



HAL
open science

Comparative Assessment of Feature Extraction for Fast Neutron Spectra Prediction with Machine Learning Algorithms using a CVD Diamond Detector

Enrica Belfiore, Pierre-Guy Allineï, Pierre Houedry, Meriem Bahhi, Simon Bartolacci, Adel Saleh, Mehdi Ben Mosbah, Rodolphe Antoni, Abdallah Lyoussi,
Jean-Emmanuel Groetz

► **To cite this version:**

Enrica Belfiore, Pierre-Guy Allineï, Pierre Houedry, Meriem Bahhi, Simon Bartolacci, et al.. Comparative Assessment of Feature Extraction for Fast Neutron Spectra Prediction with Machine Learning Algorithms using a CVD Diamond Detector. 2025. <hal-05123189>

HAL Id: hal-05123189

<https://hal.science/hal-05123189v1>

Preprint submitted on 20 Jun 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Comparative Assessment of Feature Extraction for Fast Neutron Spectra Prediction with Machine Learning Algorithms using a CVD Diamond Detector

Enrica Belfiore¹, Pierre-Guy Allineï¹, Pierre Houedry², Meriem Bahhi³, Simon Bartolacci⁴, Adel Saleh⁵, Mehdi Ben Mosbah¹, Rodolphe Antoni¹, Abdallah Lyoussi⁶ and Jean-Emmanuel Groetz⁷

¹ CEA, DES, IRESNE, DTN, Cadarache, Saint-Paul-Lez-Durance, 13108, France

² IRISA, École normale supérieure de Rennes, Rennes, 35170, France

³ Institut de Mathématiques de Bourgogne, Université de Bourgogne, Dijon, 21078, France

⁴ Institut Élie Cartan de Lorraine, Université de Lorraine, Nancy, 54012, France

⁵ Institut de Mathématique de Marseille, Aix-Marseille University, Marseille, 13009, France

⁶ CEA, DES, IRESNE, DER, Cadarache, Saint-Paul-Lez-Durance, 13108, France

⁷ Laboratoire Chrono-environnement, Université de Franche-Comté, Besançon, 25030, France

Abstract - Machine learning algorithms are recognized as effective tools for addressing challenges in nuclear methodologies and instrumentation, particularly when obtaining experimental data or precise mathematical models is difficult. One notable challenge in fast neutron spectroscopy involves unfolding neutron spectra from solid-state detectors, where conventional deconvolution techniques often suffer from instability. This study uses the PHITS code to simulate the behavior of a CVD diamond detector subjected to various neutron spectra. The simulated data serve as input to a regression model aimed at reconstructing the original spectra. In particular, we investigate the impact of four different signal processing and feature extraction techniques—Fourier transform, nodal basis projection, wavelet decomposition and agglomerative clustering-based compression—on the performance of spectrum reconstruction. These transformations are applied as preprocessing steps, and the resulting features are then used to train a supervised regression model (linear or kernel ridge, depending on the method). Model performance is evaluated using the R^2 metric against reference spectra from the IAEA Compendium database. The most effective approach is then applied to experimental data collected with a CVD diamond detector at the DANAIDES irradiation facility, part of the TOTEM platform in Cadarache. This final test, performed on 14 MeV and AmBe neutron sources, demonstrates the practical feasibility of the proposed methodology for real-world neutron spectrum reconstruction.

Keywords —Solid detectors, PHITS, neutron spectrum prediction.

