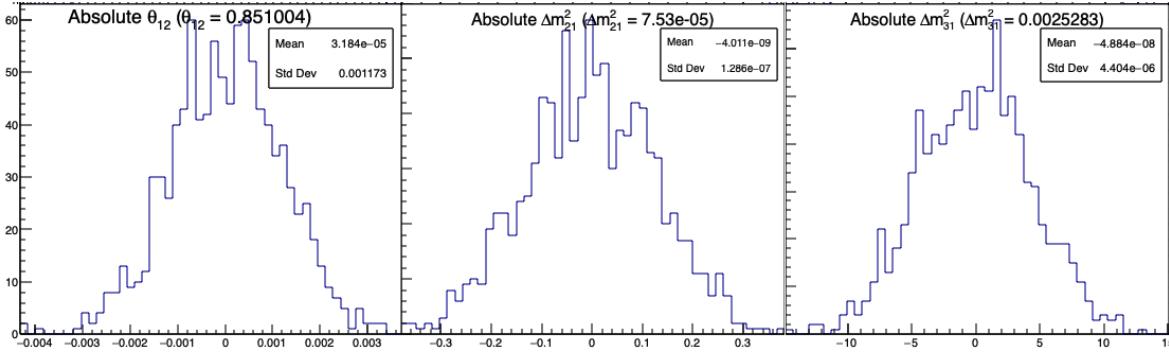


FIGURE 7.7 – Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, PearsonV χ^2 , θ_{13} fixed.

When the supplementary parameters are introduced in the Delta Joint fit, the fit is stable as shown in the results figure 7.8. The resolutions on the oscillation parameters are slightly worse in the Delta joint fit due to the supplementary freedom. As seen in the asimov studies, the resolution of the 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.

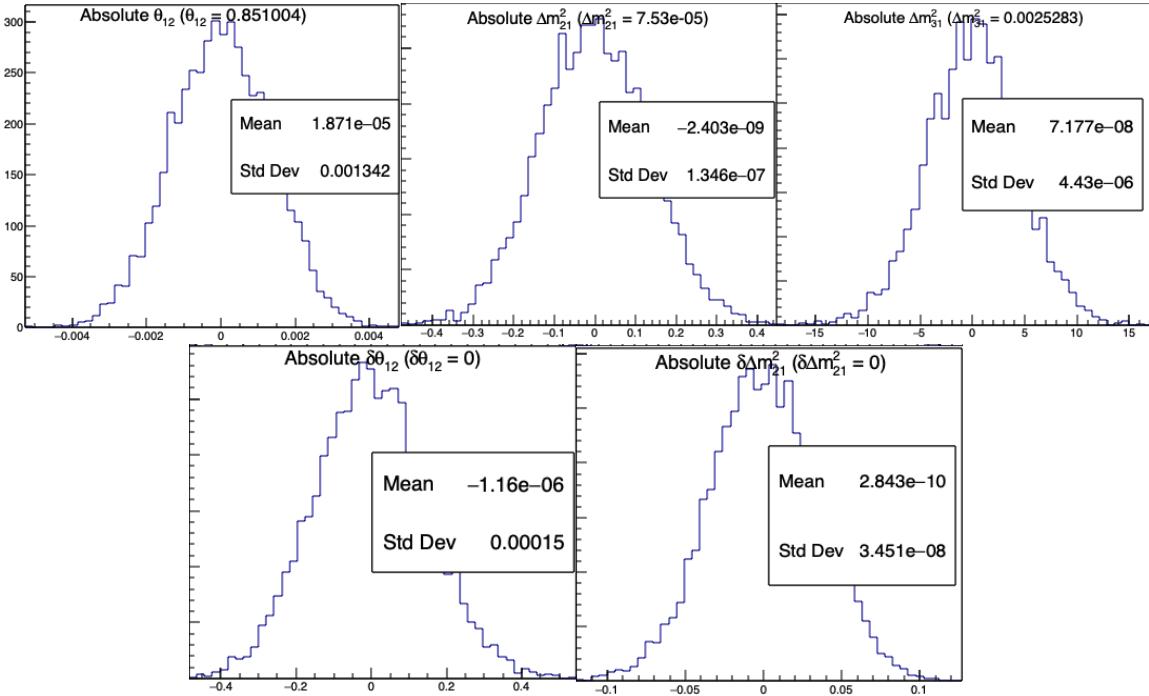


FIGURE 7.8 – Distribution of BFP - nominal value for 5000 toy Delta joint fit. 6 years exposure, all background, PearsonV χ^2 , θ_{13} fixed.

Effect of supplementary QNL on the LPMT spectrum

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

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

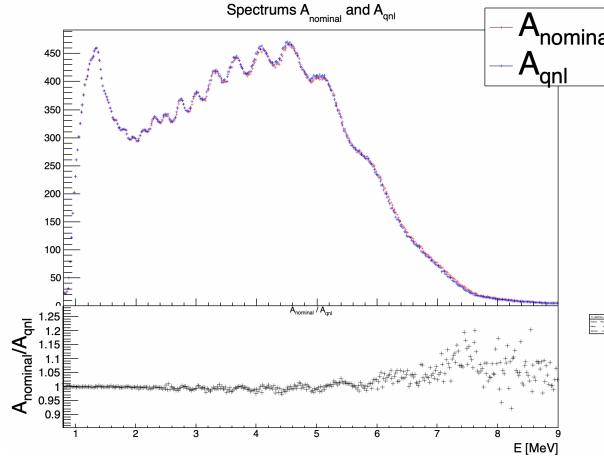


FIGURE 7.9 – **Top:** Theoretical spectrum without QNL (in red) and with $\alpha_{qnl} = 1\%$ (in blue). **Bottom:** Ratio between the theoretical spectrum with and without QNL.

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

TABLE 7.3 – Results of the different fit scenarios on QNL distorted data $\alpha_{qnl} = 1\%$. The mean value are reported subtracted from their nominal value. For SPMT 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.

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

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

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

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

2483 7.5.2 Covariance matrix

2484 The covariance matrix between the LPMT and SPMT spectra is at the heart of this study as it
 2485 was already mentioned in section 7.2 and demonstrated in section 7.5.1. In this section we discuss
 2486 the different approaches taken to estimate it. In this work we will mainly discuss the statistical
 2487 covariance matrix between the two spectra, how the number of event in a LPMT bin influence the
 2488 number of bin in the SPMT spectrum due to the resolution. We will still discuss the reconstruction
 2489 effects, mostly due to non-uniformity, in on reconstruction correlation.

2490 **Analytical method**

2491 The first method discussed is the analytical method where we propagate the resolution of the LPMT
 2492 and SPMT spectra over a non-smeared spectrum. Following the approach used in the IBD generation
 2493 in section 7.3.1, we consider the system resolution $\sigma(E)$ to be only dependent in energy. We do not
 2494 consider the position of the event.

2495 The first step is to compute the statistical uncertainty of the input spectrum while taking into account
 2496 the smearing, considering no uncertainty on the smearing. For this, using the notation of
 2497 section 39.2.5 *Propagation of errors* of PDG2020 [16] and considering an extended spectrum of 820 bins
 2498 following the binning scheme introduced in 7.2.4, the first 410 for the LPMT and the last 410, we
 2499 consider

2500 — $\theta = (\theta_0, \dots, \theta_n)$; $n = 820$ the content of the spectrum bins.

2501 — $\eta(\theta) = (\eta_0(\theta), \dots, \eta_m(\theta))$; $m = 820$ the set of smearing functions representing the PMT resolutions.

2503 η_m can thus be defined as

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

2504 where $G(i, \sigma(E_i))(j)$ is the smearing function defined as

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

2505 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
 2506 the i th bin respectively.

2507 We can then construct the transfer matrix A as

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

2508 and then compute the first part of our covariance matrix

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

2509 where V is the uncorrelated covariance matrix simply defined, under the assumption of poissonian
 2510 statistic for the bin content,

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

2511 Now we just need to consider the uncertainty on the smearing $\sigma \eta_i$, considering no uncertainty on
 2512 the unsmeared spectrum. From Eq. 7.25, the $G(i, j) \equiv G(i, \sigma(E_i))(j)$ are considered independents
 2513 from each other $\forall i, j$. This mean that this covariance matrix is diagonal, we only need $\sigma G(i, j)$. We
 2514 can derive this term from two equation:

- 2515 — The term $G(i, j)\theta_j$ represent the number of event smeared from the bin j that end up in the bin
 2516 i . This is a number, we thus assume poissonian statistic so that $\sigma[G(i, j)\theta_j] = \sqrt{G(i, j)\theta_j}$.
 2517 — Using basic error propagation we can say that $\sigma^2[G(i, j)\theta_j] = \theta_j^2\sigma^2G(i, j) + G(i, j)^2\sigma^2\theta_j$.

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

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

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

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

2518 By summing the two covariance matrix, we can extract a correlation matrix presented in figure 7.10.
 2519 The correlation between the SPMT and LPMT spectra is greater at the start of the spectrum, where
 2520 the absolute smearing is the smallest, up to 5% correlation, and diffuse as the bins are further from
 2521 each other and the absolute resolution grow.

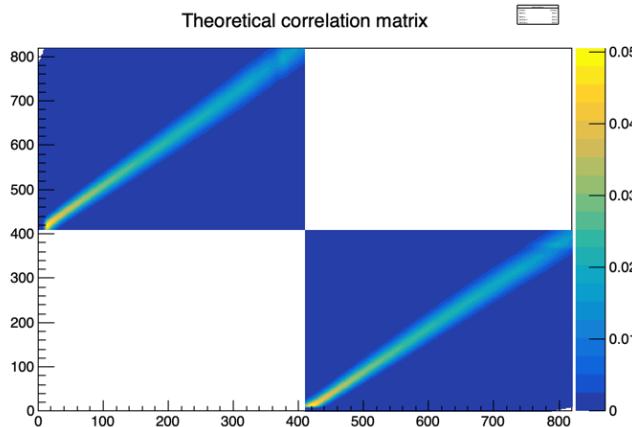


FIGURE 7.10 – Theoretical correlation matrix between the LPMT spectrum (bins 0-409) and the SPMT spectrum (410-819). The diagonal has been set to 0 (it was 1) for readability purpose.

2522 Empiric method

2523 The second method is the empiric way where we generate toys and just compute the empirical
 2524 correlation between the bin contents.

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

2525 We thus generate 10^7 event using the IBD generator presented in section 7.3.1, then produce spectra
 2526 from this finite set of events, meaning we must choose a number N of toy each composed of M event
 2527 in order to have the best estimate.

2528 Due to the nature of our estimator, the estimated correlation coefficient is subject to statistical fluctuation
 2529 as any estimator. There is no definite formula to compute the standard deviation of the correlation coefficient as suggested in this study [86] but all cited formula depend solely on the
 2530

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

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

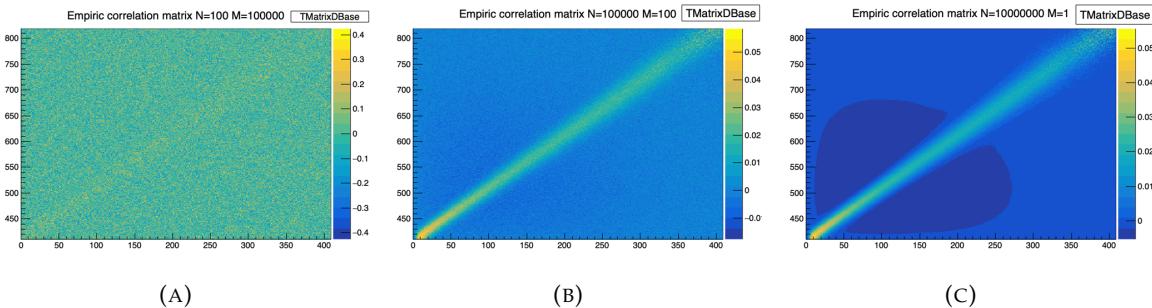


FIGURE 7.11 – Upper left corner of the estimated correlation matrix between the LPMT and SPMT spectrum for different configuration of N toy with different number of M events per toy

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

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

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

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

Empirical correlation matrix from fully simulated event

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

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

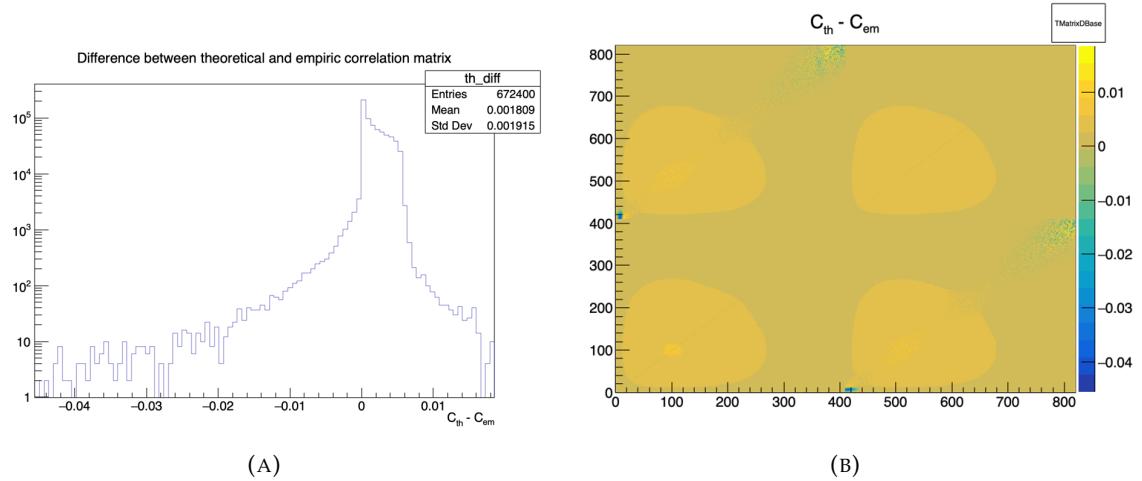


FIGURE 7.12 – Difference between the element of the theoretical and empiric correlation matrix

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

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

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

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

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

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

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

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

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

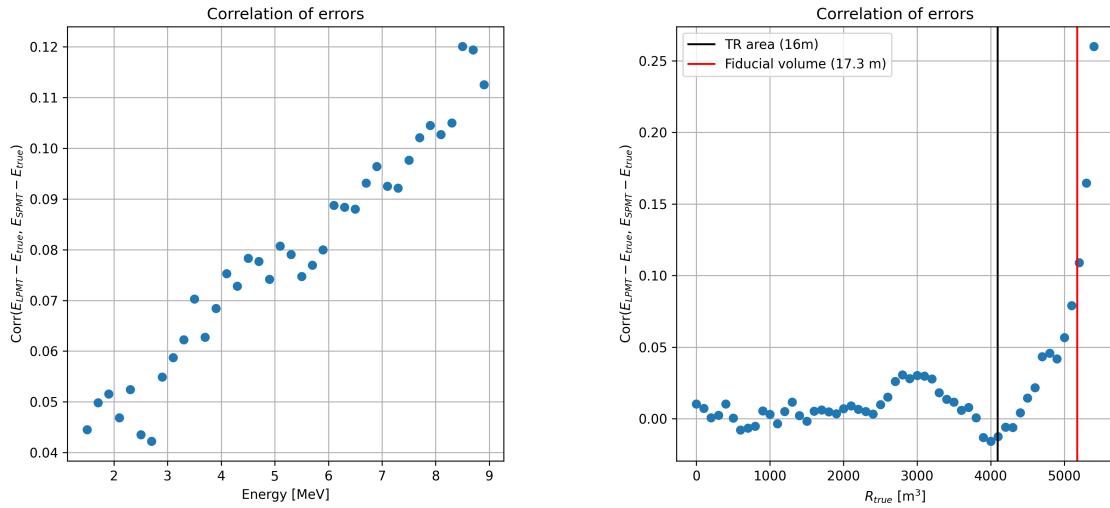


FIGURE 7.13 – Correlation on the reconstruction error between the LPMT and SPMT system as a function of (On the left) the energy, (On the right) the radius. The SPMT reconstruction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.6. To prevent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

2592 the mass ordering and parameters measurements are very sensitive to even small correlations.
 2593 We consider the binned approach that we used here, knowing that the CNN reconstruction was
 2594 deemed efficient but flawed, to be insufficient for the complete study of those effects on the fit.

2595 7.5.3 Statistical tests

2596 In this part, I present the results of the statistical tests presented in section 7.2.

2597 Test χ^2_{spe}

2598 The χ^2_{spe} is a chi-square representing the compatibility between the LPMT and SPMT spectra under
 2599 constraints of the correlation matrix between the two.

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

2600 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
 2601 calculation of V_{spe} , see section 7.2.

2602 The results for different exposures can be found in figure 7.15. To give an idea of the significance of
 2603 this test, we provide the median p-value for each test $\alpha_{qnl} \neq 0$. As expected, the power of this test
 2604 rises as the exposure does. We see significant discrimination at 6 years for $\alpha_{qnl} \geq 0.3\%$ where the
 2605 p-value for $\alpha_{qnl} = 3\%$ is 0.005 ± 0.0022 .

2606 This test relies solely on the estimated covariance matrix between the two spectra, requiring no
 2607 fitting. As a result, it is a very lightweight test that can still provide valuable indications of potential
 2608 unknown distortions between the two spectra.

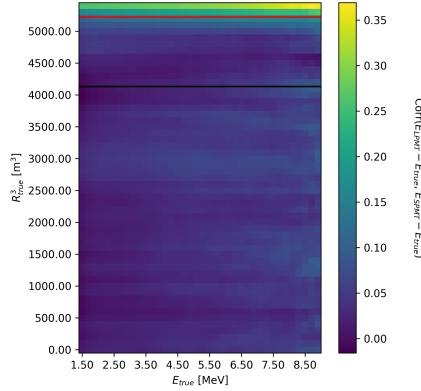


FIGURE 7.14 – Correlation on the reconstruction error between the LPMT and SPMT system as a function of the energy and the radius. The SPMT reconstruction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.6. To prevent effect due to the CNN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

2609 **Test χ_{ind}^2**

2610 The χ_{ind}^2 is the chi-square that represent the agreement between the measured oscillation parameters
 2611 θ_{12} and Δm_{21}^2 . This test is defined as

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

2612 where θ_{12}^L and $(\Delta m_{21}^2)^L$ are the oscillation parameters measured by the LPMT system. Same for θ_{12}^S
 2613 and $(\Delta m_{21}^2)^S$ for the SPMT system. We use V_{ind} computed for $\alpha_{qnl} = 0$. For more details about the
 2614 calculation of V_{ind} see section 7.2.

2615 The results are presented in figure 7.16. This test does not require any joint fit or covariance matrix
 2616 estimation between the two spectrum, it just need the estimated covariance matrix between the four
 2617 parameters. We see that the p-value are much less significant than the other tests, this is because this
 2618 test possess much less information about the relation between the LPMT and SPMT systems.

2619 This test is the most straightforward as it require only the fit of the two spectra and the estimation
 2620 of the parameters covariances, but is also the less powerful with a p value for $\alpha_{qnl} = 0.3\%$ of $0.09 \pm$
 2621 0.009 .

2622 **δ parameters significance**

2623 This test involves observing the values of the δ parameters in the Delta Joint fit and comparing them
 2624 tho their dispersion in the case where $\alpha_{qnl} = 0$. The results are shown in figures 7.17 and 7.18.

2625 We can see that the $\delta\Delta m_{21}^2$ has a very small discriminative power (figure 7.18) even at 6 years
 2626 exposure with a p-value of 0.34 ± 0.01 for $\alpha_{qnl} = 0.3\%$. On the other hand $\delta\theta_{12}$ (figure 7.17) has
 2627 much more discriminative power with a p-value for $\alpha_{qnl} = 0.3\%$ of 0.025 ± 0.005 . This test with a
 2628 single joint fit seems to be still less powerful than the χ_{spe}^2 . This can be explained as this method
 2629 only get information through the oscillation parameters θ_{12} and Δm_{21}^2 missing potential informations
 2630 contained in Δm_{31}^2 .

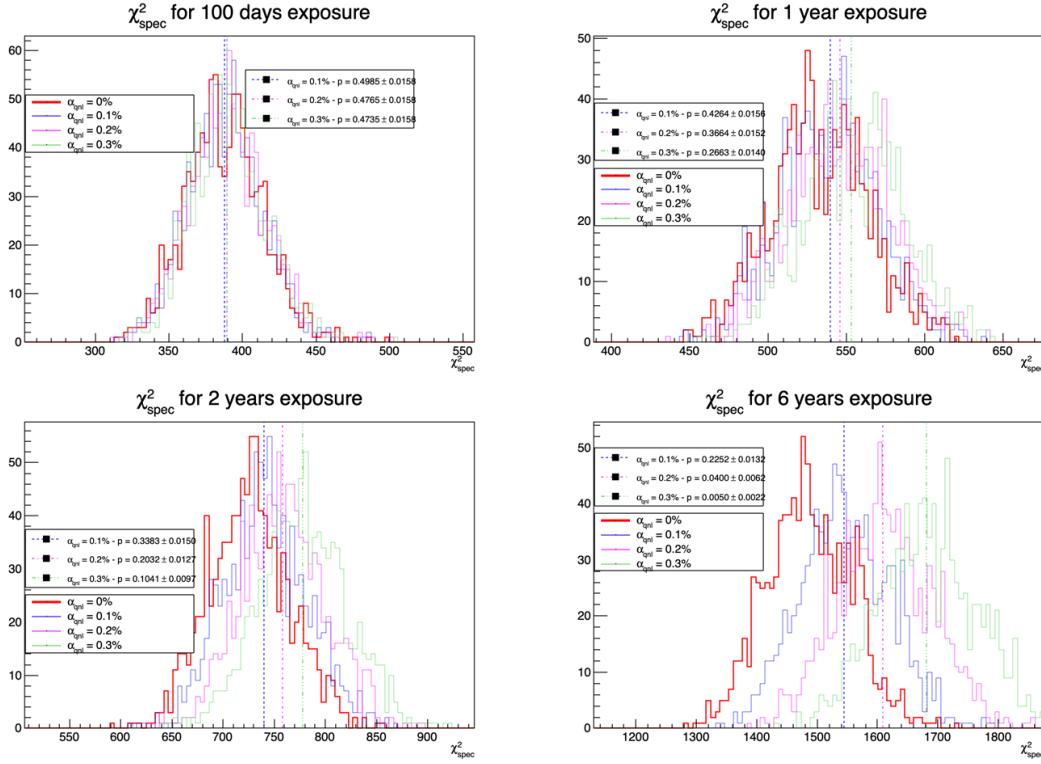


FIGURE 7.15 – Distribution of the χ^2_{spe} for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

2631 Hypothesis test

2632 In this last test we consider the two fit Standard Joint and Delta Joint as two hypothesis. The first
 2633 one, Standard Joint, is the H_0 hypothesis: we do not need supplementary parameters to describe the
 2634 energy spectrum. The second one, Delta Joint, is the H_1 hypothesis: we do need those supplementary
 2635 δ parameters to, if not correctly, approach the energy spectrum. If the δ parameter are unnecessary
 2636 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
 2637 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
 2638 exposures and α_{qnl} . The results are presented in figure 7.19.

2639 This test is the most complex, requiring two fit and the covariance matrix between the LPMT and
 2640 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
 2641 0.01 ± 0.003 .

2642 As explained in section 7.2.4, the spectra used for the fit are cut at 335 bins / 7.5 MeV to prevent
 2643 instability, while in χ^2_{spe} we use full 410 bins spectra. The χ^2_{spe} thus has more informations that the
 2644 hypothesis test leading to this difference in power.

2645 7.6 Conclusion and perspectives

2646 In this chapter, we present the development of a fit framework that allows us to fit multiple spectra
 2647 simultaneously. We also introduce a set of tools that enable us to detect potential distortions in one of
 2648 the two spectra. As an illustration of the capability of these tools, we use supplementary event-wise

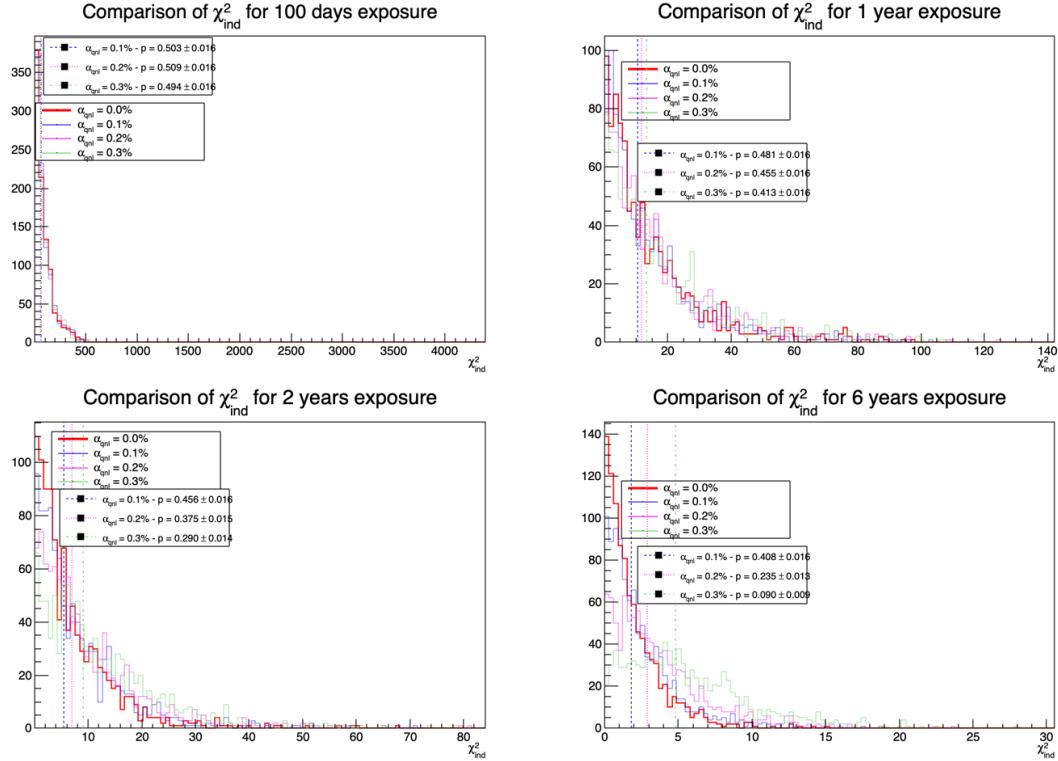


FIGURE 7.16 – Distribution of the χ^2_{Ind} for 1000 toys for different exposures. The dashed lines represent the median of the distributions and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

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.