1. INTRODUCTION

The application of machine learning (ML) in radiation detection has garnered significant interest in recent years, driven by the need for more efficient, accurate, and autonomous systems. Traditional methods for radiation detection, which often rely on manual analysis and classical statistical techniques, face limitations in terms of processing speed, scalability, and adaptability to complex environments. Machine learning, with its ability to learn patterns from large datasets and improve over time, offers promising solutions to these challenges [1]. Within neutron spectroscopy, ML approaches offer the potential to gather data in scenarios where measurements are challenging with current technologies, thus facilitating the development of innovative measurement techniques. Neutron measurement is crucial for obtaining key nuclear physics data and ensuring real-time control in future reactors. Detectors in this environment must withstand harsh conditions, including intense radiation, high temperatures, and space constraints [2]. Semiconductor detectors have prompted researchers to explore their potential application in fast neutron spectroscopy [3] and they are effectively used in the Joint European Torus (JET) and will be integrated into International Thermonuclear Experimental Reactor (ITER) [4]. However, solid-state detectors face several limitations in neutron spectrum unfolding. These methods often involve matrix inversion and deconvolution procedures, which can lead to oscillatory artefacts and numerical instabilities when over-iterated in pursuit of unattainable resolution. In such cases, regularization or smoothing techniques are typically introduced to stabilize the solution, although these may themselves introduce further uncertainties.

Researchers from CEA Cadarache are exploring new methods to reconstruct the neutron spectrum from data collected by solid-state detectors through ML techniques [5, 6]. The feasibility of this approach was initially demonstrated in a theoretical study using a Tissue Equivalent Proportional Counter gas detector [7]. However, this detector was found to be unsuitable under high neutron fluxes. Consequently, the approach was applied to a Silicon-Carbide (SiC) solid-state detector, which emerged as a promising candidate for these conditions, yielding favorable results in neutron spectrum prediction via ML, supported by theoretical studies in [8]. The challenge encountered with the SiC detector mainly lies in managing the numerous resonances in the cross sections for neutron-silicon interactions. In this work, we explore an alternative detector material: Chemical Vapour Deposition (CVD) diamond. Recent advances in CVD techniques have significantly reduced defect densities, making diamond-based detectors more viable for neutron spectrometry in harsh environments [9, 10, 11]. Using the PHITS Monte Carlo code, we simulate the energy deposition spectra from CVD diamond detectors exposed to various neutron sources. The simulated deposited energy data in the detector are then processed using four distinct signal transformation techniques—Fourier transform, nodal basis projection, wavelet

55 decomposition, and agglomerative clustering—which serve as feature extraction steps to reduce dimensionality and extract relevant
56 patterns from the raw signals. These features are used to train regression models (either linear or kernel ridge, depending on the
57 case), with the goal of reconstructing the original neutron energy spectrum. Preliminary attempts to apply regression models
58 directly to the raw input data were discarded, as such models suffered from poor generalization due to the high dimensionality and
59 redundancy in the signal, which led to overfitting and unstable predictions. Therefore, the use of feature extraction methods proved
60 essential to reduce complexity, improve robustness, and enable meaningful learning. The model performance is evaluated using
61 the coefficient of determination (R^2) against reference spectra from the IAEA Compendium database. Finally, the best-performing
62 approach is applied to experimental data acquired with a CVD diamond detector exposed to 14 MeV and Americium-Beryllium
63 (AmBe) neutron sources at the DANAIDES irradiation casemate of the TOTEM facility at CEA Cadarache, demonstrating the
64 applicability of the proposed method in real-world conditions. This paper is structured as follows: Section 2 presents a comparison
65 between CVD diamond and SiC detectors in the context of high-flux neutron environments. Section 3 describes the detector model
66 and the PHITS simulations used to generate the data. Section 4 outlines the data processing and machine learning methodology,
67 detailing the signal transformations and regression models employed. Section 5 presents the reconstruction results from both
68 simulated and experimental data. Finally, Section 6 offers conclusions and future perspectives.