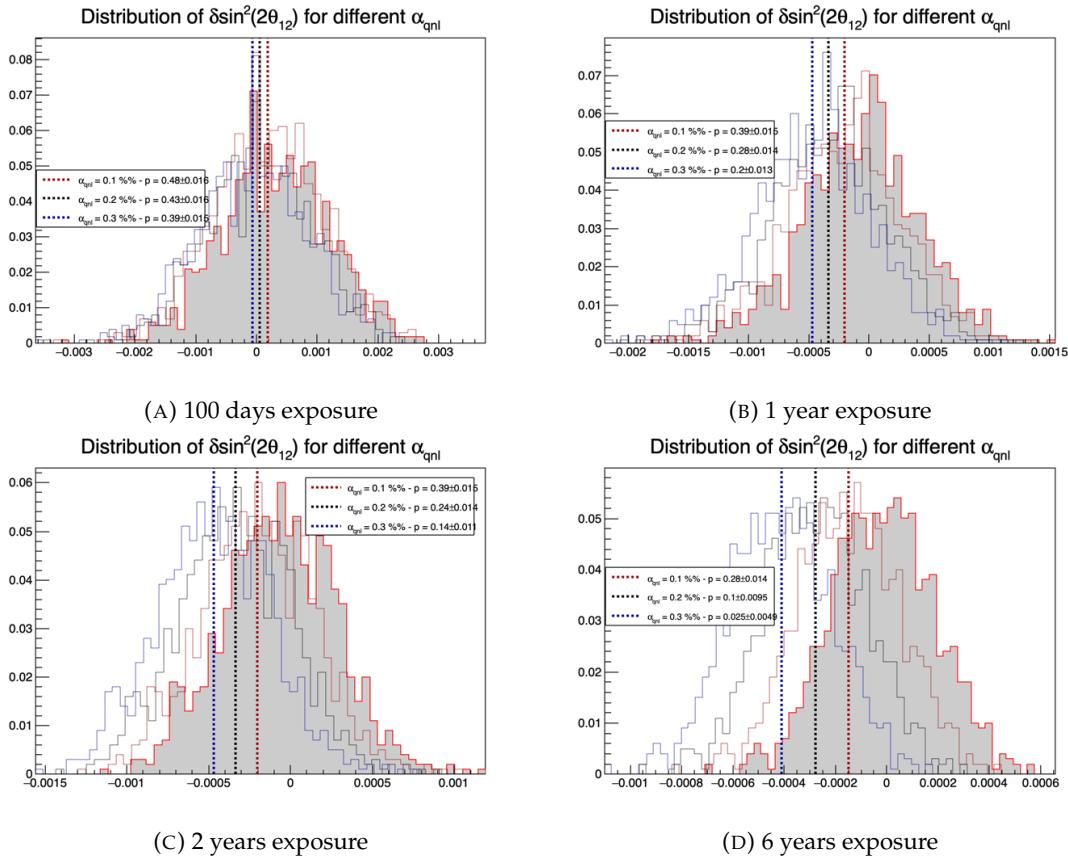


FIGURE 7.17 – Distribution of the $\delta \sin^2(2\theta_{12})$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

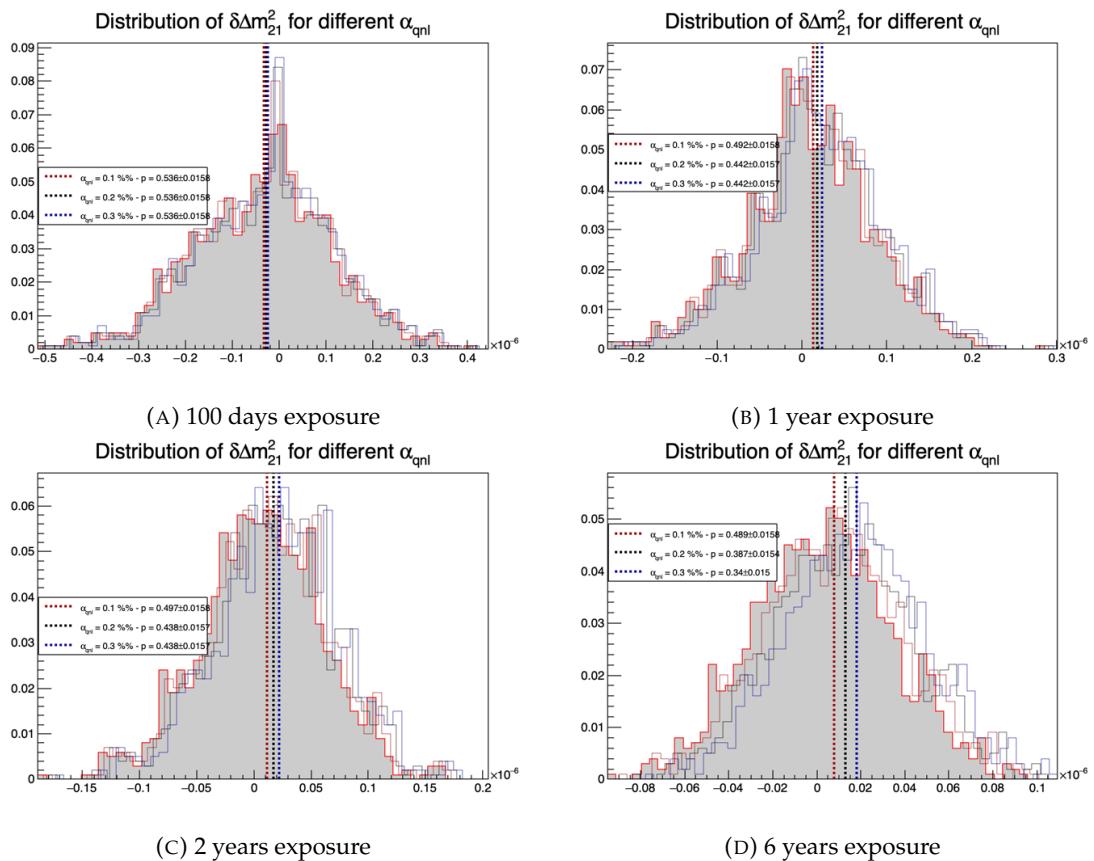


FIGURE 7.18 – Distribution of the $\delta\Delta m_{21}^2$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

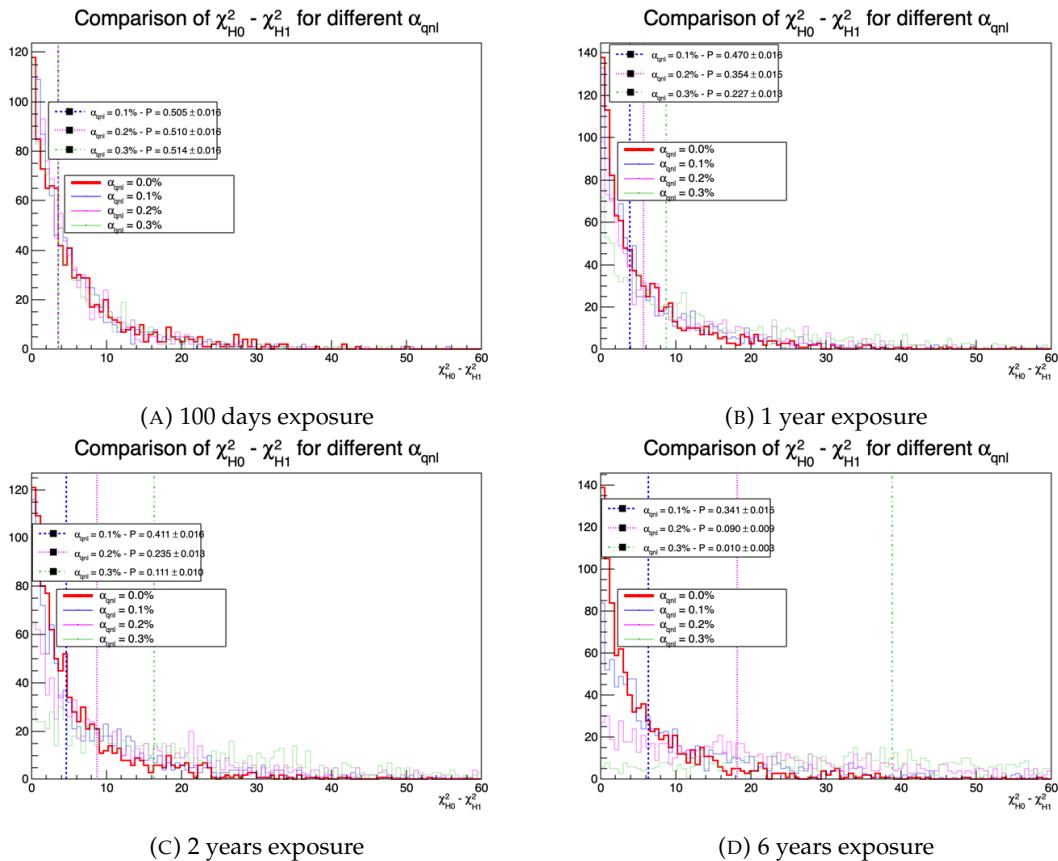


FIGURE 7.19 – Distribution of $\chi^2_{H_0} - \chi^2_{H_1}$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.

²⁶⁵⁸ **Chapter 8**

²⁶⁵⁹ **Conclusion**

²⁶⁶⁰ **Appendix A**

²⁶⁶¹ **Calculation of optimal α for estimator combination**

²⁶⁶³ This annex the details of the determination of the optimal α for estimator combination presented in
²⁶⁶⁴ section 4.3.2.

²⁶⁶⁵ As a reminder, the combined estimator $\hat{\theta}$ of X is defined as

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

²⁶⁶⁶ where θ_N and θ_C are both estimator of X .

²⁶⁶⁷ **A.1 Unbiased estimator**

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

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

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

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

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

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

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

$$(A.8)$$

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

²⁶⁶⁸ **A.2 Optimal variance estimator**

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

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

$$= \text{Var}[\alpha\theta_N] + \text{Var}[(1 - \alpha)\theta_C] + \text{Cov}[\alpha(1 - \alpha)\theta_N\theta_C] \quad (\text{A.11})$$

$$= \alpha^2\sigma_N^2 + (1 - \alpha)^2\sigma_C^2 + 2\alpha(1 - \alpha)\sigma_N\sigma_C\rho_{NC} \quad (\text{A.12})$$

²⁶⁶⁹ 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

$$\frac{d}{d\alpha} \text{Var}[\hat{\theta}] = 0 \quad (\text{A.15})$$

$$2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) - 2\sigma_C^2 + 2\sigma_N\sigma_C\rho_{NC} = 0 \quad (\text{A.16})$$

$$2\alpha(\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}) = 2\sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC} \quad (\text{A.17})$$

$$\alpha = \frac{\sigma_C^2 - \sigma_N\sigma_C\rho_{NC}}{\sigma_N^2 + \sigma_C^2 - 2\sigma_N\sigma_C\rho_{NC}} \quad (\text{A.18})$$

2670 This equation shows only one solution which is a minima. From Eq. A.18 arise two singularities:

- 2671 — $\sigma_N = \sigma_C = 0$. This is not a problem because as physicists we never measure with an absolute precision, neither us or our detectors are perfect.
- 2672 — $\sigma_N = \sigma_C$ and $\rho_{CN} = 1$. In this case θ_C and θ_N are the same estimator in term of variance thus any value for α yield the same result: an estimator with the same variance as the original ones.

2673

2674

²⁶⁷⁵ **Appendix B**

²⁶⁷⁶ **Charge spherical harmonics analysis**

²⁶⁷⁷ When looking at JUNO events we can clearly see some pattern in the charge repartition based on
²⁶⁷⁸ the event radius as illustrated in figure B.4. When dealing with identifying features and pattern on a
²⁶⁷⁹ spherical plane, the astrophysics community have been using, with success, the spherical harmonic
²⁶⁸⁰ decomposition. The principle is similar to a frequency analysis via Fourier transform. It comes to
²⁶⁸¹ saying that a function $f(r, \theta, \phi)$, here our charge repartition of the spherical plane constructed by our
²⁶⁸² PMTs, can be expressed

$$f(r, \theta, \phi) = \sum_{l=0}^{\infty} \sum_{m=-l}^l a_l^m r^l Y_l^m(\theta, \phi) \quad (\text{B.1})$$

²⁶⁸³ 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
²⁶⁸⁴ degree l and order m and P_l^m their associated Legendre Polynomials. Those harmonics are illustrated
²⁶⁸⁵ in figure B.1. By reducing the problem to the unit sphere $r = 1$, we get rid of the term r^l . The Healpix
²⁶⁸⁶ library [75] offer function to efficiently find the a_l^m factor from a given Healpix map.

²⁶⁸⁷ For the above decomposition, we will define the *Power* of an harmonic as

$$S_{ff}(l) = \frac{1}{2l+1} \sum_{m=-l}^l |a_l^m|^2 \quad (\text{B.2})$$

²⁶⁸⁸ and the *Relative Power* as:

$$P_l^h = \frac{S_{ff}(l)}{\sum_l S_{ff}(l)} \quad (\text{B.3})$$

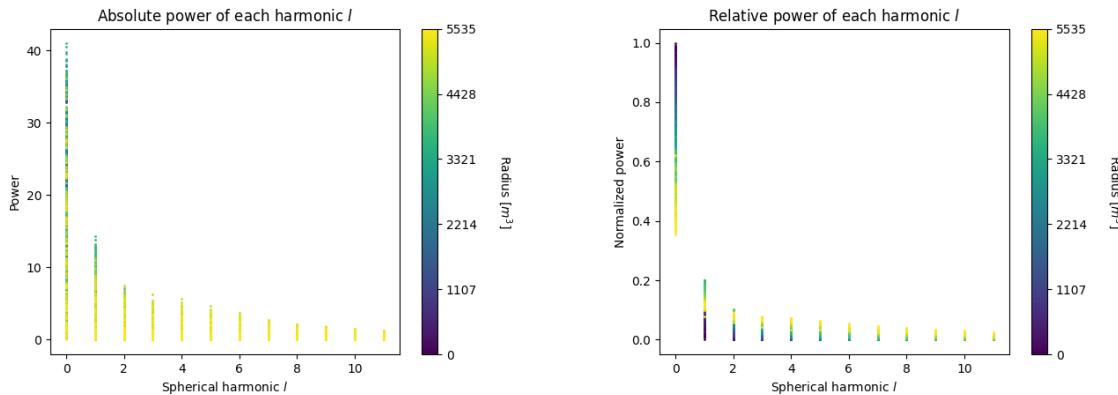
²⁶⁸⁹ For this study we will use 10k positron events with $E_{kin} \in [0; 9]$ MeV uniformly distributed in the
²⁶⁹⁰ CD from the JUNO official simulation version J23.0.1-rc8.dc1 (released the 7th January 2024). All the
²⁶⁹¹ event are *calib* level, with simulation of the physics, electronics, digitizations and triggers. We first
²⁶⁹² take a sub-set of 1k events and look at the power and relative power distribution depending on the
²⁶⁹³ radius and harmonic degree l . The results are shown in figure B.2. While don't see any pattern in
²⁶⁹⁴ absolute power, it is pretty clear that there is a correlation between the relative power of $l = 0$ and
²⁶⁹⁵ the radius of the event.

²⁶⁹⁶ When applying the same study but dependent on the energy, no clear correlation appear. The results
²⁶⁹⁷ for the $l = 0$ harmonic are presented in the figure B.5. Thus, in this study we will focus on the radial
²⁶⁹⁸ dependency of the relative power of each harmonic.

²⁶⁹⁹ In figures B.6 and B.7 are presented the distribution of the relative power of each harmonic for $l \in$
²⁷⁰⁰ $[0, 11]$. The relation between the radius and the relative power become even more clear, especially
²⁷⁰¹ for the first harmonics $l \in [0, 4]$. After that for $l > 4$ their relative power is close to 0 for central event,
²⁷⁰² thus loosing power. It also interesting to note the change of behavior in the TR area, clearly visible
²⁷⁰³ 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 $l = 0$.

2715 The minimization algorithm used are the Bent algorithm for $l = 0$ and the Bounded algorithm for
 2716 $l > 0$ provided by the Scipy library [88]. We then do the mean of the reconstructed radius from
 2717 the different harmonics. The reconstruction results are shown in figure B.3. The performance seems
 2718 correct but we see heavy fluctuation in the bias. To really be used as a reconstruction algorithm, the
 2719 method needs to be refined as discussed in the next section.

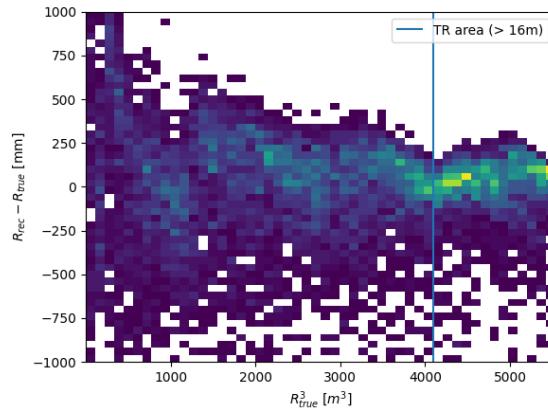


FIGURE B.3 – Error on the reconstructed radius vs the true radius by the harmonic method

Conclusion

2720 We have clearly shown in this analysis the relevance the of relative harmonic power for radius
 2721 reconstruction, and provided an erzats of a reconstruction algorithm. We will not delve further in
 2722 this thesis but if we wanted to refine this algorithm multiple paths can be explored:

- 2723 — No energy signature in the harmonics: This is surprising that there is no correlation between
 2724 the energy and the amplitude of the harmonics. We know that the energy is heavily correlated
 2725 with the total number of photoelectrons collected, it would be unintuitive that we see no
 2726 relation.
- 2727 — Localization of the event: We shown here the relation between the relative power of the har-
 2728 monic and the radius but don't get any information about the θ and ϕ spherical coordinates.
 2729 This information is probably hidden in the individual power of each order m of the degree l .
 2730 This intuition comes from the figure B.1 where in the higher degree l we see that the order m
 2731 are oriented. Intuitively, the order should be able to indicate a direction where the signal is
 2732 more powerful.
- 2733 — Combination of the degree power: Here we combined the radius reconstructed by the dif-
 2734 ferent degree via a simple mean but we shown in section 4.3.2 and annex A that this is note
 2735 the optimal way to combine estimator. A more refined algorithm probably exist to take into
 2736 account the predicting power of each order.

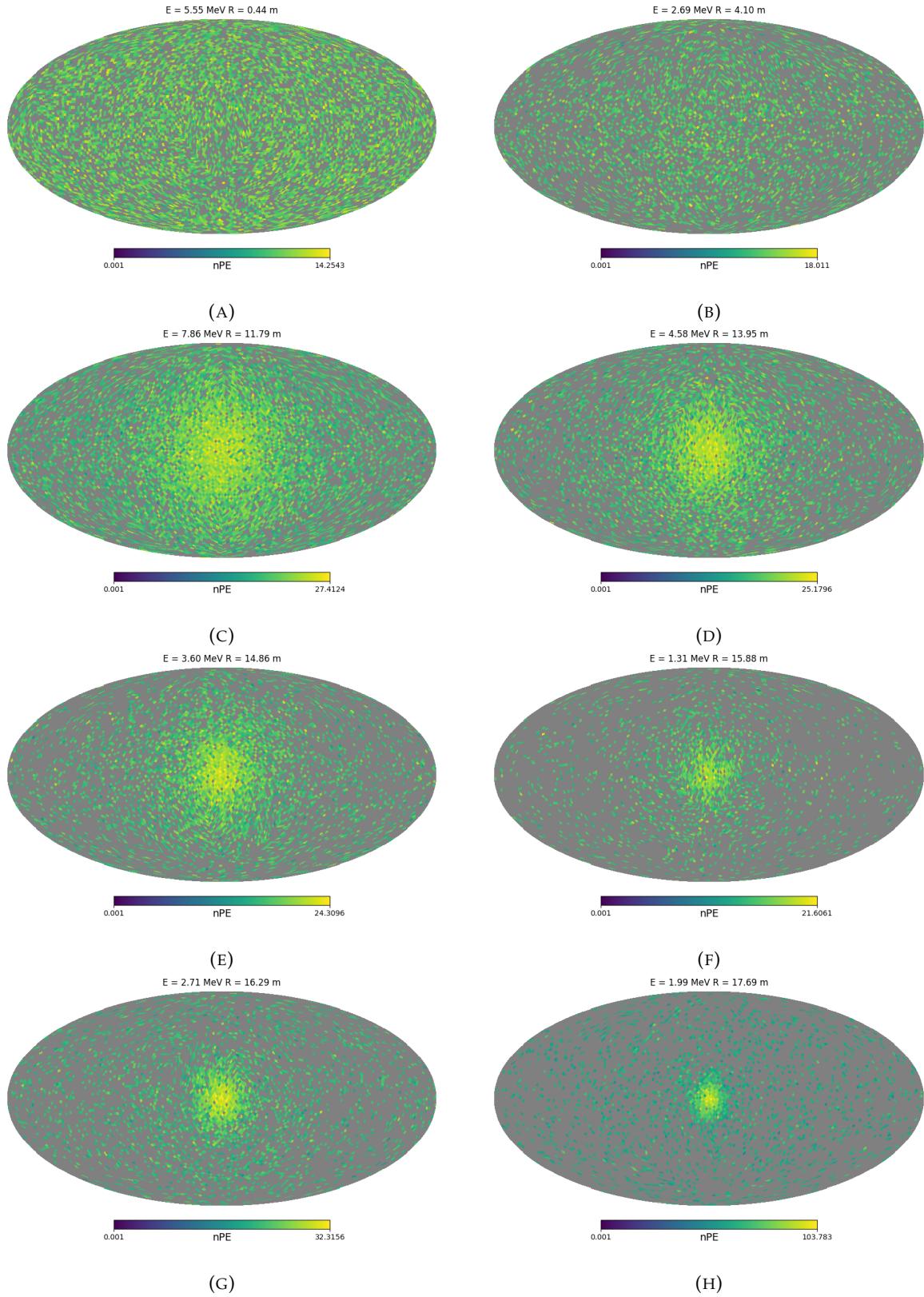


FIGURE B.4 – Charge repartition in JUNO as seen by the Healpix segmentation. Those are Healpix map of order 5 (i.e. 12288 pixels). The color represent the summed charge of the PMTs in each pixels. The color scale is logarithmic. The view have been centered to prevent event deformations.

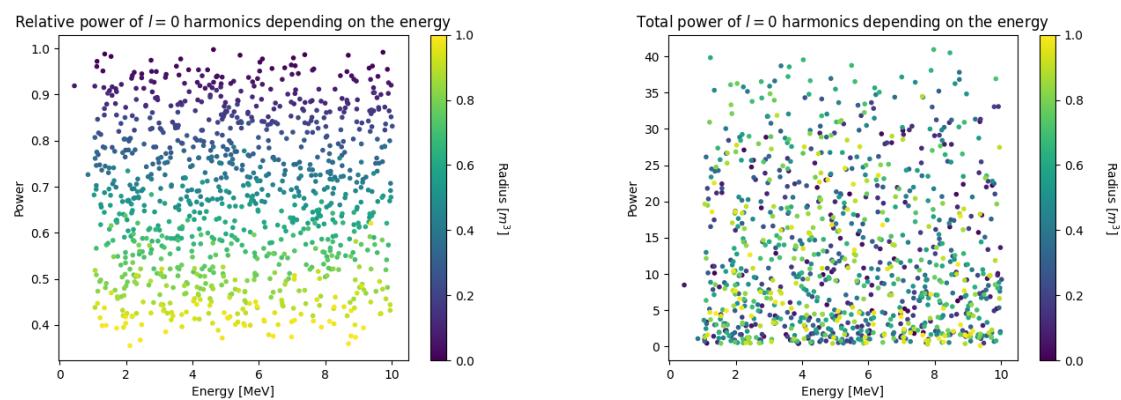


FIGURE B.5 – Scatter plot of the absolute and relative power, respectively on the left and right plot, of the $l = 0$ harmonic. The color indicate the radius of the event.

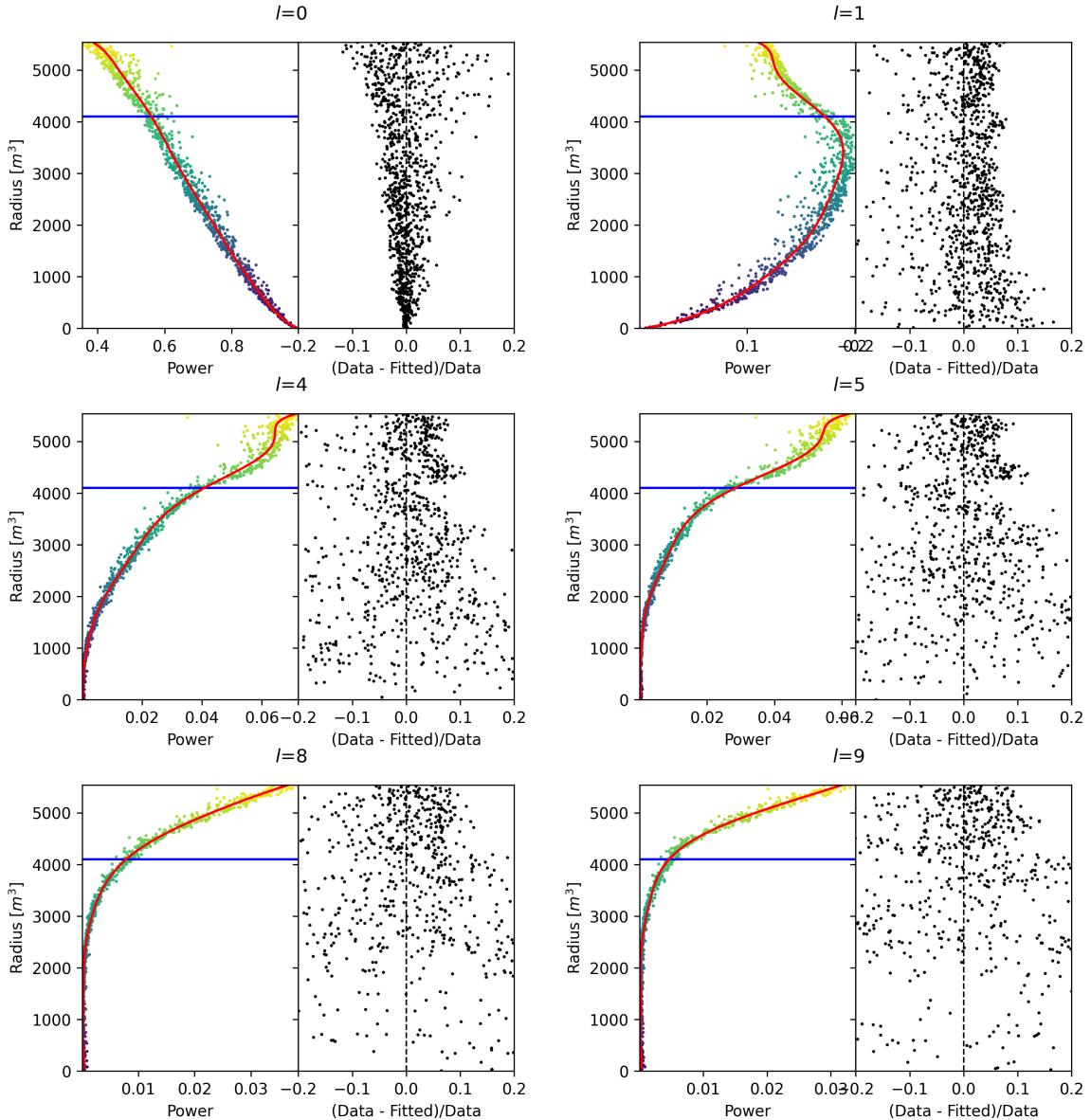


FIGURE B.6 – Plot of the distribution of the relative power of each harmonic dependent on R^3 (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. **Part 1**