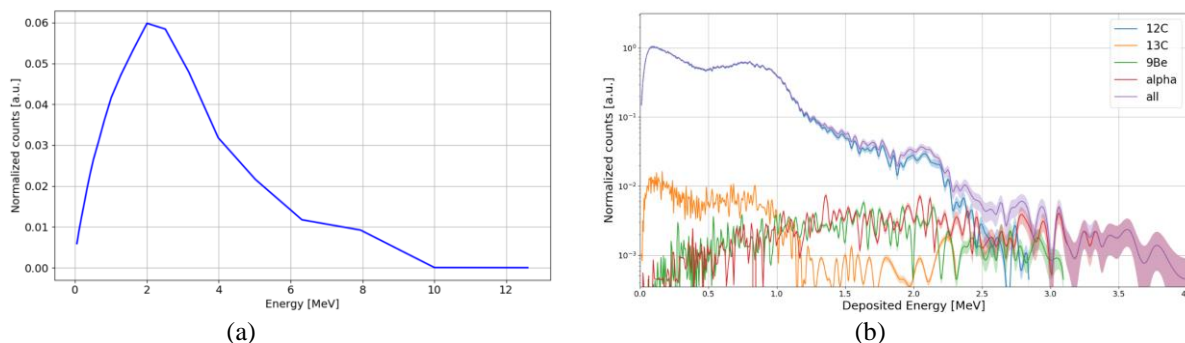
69 2. CVD DIAMOND VS SiC DETECTORS

70 Comparative studies on the performance of CVD diamond and SiC detectors have gained considerable attention in fast neutron
71 spectroscopy [12]. These detectors are among the most commonly used and commercially available semiconductor devices. Their
72 operating principle relies on the collection of charge carriers generated by ionizing particles traversing the semiconductor material.
73 Electrodes implanted within the semiconductor establish an electric drift field, which directs the movement of the ionized charges,
74 leading to the collection of these charges and the generation of a detectable signal [13]. Both sensor types are considered promising
75 candidates for use in harsh environments due to several advantageous properties. In particular, they exhibit high displacement
76 threshold energies—ranging from 20–35 eV for SiC and 40–50 eV for CVD diamond—as well as wide bandgap energies of 3.27
77 eV and 5.5 eV, respectively. Additionally, they possess excellent thermal conductivities: 4.9 W/cm·K for the former and 22
78 W/cm·K for the latter. These features make them suitable for high-temperature applications without the need for additional cooling
79 systems, and also render them more radiation resistant compared to other semiconductors such as germanium or gadolinium, etc
80 [14, 15, 16]. In the past, diamond was dismissed as a semiconductor material for detection applications due to its high cost and the
81 presence of impurities and defects in its crystal structure, which could lead to charge carrier issues and polarization effects.
82 However, with advancements in manufacturing techniques such as CVD [17], which significantly reduce defect density, diamond-
83 based sensors have become more accessible for industrial applications. In the analysis in [12], a comparison between the 4H-SiC
84 p+-n diode and a single crystal CVD diamond detector is reported, considering 14 MeV energy neutron flux generated by a D-T
85 neutron generator at the Technical University of Dresden. This study highlights that the selection of the most suitable sensor
86 among these two is heavily influenced by the specific application, as both detectors in this case exhibit well-characterized pulse-
87 height spectra with a clearly defined peak resulting from the $^{12}\text{C}(n, \alpha)^9\text{Be}$ reaction. While the energy resolution of this peak is
88 marginally superior for the 4H-SiC-based detector, the diamond-based detector demonstrates a higher count rate.
89 In our application, which involves modelling simulations to populate the database used by the ML algorithm for neutron spectrum
90 prediction, the SiC detector appears to be a more promising candidate. This is due to the inclusion of interactions with silicon,
91 which enhance the variability of the system and increase the range of cases that can be addressed. However, previous studies have
92 highlighted that discretizing the energy grid poses greater challenges because it must contend with multiple resonances in the cross
93 sections for Si-n interactions. In this work, we therefore examine the CVD diamond detector for the same application, as it does
94 not require such fine discretization of the energy grid.