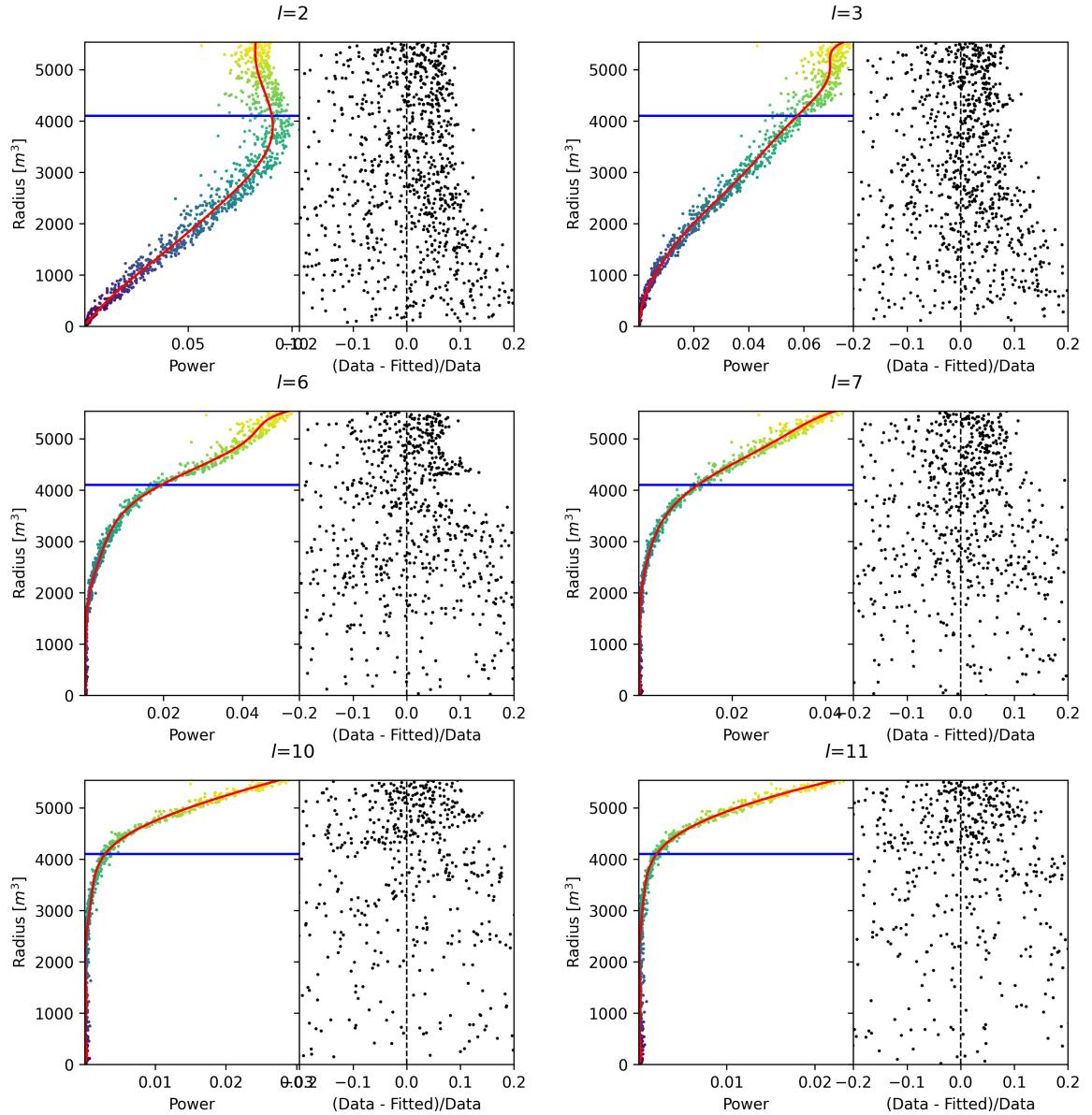


FIGURE B.7 – Plot of the distribution of the relative power of each harmonic dependent on R^3 (on the left). The Total Reflection (TR) area is represented by the horizontal blue line. The distribution are fitted using a 9th degree polynomial (red curve). The relative power error between the distribution and the fit is represented on the left. **Part 2**

²⁷³⁸ **Appendix C**

²⁷³⁹ **Additional spectrum smearing**

²⁷⁴⁰ In this section we demonstrate that a spectrum S smeared by a gaussian G parametrized by its
²⁷⁴¹ varianse σ_1^2 can be smeared by a gaussian parametrized by the variance σ_2^2 from the the smeared
²⁷⁴² 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

$$K'(E, \sigma_2) = S(E) \star G(E, \sigma_1) \star G^{-1}(E, \sigma_1) \star G(E, \sigma_2) \quad (\text{C.1})$$

$$= K(E, \sigma_1) \star G^{-1}(E, \sigma_1) \star G(E, \sigma_2) \quad (\text{C.2})$$

²⁷⁴³ 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

$$G(E, \sigma_1) \star G^{-1}(E, \sigma_1) = \delta(E) \quad (\text{C.3})$$

$$F[G(E, \sigma_1)](\nu) \times F[G^{-1}(E, \sigma_1)](\nu) = 1 \quad (\text{C.4})$$

²⁷⁴⁴ 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

$$F[G^{-1}(E, \sigma_1) \star G(E, \sigma_2)] = e^{\frac{\sigma_1^2(2\pi)^2}{2}\nu^2} \times e^{-\frac{\sigma_2^2(2\pi)^2}{2}\nu^2} \quad (\text{C.8})$$

$$= e^{\frac{(2\pi)^2}{2}(\sigma_1^2 - \sigma_2^2)\nu^2} \quad (\text{C.9})$$

$$= e^{\frac{(2\pi)^2}{2}\Delta\sigma^2\nu^2}; \Delta\sigma^2 = (\sigma_1^2 - \sigma_2^2) \quad (\text{C.10})$$

²⁷⁴⁵ 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

$$G^{-1}(E, \sigma_1) \star G(E, \sigma_2) = \frac{1}{\sqrt{|\Delta\sigma^2|}\sqrt{2\pi}} e^{-\frac{E^2}{2|\Delta\sigma^2|}} \quad (\text{C.11})$$

²⁷⁴⁶ **Appendix D**

²⁷⁴⁷ **Correction of E_{vis} bias**

²⁷⁴⁸ The reconstruction algorithms that are presented in this thesis in Chapters 4 and 5 do not reconstruct
²⁷⁴⁹ the same energy as the classical algorithms presented in section 2.6. Our algorithms reconstruct the
²⁷⁵⁰ deposited energy E_{dep} while the classical algorithms reconstruct a visible energy E_{vis} .

To understand this phenomena, let's look at the equation 2.15:

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

²⁷⁵¹ which define the expected N_{pe}/E . This define a linear relation between the number of photoelectrons
²⁷⁵² and the energy. However we discussed in sections 2.2.2 and 2.3 that the number of photoelectrons
²⁷⁵³ collected by the LPMT system do not follow a linear relationship. Thus this visible energy is not
²⁷⁵⁴ linear with the deposited energy. This effect is corrected in physics analysis and in Chapter 7 by
²⁷⁵⁵ applying the calibrated non-linearity profile the energy spectrum.

²⁷⁵⁶ When we need to compare our algorithm that reconstruct the deposited energy to the classical
²⁷⁵⁷ algorithms we need to correct this non-linearity. For this we fit the systematic bias of the classical
²⁷⁵⁸ algorithm using a 5th degree polynomial

$$\frac{E_{dep}}{E_{vis}} = \sum_{i=0}^5 P_i E_{dep}^i \quad (D.1)$$

²⁷⁵⁹ The fitted distribution and the corresponding fit is presented in figure D.1. The value fitted for this
²⁷⁶⁰ 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.

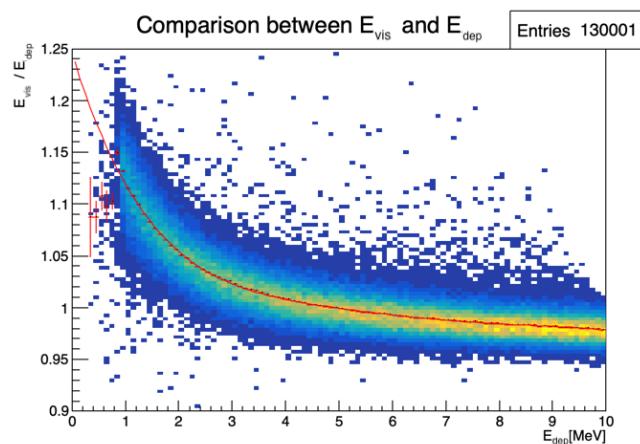


FIGURE D.1 – Comparison between Omilrec reconstructed E_{vis} and the deposited energy E_{dep} . The profile of the distribution E_{vis}/E_{dep} vs E_{dep} is fitted with a 5th degree polynomial.

List of Tables

2762	2.1	Characteristics of the nuclear power plants observed by JUNO.	14
2763	2.2	A summary of precision levels for the oscillation parameters. The reference value (PDG 2020 [16]) is compared with 100 days, 6 years and 20 years of JUNO data taking.	15
2764	2.3	Detectable neutrino signal in JUNO and the expected signal rates and major back- ground sources	16
2765	2.4	List of sources and their process considered for the energy scale calibration	24
2766	2.5	Calibration program of the JUNO experiment	26
2767	2.6	Features used by the BDT for vertex reconstruction	36
2768	2.7	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	37
2771	4.1	Sets of hyperparameters values considered in this study	58
2774	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. (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 <i>m</i> . (X_m, Y_m, Z_m) is 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 P_l^h is the relative power of the <i>l</i> th harmonic. See annex B for details.	76
2775	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 <i>m1</i> and the mesh <i>m2</i> . The features A and B are detailed in section 5.1	77
2779	7.1	Nominal PDG2020 value [16]. All value are reported assuming Normal Ordering.	101
2780	7.2	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 the error is estimated using the data spectrum	107
2781	7.3	Results of the different fit scenarios on QNL distorted data $\alpha_{qnl} = 1\%$. The mean value are reported subtracted from their nominal value. For SPMT 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.	109
2793	D.1	Parameters of the 5th degree polynomial used to correct Omilrec reconstructed energy.	135

List of Figures

2794	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
2796	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
2797		
2798		
2799		
2800		
2801		
2802		
2803		
2804		
2805		
2806		
2807		
2808		
2809		
2810		
2811		
2812		
2813		
2814		
2815		
2816		
2817		
2818		
2819		
2820		
2821		
2822		
2823		
2824		
2825		
2826		
2827		
2828		
2829		
2830		
2831		
2832		
2833		
2834		
2835		
2836		
2837		
2838		
2839		
2840		
2841		
2842		
2843		
2844		
2845		
2846		
2847		
2848		
2849		
2850		
2851		
2852		
2853		
2854		
2855		
2856		
2857		
2858		
2859		
2860		
2861		
2862		
2863		
2864		
2865		
2866		
2867		
2868		
2869		
2870		
2871		
2872		
2873		
2874		
2875		
2876		
2877		
2878		
2879		
2880		
2881		
2882		
2883		
2884		
2885		
2886		
2887		
2888		
2889		
2890		
2891		
2892		
2893		
2894		
2895		
2896		
2897		
2898		
2899		
2900		
2901		
2902		
2903		
2904		
2905		
2906		
2907		
2908		
2909		
2910		
2911		
2912		
2913		
2914		
2915		
2916		
2917		
2918		
2919		
2920		
2921		
2922		
2923		
2924		
2925		
2926		
2927		
2928		
2929		
2930		
2931		
2932		
2933		
2934		
2935		
2936		
2937		
2938		
2939		
2940		
2941		
2942		
2943		
2944		
2945		
2946		
2947		
2948		
2949		
2950		
2951		
2952		
2953		
2954		
2955		
2956		
2957		
2958		
2959		
2960		
2961		
2962		
2963		
2964		
2965		
2966		
2967		
2968		
2969		
2970		
2971		
2972		
2973		
2974		
2975		
2976		
2977		
2978		
2979		
2980		
2981		
2982		
2983		
2984		
2985		
2986		
2987		
2988		
2989		
2990		
2991		
2992		
2993		
2994		
2995		
2996		
2997		
2998		
2999		
3000		
3001		
3002		
3003		
3004		
3005		
3006		
3007		
3008		
3009		
3010		
3011		
3012		
3013		
3014		
3015		
3016		
3017		
3018		
3019		
3020		
3021		
3022		
3023		
3024		
3025		
3026		
3027		
3028		
3029		
3030		
3031		
3032		
3033		
3034		
3035		
3036		
3037		
3038		
3039		
3040		
3041		
3042		
3043		
3044		
3045		
3046		
3047		
3048		
3049		
3050		
3051		
3052		
3053		
3054		
3055		
3056		
3057		
3058		
3059		
3060		
3061		
3062		
3063		
3064		
3065		
3066		
3067		
3068		
3069		
3070		
3071		
3072		
3073		
3074		
3075		
3076		
3077		
3078		
3079		
3080		
3081		
3082		
3083		
3084		
3085		
3086		
3087		
3088		
3089		
3090		
3091		
3092		
3093		
3094		
3095		
3096		
3097		
3098		
3099		
3100		
3101		
3102		
3103		
3104		
3105		
3106		
3107		
3108		
3109		
3110		
3111		
3112		
3113		
3114		
3115		
3116		
3117		
3118		
3119		
3120		
3121		
3122		
3123		
3124		
3125		
3126		
3127		
3128		
3129		
3130		
3131		
3132		
3133		
3134		
3135		
3136		
3137		
3138		
3139		
3140		
3141		
3142		
3143		
3144		
3145		
3146		
3147		
3148		
3149		
3150		
3151		
3152		
3153		
3154		
3155		
3156		
3157		
3158		
3159		
3160		
3161		
3162		
3163		
3164		
3165		
3166		
3167		
3168		
3169		
3170		
3171		
3172		
3173		
3174		
3175		
3176		
3177		
3178		
3179		
3180		
3181		
3182		
3183		
3184		
3185		
3186		
3187		
3188		
3189		
3190		
3191		
3192		
3193		
3194		
3195		
3196		
3197		
3198		
3199		
3200		
3201		
3202		
3203		
3204		
3205		
3206		
3207		
3208		
3209		
3210		
3211		
3212		
3213		
3214		
3215		
3216		
3217		
3218		
3219		
3220		
3221		
3222		
3223		
3224		
3225		
3226		
3227		
3228		
3229		
3230		
3231		
3232		
3233		
3234		
3235		
3236		
3237		
3238		
3239		
3240		
3241		
3242		
3243		
3244		
3245		
3246		
3247		
3248		
3249		
3250		
3251		
3252		
3253		
3254		
3255		
3256		
3257		
3258		
3259		
3260		
3261		
3262		
3263		
3264		
3265		
3266		
3267		
3268		
3269		
3270		
3271		
3272		
3273		
3274		
3275		
3276		
3277		
3278		
3279		
3280		
3281		
3282		
3283		
3284		
3285		
3286		
3287		
3288		
328		

2838	a	Schematic of the TAO satellite detector	28
2839	b	Schematic of the OSIRIS satellite detector	28
2840	2.16	.	29
2841	a	Illustration of the different optical photons reflection scenarios. 1 is the reflection of the photon at the interface LS-acrylic or acrylic-water. 2 is the transmission of the photons through the interfaces. 3 is the conduction of the photon in the acrylic.	29
2842	b	Heatmap of R_{rec} and $R_{rec} - R_{true}$ as a function of R_{true} for 4MeV prompt signals uniformly distributed in the detector calculated by the charge based algorithm	29
2843	2.17	.	30
2844	a	Δt distribution at different iterations step j	30
2845	b	Heatmap of R_{rec} and $R_{rec} - R_{true}$ as a function of R_{true} for 4MeV prompt signals uniformly distributed in the detector calculated by the time based algorithm	30
2846	2.18	Bias of the reconstructed radius R (left), θ (middle) and ϕ (right) for multiple energies by the time likelihood algorithm	31
2847	2.19	On the left: Resolution of the reconstructed R as a function of the energy in the TR area ($R^3 > 4000\text{m}^3 \equiv R > 16\text{m}$) by the charge and time likelihood algorithms. On the right: Bias of the reconstructed R in the TR area for different energies by the charge likelihood algorithm	32
2848	2.20	Radial resolution of the different vertex reconstruction algorithms as a function of the energy	33
2849	2.21	.	33
2850	a	Spherical coordinate system used in JUNO for reconstruction	33
2851	b	Definition of the variables used in the energy reconstruction	33
2852	2.22	.	35
2853	a	Radial resolutions of the likelihood-based algorithm TMLE, QMLE and QTMLE	35
2854	b	Energy resolution of QMLE and QTMLE using different vertex resolutions	35
2855	2.23	Projection of the LPMTs in JUNO on a 2D plane. (a) Show the distribution of all PMTs and (b) and (c) are example of what the charge and time channel looks like respectively	37
2856	2.24	Radial (left) and energy (right) resolutions of different ML algorithms. The results presented here are from [42]. DNN is a deep neural network, BDT is a BDT, ResNet-J and VGG-J are CNN and GNN-J is a GNN.	38
2857	3.1	Example of a BDT that determine if the given object is a duck	42
2858	3.2	Schema of a simple neural network	43
2859	3.3	Illustration of the training lifecycle	45
2860	3.4	.	46
2861	a	Illustration of SGD falling into a local minima	46
2862	b	Illustration of the Adam momentum allowing it to overcome local minima	46
2863	3.5	Illustration of the SGD optimizer. In blue is the value of the loss function, orange, green and red are the path taken by the optimized parameter during the training for different LR.	47
2864	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.	47
2865	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.	47
2866	3.6	.	48
2867	a	Illustration of overtraining. The task at hand is to determine depending on two input variable x and y if the data belong to the dataset A or the dataset B . The expected boundary between the two dataset is represented in grey. A possible boundary learnt by overtraining is represented in brown.	48
2868	b	Illustration of a very simple NN	48

2890	3.7	Illustration of the ResNet framework	49
2891	3.8	Illustration of the gradient explosion. Here it can be solved with a lower learning rate but its not always the case.	49
2892			
2893	3.9	50
2894	a	Schema of a FCDNN	50
2895	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	50
2896			
2897	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].	51
2898			
2899	3.11	52
2900	a	Example of images in the MNIST dataset	52
2901	b	Schema of the CNN used in Pytorch example to process the MNIST dataset	52
2902			
2903	3.12	Illustration of a graph and its tensor representation.	53
2904	3.13	Illustration of the message passing algorithm. The detailed explanation can be found in section 3.2.3	53
2905			
2906			
2907	4.1	Graphic representation of the VGG-16 architecture, presenting the different kind of layer composing the architecture.	56
2908			
2909	4.2	61
2910	a	Spherical coordinate system used in JUNO for reconstruction	61
2911	b	Repartition of SPMTs in the image projection. The color scale is the number of SPMTs per pixel	61
2912			
2913	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.	61
2914			
2915			
2916			
2917	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.	62
2918			
2919			
2920			
2921	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.	62
2922			
2923			
2924			
2925	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.	63
2926			
2927			
2928			
2929	4.7	64
2930	a	Distribution of PE/MeV in the J23 Dataset. This distribution is profiled and fitted using equation 4.6	64
2931	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	64
2932			
2933	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.	65
2934			
2935			
2936			
2937	a	Resolution and bias of energy reconstruction vs energy	65
2938	b	Resolution and bias of energy reconstruction vs radius	65
2939	c	Resolution and bias of radius reconstruction vs energy	65
2940	d	Resolution and bias of radius reconstruction vs radius	65
2941	e	Resolution and bias of radius reconstruction vs θ	65

2942	f	Resolution and bias of radius reconstruction vs ϕ	65
2943	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 .	66
2944	a	Distribution of the error on reconstructed x by Gen ₃₀	66
2945	b	Distribution of the error on reconstructed y by Gen ₃₀	66
2946	c	Distribution of the error on reconstructed z by Gen ₃₀	66
2947	d	Distribution of the error on reconstructed R by Gen ₃₀	66
2948	e	Distribution of the error on reconstructed θ by Gen ₃₀	66
2949	f	Distribution of the error on reconstructed ϕ by Gen ₃₀	66
2950	4.10		67
2951	a	Distribution of Gen ₃₀ reconstructed energy and true energy of the analysis dataset (J21)	67
2952	b	Distribution of Gen ₄₂ reconstructed energy and true energy of the analysis dataset (J23)	67
2953	4.11	Radius bias (on the left) and resolution (on the right) of the classical algorithm in a E , R^3 grid	68
2954	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.	69
2955	a	Resolution and bias of energy reconstruction vs energy	69
2956	b	Resolution and bias of energy reconstruction vs radius	69
2957	c	Resolution and bias of radius reconstruction vs energy	69
2958	d	Resolution and bias of radius reconstruction vs radius	69
2959	e	Resolution and bias of radius reconstruction vs θ	69
2960	f	Resolution and bias of radius reconstruction vs ϕ	69
2961	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	70
2962	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.	71
2963	a	Resolution and bias of energy reconstruction vs energy	71
2964	b	Resolution and bias of energy reconstruction vs radius	71
2965	c	Resolution and bias of radius reconstruction vs energy	71
2966	d	Resolution and bias of radius reconstruction vs radius	71
2967	e	Resolution and bias of radius reconstruction vs θ	71
2968	f	Resolution and bias of radius reconstruction vs ϕ	71
2969	5.1		75
2970	a	Illustration of the different nodes in our graphs and their relations	75
2971	b	Illustration of what a dense adjacency matrix would looks like and the part we are really interested in. Because Fired \rightarrow Mesh and Mesh \rightarrow I/O relations are undirected, we only consider in practice the top right part of the matrix for those relations	75
2972	5.2	Illustration of the Healpix segmentation. On the left: A segmentation of order 0. On the right: A segmentation of order 1	75
2973	5.3	Illustration of the different update function needed by our GNN	78
2974	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	79
2975	a		79
2976	b		79
2977	c		79

2993	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 $E_{true} \in [0, 9]$	79
2994	a	79	
2995	b	79	
2996			
2997	5.6	Schema of the JWGv8.4.0 architecture, the colored triplet is the graph configuration after each JWG layers	81
2998			
2999	5.7	Energy reconstruction depending on the true energy for samples of the different versions of the GNN	82
3000			
3001	5.8	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.6, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.	84
3002	a	Resolution and bias of energy reconstruction vs energy	84
3003	b	Resolution and bias of energy reconstruction vs radius	84
3004			
3005	5.9	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.6, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.	85
3006	a	Resolution and bias of radius reconstruction vs energy	85
3007	b	Resolution and bias of radius reconstruction vs radius	85
3008			
3009	5.10	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.6, JWGv8.4 presented in this chapter. The top part of each plot is the resolution and the bottom part is the bias.	86
3010	a	Resolution and bias of radius reconstruction vs θ	86
3011	b	Resolution and bias of radius reconstruction vs ϕ	86
3012			
3013	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.	87
3014	a	Resolution and bias of energy reconstruction vs energy	87
3015	b	Resolution and bias of energy reconstruction vs radius	87
3016			
3017	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.	88
3018	a	Resolution and bias of radius reconstruction vs energy	88
3019	b	Resolution and bias of radius reconstruction vs radius	88
3020			
3021	5.13	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.6, 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.	89
3022	a	Resolution and bias of energy reconstruction vs energy	89
3023	b	Resolution and bias of energy reconstruction vs radius	89
3024			
3025	5.14	Reconstruction performance of the Omilrec algorithm based on QTML presented in section 2.6, 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.	89
3026	a	Resolution and bias of radius reconstruction vs energy	89
3027	b	Resolution and bias of radius reconstruction vs radius	89
3028			
3029	6.1	Schema of the method to discover vulnerabilities in the reconstruction methods	93
3030			
3031	7.1	Two oscillated spectra of $1e7$ event expected in JUNO. In red the spectrum without supplementary QNL. In blue the same spectrum but where an event-wise QNL $\alpha_{qnl} = 10\%$ is introduced.	97
3032			
3033	7.2	98
3034	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.	98
3035			
3036			
3037			
3038			
3039			
3040			
3041			
3042			
3043			
3044			

3045 b	Ratio of collected nPE after the additional QNL over the number of nPE that would be collected at different energies. We select event with an interaction radius $R < 4\text{m}$ to not be affected by the non-uniformity. The dots represent the mean of the distributions in figure 7.2a and the dashed line are the equivalent event-wise non-linearity from eq 7.2. The hatched zone is the residual non- linearity expected after calibration [29].	98
3051 7.3	Theoretical LPMT spectrum at nominal oscillation values binned using 410 bins from 0.8 to 9 MeV. It is rescaled to 6 years statistic. The black line represent the 335 bin cut .	102
3053 7.4	Schematic description of the fit framework	103
3054 7.5	Relative (On the left) and absolute (On the right) resolutions of the LPMT and SPMT systems used in this study. The number in parenthesis are the parameter A , B and C respectively for each systems.	104
3056 7.6	Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, Pearson χ^2 , θ_{13} fixed.	107
3059 7.7	Distribution of BFP - nominal value for 1000 toy Standard joint fit. 6 years exposure, all background, PearsonV χ^2 , θ_{13} fixed.	108
3061 7.8	Distribution of BFP - nominal value for 5000 toy Delta joint fit. 6 years exposure, all background, PearsonV χ^2 , θ_{13} fixed.	108
3063 7.9	Top: Theoretical spectrum without QNL (in red) and with $\alpha_{qnl} = 1\%$ (in blue). Bottom: Ratio between the theoretical spectrum with and without QNL.	109
3064 7.10	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. .	111
3067 7.11	Upper left corner of the estimated correlation matrix between the LPMT and SPMT spectrum for different configuration of N toy with different number of M events per toy	112
3069 a	112
3070 b	112
3071 c	112
3072 7.12	Difference between the element of the theoretical and empiric correlation matrix . .	113
3073 a	113
3074 b	113
3075 7.13	Correlation on the reconstruction error between the LPMT and SPMT system as a function of (On the left) the energy, (On the right) the radius. The SPMT recon- struction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.6. To prevent effect due to the CNN bad reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.	114
3080 7.14	Correlation on the reconstruction error between the LPMT and SPMT system as a function of the energy and the radius. The SPMT reconstruction comes from the NN presented in Chapter 4 and the LPMT reconstruction comes from OMILREC presented in section 2.6. To prevent effect due to the CNN bad reconstruction, we select the event with $1 < E_{dep} < 9$ MeV.	115
3085 7.15	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.	116
3088 7.16	Distribution of the χ^2_{ind} for 1000 toys for different exposures. The dashed lines repre- sent the median of the distributions and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.	117
3091 7.17	Distribution of the $\delta \sin^2(2\theta_{12})$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.	118
3094 a	100 days exposure	118
3095 b	1 year exposure	118
3096 c	2 years exposure	118
3097 d	6 years exposure	118