95 3. CVD DIAMOND DETECTOR

96 To train and validate the machine learning models used in this study, a large dataset was generated via Monte Carlo simulations
97 using the PHITS code [20]. This versatile Monte Carlo particle transport simulation software is capable of accurately modelling
98 particle transport over a wide energy range, employing multiple nuclear reaction models and nuclear data libraries. The simulated
99 detector is a single-crystal CVD diamond sensor with a sensitive volume of $1.4 \times 1.4 \text{ mm}^2$ surface area and $320 \text{ }\mu\text{m}$ thickness. The
100 design of the CVD diamond detector was analyzed in [21]. The purpose of the simulation is to produce an energy-dependent dataset
101 that mimics the deposited energy in the detector, obtaining a set of realistic deposited energy patterns corresponding to various
102 incident neutron spectra. These data serve as input for supervised regression models aiming to reconstruct the original neutron
103 energy distribution. The detector was therefore modelled as a parallelepiped with a thickness of $320 \text{ }\mu\text{m}$ and a width of 1.4 mm,
104 considering only the sensitive volume of the detector. It is assumed that excluding other parts in the detector design does not
105 influence the output of the simulations since C-n interactions mainly occur within the sensitive volume of the detector. The source
106 geometry is a rectangular solid, positioned in front of the detector. The neutrons are emitted in the z-direction, with an emission
107 angle of 0° (perpendicular to the detector surface, along the +z-axis). Ten batches were set in the simulations, with 50,000,000
108 particles for each batch. In defining the materials, it was decided to consider a detector made of pure carbon, without doping or
109 impurities. The carbon composition used is 98.94 % ^{12}C and 1.06% ^{13}C , in line with standard natural abundance [22]. All the
110 nuclear data were taken from the JENDL-4.0 nuclear data library [23]. The spectrum was discretized in 600 energy bins, log-equal
111 spaced from 100 keV to 15 MeV. The results of the simulations will provide the amount of energy deposited in the detector at each

112 point of the energy grid. These energy deposition values will serve as the features in the input vector for the machine learning
 113 models, as described in the following section.
 114 To illustrate the purpose of the detector simulation, Figure 1 shows an example of a neutron energy spectrum taken from the IAEA
 115 Compendium database (Turkey Point reactor, fuel element at 26 in) and the corresponding deposited energy spectrum obtained
 116 from the PHITS simulation of the CVD diamond detector. These input-output pairs represent the data used to train the supervised
 117 regression model, which aims to learn the inverse mapping from the deposited energy signal (output of the simulation) to the
 118 original neutron spectrum (input of the simulation).



119
 120
 121 **Fig. 1.** (a) Neutron energy spectrum from the IAEA Compendium database (Turkey Point reactor, fuel element at 26 in), (b)
 122 Corresponding deposited energy spectrum simulated with PHITS for a CVD diamond detector.

123 Figure 1a shows the input neutron energy spectrum for the simulation taken from the IAEA Compendium database, corresponding
 124 to a measurement at the Turkey Point reactor (fuel element at 26 inches). As evident, the spectrum contains a relatively low number
 125 of neutrons with energies exceeding 6 MeV. This energy marks the approximate threshold for the activation of the $^{12}\text{C}(n, \alpha)^9\text{Be}$
 126 reaction, which is one of the few channels leading to heavy charged particle production (alphas and beryllium nuclei) in carbon.
 127 As a result, the corresponding deposited energy spectrum obtained through PHITS simulation of the CVD diamond detector (Figure
 128 1b) is dominated by the contribution from neutron elastic scattering with carbon atoms (shown in blue). The contributions from
 129 alpha particles (in red) and beryllium nuclei (in green), resulting from the $^{12}\text{C}(n, \alpha)^9\text{Be}$ reaction, are present but significantly
 130 smaller due to the limited flux above threshold energy. The total deposited energy spectrum (shown in violet), which includes all
 131 particle contributions, is incorporated in the testing database for the machine learning model described in the following section.