3098	7.18	Distribution of the $\delta\Delta m_{21}^2$ for 1000 toys for different exposure. The dashed line represent the median of the distribution and the p-value are the percentage of the $\alpha_{qnl} = 0$ distribution that are greater than those medians.	119
3099	a	100 days exposure	119
3100	b	1 year exposure	119
3101	c	2 years exposure	119
3102	d	6 years exposure	119
3105	7.19	Distribution of $\chi_{H_0}^2 - \chi_{H_1}^2$ 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.	120
3106	a	100 days exposure	120
3107	b	1 year exposure	120
3108	c	2 years exposure	120
3109	d	6 years exposure	120
3111			
3112	B.1	Illustration of the real part of the spherical harmonics	126
3113	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.	126
3114	B.3	Error on the reconstructed radius vs the true radius by the harmonic method	127
3115	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.	128
3119	a	128
3120	b	128
3121	c	128
3122	d	128
3123	e	128
3124	f	128
3125	g	128
3126	h	128
3127			
3128	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.	129
3129	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	130
3130	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	131
3138	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.	136
3139			

³¹⁴⁰ 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> (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> (visited on 07/28/2023).
- [3] Liang Zhan, Yifang Wang, Jun Cao, and Liangjian Wen. "Experimental Requirements to Determine the Neutrino Mass Hierarchy Using Reactor Neutrinos". *Physical Review D* 79.7 (Apr. 14, 2009), 073007. ISSN: 1550-7998, 1550-2368. DOI: [10.1103/PhysRevD.79.073007](https://doi.org/10.1103/PhysRevD.79.073007). eprint: [0901.2976\[hep-ex\]](https://arxiv.org/abs/0901.2976). URL: <http://arxiv.org/abs/0901.2976> (visited on 09/18/2023).
- [4] A. A. Hahn, K. Schreckenbach, W. Gelletly, F. von Feilitzsch, G. Colvin, and B. Krusche. "Antineutrino spectra from 241Pu and 239Pu thermal neutron fission products". *Physics Letters B* 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> (visited on 01/16/2024).
- [6] F. von Feilitzsch, A. A. Hahn, and K. Schreckenbach. "Experimental beta-spectra from 239Pu and 235U 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 235U 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> (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> (visited on 01/16/2024).

- [11] JUNO Collaboration et al. "Sub-percent Precision Measurement of Neutrino Oscillation Parameters with JUNO". *Chinese Physics C* 46.12 (Dec. 1, 2022), 123001. ISSN: 1674-1137, 2058-6132. DOI: [10.1088/1674-1137/ac8bc9](https://doi.org/10.1088/1674-1137/ac8bc9). eprint: [2204.13249\[hep-ex\]](https://arxiv.org/abs/2204.13249). URL: <http://arxiv.org/abs/2204.13249> (visited on 08/11/2023).
- [12] JUNO Collaboration et al. *TAO Conceptual Design Report: A Precision Measurement of the Reactor Antineutrino Spectrum with Sub-percent Energy Resolution*. May 18, 2020. DOI: [10.48550/arXiv.2005.08745](https://doi.org/10.48550/arXiv.2005.08745). eprint: [2005.08745\[hep-ex, physics:nucl-ex, physics:physics\]](https://arxiv.org/abs/2005.08745). URL: <http://arxiv.org/abs/2005.08745> (visited on 01/18/2024).
- [13] G. Mention, M. Fechner, Th. Lasserre, Th. A. Mueller, D. Lhuillier, M. Cribier, and A. Letourneau. "Reactor antineutrino anomaly". *Physical Review D* 83.7 (Apr. 29, 2011). Publisher: American Physical Society, 073006. DOI: [10.1103/PhysRevD.83.073006](https://doi.org/10.1103/PhysRevD.83.073006). URL: <https://link.aps.org/doi/10.1103/PhysRevD.83.073006> (visited on 03/05/2024).
- [14] V. Kopeikin, M. Skorokhvatov, and O. Titov. "Reevaluating reactor antineutrino spectra with new measurements of the ratio between ^{235}U and ^{239}Pu β^- spectra". *Physical Review D* 104.7 (Oct. 25, 2021), L071301. ISSN: 2470-0010, 2470-0029. DOI: [10.1103/PhysRevD.104.L071301](https://doi.org/10.1103/PhysRevD.104.L071301). eprint: [2103.01684\[hep-ph, physics:nucl-ex, physics:nucl-th\]](https://arxiv.org/abs/2103.01684). URL: <http://arxiv.org/abs/2103.01684> (visited on 01/18/2024).
- [15] A. Letourneau et al. "On the origin of the reactor antineutrino anomalies in light of a new summation model with parameterized β^- transitions". *Physical Review Letters* 130.2 (Jan. 10, 2023), 021801. ISSN: 0031-9007, 1079-7114. DOI: [10.1103/PhysRevLett.130.021801](https://doi.org/10.1103/PhysRevLett.130.021801). eprint: [2205.14954\[hep-ex, physics:hep-ph\]](https://arxiv.org/abs/2205.14954). URL: <http://arxiv.org/abs/2205.14954> (visited on 01/16/2024).
- [16] Particle Data Group et al. "Review of Particle Physics". *Progress of Theoretical and Experimental Physics* 2020.8 (Aug. 14, 2020), 083C01. ISSN: 2050-3911. DOI: [10.1093/ptep/ptaa104](https://doi.org/10.1093/ptep/ptaa104). URL: <https://doi.org/10.1093/ptep/ptaa104> (visited on 12/04/2023).
- [17] Super-Kamiokande Collaboration et al. "Diffuse Supernova Neutrino Background Search at Super-Kamiokande". *Physical Review D* 104.12 (Dec. 10, 2021), 122002. ISSN: 2470-0010, 2470-0029. DOI: [10.1103/PhysRevD.104.122002](https://doi.org/10.1103/PhysRevD.104.122002). eprint: [2109.11174\[astro-ph, physics:hep-ex\]](https://arxiv.org/abs/2109.11174). URL: <http://arxiv.org/abs/2109.11174> (visited on 02/28/2024).
- [18] JUNO Collaboration et al. "JUNO Sensitivity on Proton Decay $p \rightarrow \bar{\nu}K^+$ Searches". *Chinese Physics C* 47.11 (Nov. 1, 2023), 113002. ISSN: 1674-1137, 2058-6132. DOI: [10.1088/1674-1137/ace9c6](https://doi.org/10.1088/1674-1137/ace9c6). eprint: [2212.08502\[hep-ex, physics:hep-ph\]](https://arxiv.org/abs/2212.08502). URL: <http://arxiv.org/abs/2212.08502> (visited on 08/09/2024).
- [19] Alessandro Strumia and Francesco Vissani. "Precise quasielastic neutrino/nucleon cross section". *Physics Letters B* 564.1 (July 2003), 42–54. ISSN: 03702693. DOI: [10.1016/S0370-2693\(03\)00616-6](https://doi.org/10.1016/S0370-2693(03)00616-6). eprint: [astro-ph/0302055](https://arxiv.org/abs/astro-ph/0302055). URL: <http://arxiv.org/abs/astro-ph/0302055> (visited on 01/16/2024).
- [20] Daya Bay et al. *Optimization of the JUNO liquid scintillator composition using a Daya Bay antineutrino detector*. July 1, 2020. DOI: [10.48550/arXiv.2007.00314](https://doi.org/10.48550/arXiv.2007.00314). eprint: [2007.00314\[hep-ex, physics:physics\]](https://arxiv.org/abs/2007.00314). URL: <http://arxiv.org/abs/2007.00314> (visited on 07/26/2023).
- [21] J. B. Birks. "CHAPTER 3 - THE SCINTILLATION PROCESS IN ORGANIC MATERIALS—I". *The Theory and Practice of Scintillation Counting*. Ed. by J. B. Birks. International Series of Monographs in Electronics and Instrumentation. Jan. 1, 1964, 39–67. ISBN: 978-0-08-010472-0. DOI: [10.1016/B978-0-08-010472-0.50008-2](https://doi.org/10.1016/B978-0-08-010472-0.50008-2). URL: <https://www.sciencedirect.com/science/article/pii/B9780080104720500082> (visited on 02/07/2024).
- [22] Photomultiplier tube R12860 | Hamamatsu Photonics. URL: https://www.hamamatsu.com/eu/en/product/optical-sensors/pmt/pmt_tube-alone/head-on-type/R12860.html (visited on 02/08/2024).
- [23] Yan Zhang, Ze-Yuan Yu, Xin-Ying Li, Zi-Yan Deng, and Liang-Jian Wen. "A complete optical model for liquid-scintillator detectors". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 967 (July 2020), 163860. ISSN: 01689002. DOI: [10.1016/j.nima.2020.163860](https://doi.org/10.1016/j.nima.2020.163860). eprint: [2003.12212\[physics\]](https://arxiv.org/abs/2003.12212). URL: <http://arxiv.org/abs/2003.12212> (visited on 02/07/2024).

- [24] Hai-Bo Yang et al. "Light Attenuation Length of High Quality Linear Alkyl Benzene as Liquid Scintillator Solvent for the JUNO Experiment". *Journal of Instrumentation* 12.11 (Nov. 27, 2017), T11004–T11004. ISSN: 1748-0221. DOI: [10.1088/1748-0221/12/11/T11004](https://doi.org/10.1088/1748-0221/12/11/T11004). eprint: [1703.01867](https://arxiv.org/abs/1703.01867) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/1703.01867> (visited on 07/28/2023).
- [25] JUNO Collaboration et al. *The Design and Sensitivity of JUNO's scintillator radiopurity pre-detector OSIRIS*. Mar. 31, 2021. DOI: [10.48550/arXiv.2103.16900](https://doi.org/10.48550/arXiv.2103.16900). eprint: [2103.16900](https://arxiv.org/abs/2103.16900) [physics]. URL: <http://arxiv.org/abs/2103.16900> (visited on 02/07/2024).
- [26] Angel Abusleme et al. "Mass Testing and Characterization of 20-inch PMTs for JUNO". *The European Physical Journal C* 82.12 (Dec. 24, 2022), 1168. ISSN: 1434-6052. DOI: [10.1140/epjc/s10052-022-11002-8](https://doi.org/10.1140/epjc/s10052-022-11002-8). eprint: [2205.08629](https://arxiv.org/abs/2205.08629) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/2205.08629> (visited on 02/08/2024).
- [27] Yang Han. "Dual Calorimetry for High Precision Neutrino Oscillation Measurement at JUNO Experiment". AstroParticule et Cosmologie, France, Paris U. VII, APC, June 2021.
- [28] R. Acquaferredda et al. "The OPERA experiment in the CERN to Gran Sasso neutrino beam". *Journal of Instrumentation* 4.4 (Apr. 2009), P04018. ISSN: 1748-0221. DOI: [10.1088/1748-0221/4/04/P04018](https://doi.org/10.1088/1748-0221/4/04/P04018). URL: <https://dx.doi.org/10.1088/1748-0221/4/04/P04018> (visited on 02/29/2024).
- [29] JUNO collaboration et al. "Calibration Strategy of the JUNO Experiment". *Journal of High Energy Physics* 2021.3 (Mar. 2021), 4. ISSN: 1029-8479. DOI: [10.1007/JHEP03\(2021\)004](https://doi.org/10.1007/JHEP03(2021)004). eprint: [2011.06405](https://arxiv.org/abs/2011.06405) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/2011.06405> (visited on 08/10/2023).
- [30] Hans Th J. Steiger. *TAO – The Taishan Antineutrino Observatory*. Sept. 21, 2022. DOI: [10.48550/arXiv.2209.10387](https://doi.org/10.48550/arXiv.2209.10387). eprint: [2209.10387](https://arxiv.org/abs/2209.10387) [physics]. URL: <http://arxiv.org/abs/2209.10387> (visited on 01/16/2024).
- [31] Tao Lin et al. "The Application of SNiPER to the JUNO Simulation". *Journal of Physics: Conference Series* 898.4 (Oct. 2017). Publisher: IOP Publishing, 042029. ISSN: 1742-6596. DOI: [10.1088/1742-6596/898/4/042029](https://doi.org/10.1088/1742-6596/898/4/042029). URL: <https://dx.doi.org/10.1088/1742-6596/898/4/042029> (visited on 02/27/2024).
- [32] S. Agostinelli et al. "Geant4—a simulation toolkit". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 506.3 (July 1, 2003), 250–303. ISSN: 0168-9002. DOI: [10.1016/S0168-9002\(03\)01368-8](https://doi.org/10.1016/S0168-9002(03)01368-8). URL: <https://www.sciencedirect.com/science/article/pii/S0168900203013688> (visited on 02/27/2024).
- [33] J. Allison et al. "Geant4 developments and applications". *IEEE Transactions on Nuclear Science* 53.1 (Feb. 2006). Conference Name: IEEE Transactions on Nuclear Science, 270–278. ISSN: 1558-1578. DOI: [10.1109/TNS.2006.869826](https://doi.org/10.1109/TNS.2006.869826). URL: <https://ieeexplore.ieee.org/document/1610988?isnumber=33833&arnumber=1610988&count=33&index=7> (visited on 02/27/2024).
- [34] J. Allison et al. "Recent developments in Geant4". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 835 (Nov. 1, 2016), 186–225. ISSN: 0168-9002. DOI: [10.1016/j.nima.2016.06.125](https://doi.org/10.1016/j.nima.2016.06.125). URL: <https://www.sciencedirect.com/science/article/pii/S0168900216306957> (visited on 02/27/2024).
- [35] Wenjie Wu, Miao He, Xiang Zhou, and Haoxue Qiao. "A new method of energy reconstruction for large spherical liquid scintillator detectors". *Journal of Instrumentation* 14.3 (Mar. 8, 2019), P03009–P03009. ISSN: 1748-0221. DOI: [10.1088/1748-0221/14/03/P03009](https://doi.org/10.1088/1748-0221/14/03/P03009). eprint: [1812.01799](https://arxiv.org/abs/1812.01799) [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/1812.01799> (visited on 07/28/2023).
- [36] Guihong Huang et al. "Improving the energy uniformity for large liquid scintillator detectors". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1001 (June 11, 2021), 165287. ISSN: 0168-9002. DOI: [10.1016/j.nima.2021.165287](https://doi.org/10.1016/j.nima.2021.165287). URL: <https://www.sciencedirect.com/science/article/pii/S0168900221002710> (visited on 03/01/2024).
- [37] Ziyuan Li et al. "Event vertex and time reconstruction in large volume liquid scintillator detector". *Nuclear Science and Techniques* 32.5 (May 2021), 49. ISSN: 1001-8042, 2210-3147. DOI:

- 3291 [10.1007/s41365-021-00885-z](https://doi.org/10.1007/s41365-021-00885-z). eprint: 2101.08901 [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/2101.08901> (visited on 07/28/2023).
- 3292
- 3293 [38] Gioacchino Ranucci. "An analytical approach to the evaluation of the pulse shape discrimination properties of scintillators". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 354.2 (Jan. 30, 1995), 389–399. ISSN: 0168-9002. DOI: 10.1016/0168-9002(94)00886-8. URL: <https://www.sciencedirect.com/science/article/pii/0168900294008868> (visited on 03/07/2024).
- 3294
- 3295
- 3296
- 3297
- 3298 [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. URL: <https://www.sciencedirect.com/science/article/pii/S0168900206013519> (visited on 03/07/2024).
- 3299
- 3300
- 3301
- 3302
- 3303 [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. URL: <https://www.sciencedirect.com/science/article/pii/S0029554X79901708> (visited on 03/07/2024).
- 3304
- 3305
- 3306
- 3307 [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. URL: <https://doi.org/10.1007/s41365-023-01240-0> (visited on 08/17/2023).
- 3308
- 3309
- 3310
- 3311 [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. eprint: 2101.04839 [hep-ex, physics:physics]. URL: <http://arxiv.org/abs/2101.04839> (visited on 07/24/2023).
- 3312
- 3313
- 3314
- 3315
- 3316 [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. eprint: 2206.09040 [physics]. URL: <http://arxiv.org/abs/2206.09040> (visited on 07/24/2023).
- 3317
- 3318
- 3319
- 3320
- 3321 [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. URL: <https://dx.doi.org/10.1088/1748-0221/17/11/P11003> (visited on 04/04/2024).
- 3322
- 3323
- 3324
- 3325 [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. URL: <https://dx.doi.org/10.1088/1748-0221/16/10/C10011> (visited on 04/04/2024).
- 3326
- 3327
- 3328
- 3329
- 3330 [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. URL: <https://dx.doi.org/10.1088/1748-0221/16/07/P07041> (visited on 04/04/2024).
- 3331
- 3332
- 3333
- 3334 [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. URL: <https://link.aps.org/doi/10.1103/PhysRevD.102.092003> (visited on 04/04/2024).
- 3335
- 3336
- 3337
- 3338 [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. URL: <https://ui.adsabs.harvard.edu/abs/2005ApJ...622..759G> (visited on 04/04/2024).
- 3339
- 3340
- 3341
- 3342

- 3343 [49] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. *Convolutional Neural Networks*
3344 on Graphs with Fast Localized Spectral Filtering. Feb. 5, 2017. DOI: [10.48550/arXiv.1606.09375](https://doi.org/10.48550/arXiv.1606.09375).
3345 eprint: [1606.09375\[cs, stat\]](https://arxiv.org/abs/1606.09375). URL: <http://arxiv.org/abs/1606.09375> (visited on
3346 04/04/2024).
- 3347 [50] JUNO Collaboration et al. "JUNO Physics and Detector". *Progress in Particle and Nuclear Physics*
3348 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).
- 3350 [51] Leo Breiman, Jerome Friedman, R. A. Olshen, and Charles J. Stone. *Classification and Regression*
3351 *Trees*. New York: Chapman and Hall/CRC, Oct. 25, 2017. 368 pp. ISBN: 978-1-315-13947-0. DOI:
3352 [10.1201/9781315139470](https://doi.org/10.1201/9781315139470).
- 3353 [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:
3354 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).
- 3355 [53] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. Jan. 29, 2017.
3356 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).
- 3357 [54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep Residual Learning for Image
3358 Recognition". *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016
3359 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). ISSN: 1063-6919. June
3360 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).
- 3361 [55] Olga Russakovsky et al. *ImageNet Large Scale Visual Recognition Challenge*. Jan. 29, 2015. DOI:
3362 [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>
3363 (visited on 05/17/2024).
- 3364 [56] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image*
3365 *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:
3366 <http://arxiv.org/abs/1409.1556> (visited on 05/17/2024).
- 3367 [57] Anna Allen. *generic-github-user/Image-Convolution-Playground*. original-date: 2018-09-28T22:42:55Z.
3368 July 15, 2024. URL: <https://github.com/generic-github-user/Image-Convolution-Playground> (visited on 07/16/2024).
- 3369 [58] Jason Ansel et al. *PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Trans-*
3370 *formation and Graph Compilation*. Publication Title: 29th ACM International Conference on Ar-
3371 *chitectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS*
3372 '24) original-date: 2016-08-13T05:26:41Z. Apr. 2024. DOI: [10.1145/3620665.3640366](https://doi.org/10.1145/3620665.3640366). URL:
3373 <https://pytorch.org/assets/pytorch2-2.pdf> (visited on 07/16/2024).
- 3374 [59] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document
3375 recognition". *Proceedings of the IEEE* 86.11 (Nov. 1998). Conference Name: Proceedings of the
3376 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).
- 3377 [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).
- 3378 [61] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. *Neural*
3379 *Message Passing for Quantum Chemistry*. June 12, 2017. DOI: [10.48550/arXiv.1704.01212](https://doi.org/10.48550/arXiv.1704.01212).
3380 eprint: [1704.01212\[cs\]](https://arxiv.org/abs/1704.01212). URL: <http://arxiv.org/abs/1704.01212> (visited on 05/22/2024).
- 3381 [62] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. *Diffusion Convolutional Recurrent Neural*
3382 *Network: Data-Driven Traffic Forecasting*. Feb. 22, 2018. DOI: [10.48550/arXiv.1707.01926](https://doi.org/10.48550/arXiv.1707.01926).
3383 eprint: [1707.01926\[cs, stat\]](https://arxiv.org/abs/1707.01926). URL: <http://arxiv.org/abs/1707.01926> (visited on
3384 05/22/2024).
- 3385 [63] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
3386 Ozair, Aaron Courville, and Yoshua Bengio. *Generative Adversarial Networks*. June 10, 2014. DOI:
3387 [10.48550/arXiv.1406.2891](https://doi.org/10.48550/arXiv.1406.2891).

- 3395 10 . 48550 / arXiv . 1406 . 2661 . eprint: 1406 . 2661 [cs , stat] . URL: <http://arxiv.org/abs/1406.2661> (visited on 05/29/2024).
- 3396
- 3397 [64] Anatael Cabrera et al. *Multi-Calorimetry in Light-based Neutrino Detectors*. Dec. 20, 2023. DOI: 10 . 48550 / arXiv . 2312 . 12991 . eprint: 2312 . 12991 [hep-ex , physics : physics] . URL: <http://arxiv.org/abs/2312.12991> (visited on 08/19/2024).
- 3398
- 3399 [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).
- 3400
- 3401
- 3402
- 3403 [66] Dan Ciresan, Ueli Meier, and Juergen Schmidhuber. *Multi-column Deep Neural Networks for Image Classification*. version: 1. Feb. 13, 2012. DOI: 10 . 48550 / arXiv . 1202 . 2745 . eprint: 1202 . 2745 [cs] . URL: <http://arxiv.org/abs/1202.2745> (visited on 06/27/2024).
- 3404
- 3405
- 3406 [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 . eprint: 2101 . 11589 [hep-ex] . URL: <http://arxiv.org/abs/2101.11589> (visited on 06/27/2024).
- 3407
- 3408
- 3409
- 3410 [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 . URL: <https://dx.doi.org/10.1088/1475-7516/2021/11/051> (visited on 06/27/2024).
- 3411
- 3412
- 3413
- 3414
- 3415 [69] Taco S. Cohen, Mario Geiger, Jonas Koehler, and Max Welling. *Spherical CNNs*. Feb. 25, 2018. DOI: 10 . 48550 / arXiv . 1801 . 10130 . eprint: 1801 . 10130 [cs , stat] . URL: <http://arxiv.org/abs/1801.10130> (visited on 07/13/2024).
- 3416
- 3417
- 3418 [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).
- 3419
- 3420 [71] NVIDIA V100. NVIDIA. URL: <https://www.nvidia.com/en-gb/data-center/v100/> (visited on 08/06/2024).
- 3421
- 3422 [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).
- 3423
- 3424 [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 / IEEEESTD . 2019 . 8766229 . URL: <https://ieeexplore.ieee.org/document/8766229> (visited on 07/03/2024).
- 3425
- 3426
- 3427
- 3428 [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 . eprint: 2102 . 11538 [hep-ex , physics : physics] . URL: <http://arxiv.org/abs/2102.11538> (visited on 02/08/2024).
- 3429
- 3430
- 3431
- 3432
- 3433 [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 . eprint: astro-ph/0409513 . URL: <http://arxiv.org/abs/astro-ph/0409513> (visited on 08/10/2023).
- 3434
- 3435
- 3436
- 3437
- 3438 [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 . URL: <https://dx.doi.org/10.1088/1674-1137/41/6/066201> (visited on 08/16/2024).
- 3439
- 3440
- 3441
- 3442
- 3443 [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).
- 3444
- 3445 [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*
- 3446

- 3447 *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](https://arxiv.org/abs/1804.09456) [physics]. URL: <http://arxiv.org/abs/1804.09456> (visited on 09/17/2024).
- 3448
- 3449 [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). URL: <https://dx.doi.org/10.1088/1742-6596/2145/1/012017> (visited on 09/17/2024).
- 3450
- 3451 [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).
- 3452
- 3453 [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).
- 3454
- 3455 [82] Angel Abusleme et al. "Potential to Identify the Neutrino Mass Ordering with Reactor Antineutrinos in JUNO" (May 2024). eprint: 2405.18008.
- 3456
- 3457 [83] Rene Brun et al. *root-project/root: v6.26/06*. Version v6-26-06. Mar. 3, 2022. DOI: [10.5281/zenodo.3895860](https://doi.org/10.5281/zenodo.3895860). URL: <https://zenodo.org/records/3895860> (visited on 09/05/2024).
- 3458
- 3459 [84] 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).
- 3460
- 3461 [85] 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).
- 3462
- 3463 [86] 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).
- 3464
- 3465 [87] "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).
- 3466
- 3467 [88] 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).
- 3468
- 3469
- 3470
- 3471
- 3472
- 3473
- 3474
- 3475
- 3476
- 3477
- 3478
- 3479
- 3480
- 3481
- 3482
- 3483
- 3484
- 3485

3486

3487

Titre : Méthode Deep Learning and analyse Double Calorimétrique pour la mesure de haute précision des paramètres d'oscillation des neutrinos dans JUNO

Mot clés : Neutrinos; expérience JUNO; Deep Learning; reconstruction d'IBD; oscillations des neutrinos; double calorimetrie

Résumé : JUNO est un observatoire de neutrinos à scintillateur liquide, polyvalent et medium baseline (environ 52 km), situé en Chine. Ses principaux objectifs sont de mesurer les paramètres d'oscillation θ_{12} , Δm_{21}^2 et Δm_{31}^2 avec une précision au pour-mille et de déterminer l'ordre des masses des neutrinos avec un niveau de confiance de 3σ . Atteindre ces objectifs nécessite une résolution énergétique sans précédent de $3\%/\sqrt{E(\text{MeV})}$ avec cette technologie. Cela demande une compréhension approfondie des divers effets au sein du détecteur.

Le système de double calorimetrie, composé de deux systèmes de mesure distincts observant le même événement, permet non seulement une calibration mais aussi une détection des effets du détecteur avec une grande précision, comme démontré dans cette thèse. Le Deep Learning, un outil de plus en plus utilisé en physique expérimentale, joue un rôle crucial dans cet effort. Dans cette thèse, je présente le développement, l'application et l'analyse des techniques de Deep Learning pour la reconstruction d'évènements dans l'expérience JUNO.

3519

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

Keywords: Neutrinos; JUNO experiment; Deep learning; IBD reconstruction; neutrinos Oscillation; dual Calorimetry

Abstract: JUNO is a multipurpose, medium baseline (~ 52 km) liquid scintillator neutrino observatory located in China. Its primary objectives are to measure the oscillation parameters θ_{12} , Δm_{21}^2 , and Δm_{31}^2 with per mil precision and to determine the neutrino mass ordering at a 3σ confidence level. Achieving these goals requires an unprecedented energy resolution of $3\%/\sqrt{E(\text{MeV})}$ with this technology. This demands a comprehensive understanding of the various effects within the

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.