132 4. MACHINE LEARNING APPROACH

133 This section outlines the methodology adopted to reconstruct the incident neutron spectrum from the simulated energy deposition
 134 in a CVD diamond detector using supervised regression techniques. Unlike classical unfolding approaches that rely on detector
 135 response functions, this method directly uses the simulated energy deposition profiles as input to train regression models capable
 136 of inferring the original neutron spectra. The use of machine learning in nuclear spectroscopy is gaining traction, with several
 137 studies exploring the application of regression and deep learning models for neutron and gamma-ray analysis [18, 19].

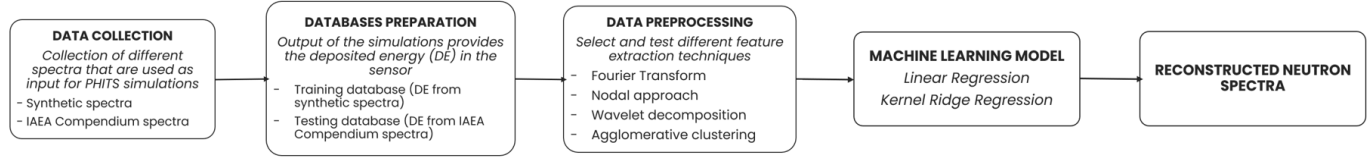
138 The training dataset used in this study was generated using the PHITS Monte Carlo code by simulating the response of the detector
 139 to 544 synthetic neutron spectra, designed to resemble realistic shapes as described in [8]. For each simulated case, the total energy
 140 deposited in the sensitive volume of the detector was recorded, producing spectra segmented into 600 logarithmically spaced bins
 141 from 0.1 MeV to 15 MeV. These spectra constitute the explanatory variables (input features) for the regression task.

142 The target variables are the corresponding neutron spectra used in the simulation, originally segmented into 55 bins across the
 143 energy range from 2×10^{-8} MeV to 44 MeV, in accordance with the IAEA Neutron Compendium format [24]. However, many bins
 144 in the energy range below 0.05 MeV contain negligible information for the CVD detector, which has low sensitivity in this region.
 145 For this reason, only the energy range from 0.05 MeV to 15 MeV, corresponding to 25 bins, was retained for training and
 146 evaluation. To address the high dimensionality of the input (600 bins), four signal processing techniques were adopted for feature
 147 extraction and dimensionality reduction prior to regression:

- 148 • **Fourier Transform**, which captures global spectral trends and periodicities;
- 149 • **Nodal basis projection**, which approximates the spectra using orthonormal polynomial bases tailored to spatially
 150 distributed signals;
- 151 • **Wavelet decomposition**, which captures localized features and discontinuities at multiple resolutions;
- 152 • **Agglomerative clustering**, which groups adjacent input bins based on their similarity, effectively reducing
 153 dimensionality while preserving local structure and correlations in the deposited energy signal.

154 These methods were selected to explore complementary perspectives on the input data, offering different trade-offs between
 155 localization, smoothness, sparsity, and interpretability in the reduced feature space. While Fourier and nodal decompositions are
 156 most suitable for capturing globally smooth patterns, wavelets and agglomerative clustering offer advantages in preserving
 157 localized or hierarchical structures. The use of feature extraction is crucial in this context, as it reduces computational complexity,

158 mitigates the curse of dimensionality, and enhances the model’s ability to generalize to unseen data. Initial attempts to apply
 159 regression models directly to the raw 600-dimensional input data resulted in poor predictive performance and high computational
 160 cost. These outcomes can be attributed to the high dimensionality and redundancy of the input features, which introduce noise and
 161 multicollinearity, thereby impairing the learning process and leading to unstable or biased estimators.
 162 All models were trained using a supervised learning framework implemented with scikit-learn [25], and the predicted neutron
 163 spectra were post-processed to enforce two physical constraints: non-negativity and normalization (i.e., the integral of each
 164 spectrum must equal one). An overview of the complete described reconstruction steps—from simulated neutron source to
 165 spectrum estimation—is presented in Figure 2.
 166



167
 168 **Fig. 2.** Steps of the ML methodology for the reconstruction of neutron spectra with CVD diamond sensor.
 169

170 4.1 FOURIER SERIES

171 In the context of signal processing, it was considered to explore a Fourier-based approach. The idea involves treating each deposited
 172 energy spectrum as a function $x(t)$ on $[0,1]$, which is sampled at 600 points (corresponding to the number of energy bins in which
 173 the energy grid is discretized in the simulations). Subsequently, the Fourier series of this function can be computed.
 174 For $n \in \mathbb{Z}$, the n -order Fourier series of x is expressed as follows:
 175

$$176 S_n(x(t)) = \frac{1}{2} a_0(x) + \sum_{k=1}^n (a_k(x) \cos(kt2\pi) + b_k(x) \sin(kt2\pi)) \quad (1)$$

177 where:

$$178 - a_0(x) = \int_0^1 x(t) dt;$$

$$179 - b_0(x) = 0;$$

$$180 - a_n(x) = 2 \int_{[0,1]} x(t) \cos(nt2\pi) dt, \text{ for } n > 0;$$

$$181 - b_n(x) = 2 \int_{[0,1]} x(t) \sin(nt2\pi) dt, \text{ for } n > 0.$$

182 With X denoting the deposited energy spectrum data and Y the target neutron spectrum data.

183 The Fourier coefficients were calculated up to the order 30 for each deposited energy spectrum $x \in X$, and up to the order 5 for
 184 each neutron spectrum $y \in Y$. This resulted in two new datasets: $X_{Fourier}$ of size 544x60 and $Y_{Fourier}$ of size 544x10. Subsequently,
 185 a multi-output linear regression (with L2-penalization, $\alpha = 10^{-6}$) was trained to map $X_{Fourier} \rightarrow Y_{Fourier}$.
 186 Finally, using the inverse of the Fourier transform equation (1) it is possible to reconstruct the prediction of Y .
 187

188 4.2 NODAL APPROACH

189 The second method proposed relies similarly on regarding the deposited energy spectrum as a function $X(f)$, where f represents
 190 a energy bin within the feasible energy range $[f_{min}, f_{max}]$. As previously stated, each row of the dataset comprises observations of
 191 the deposited energy X at bins $f_{min} = f_1 < f_2 < \dots < f_N = f_{max}$, where N represents the total number of energy bins, i.e. 600, in
 192 the considered CVD diamond detector model. Consequently, each row generates a linear spline interpolation function $X_N(f)$ of X
 193 according to:
 194

$$195 X_N(f) = \sum_{k=1}^N X_k \varphi_k(f), \quad (2)$$

196 where φ_k satisfies the following conditions:

- 197 - it is continuous on $[f_{min}, f_{max}]$;
- 198 - it is piecewise linear on each $[f_k, f_{k+1}]$;
- 199 - it satisfies $\varphi_k = \delta_{kj}$ for all $k, j = 1, \dots, N$.

200 In equation (2) X_k represents the observed value of $X(f_k)$, i.e. the value in the k -th column, which is assumed to be a good
 201 approximation of $X(f_k)$. Instead of applying linear regression on X_N , it was decided to first select a significantly smaller number
 202 of bins $f_{min} = \tilde{f}_1 < \dots < \tilde{f}_n = f_{max}$ with $n \ll N$ (arbitrary chosen), and then to define the vector:
 203
 204

$$205 X_n = (x_1, \dots, x_n) := (X_N(\tilde{f}_1), X_N(\tilde{f}_2), \dots, X_N(\tilde{f}_n)). \quad (3)$$

206 We denote X_{Nod} as the new dataset of deposited energies of size 544xn, where each row is represented by X_n , and Y remains
 207

208 unchanged with a size of 544x25. Linear regression (with L2-penalization, $\alpha=10^{-6}$) is then applied from X_{Nod} to Y . Furthermore,
 209 \tilde{f}_k lies between two consecutive bins f_j and f_{j+1} . Defining the parameter $\theta_k \in [0, 1]$, it can be expressed as $\tilde{f}_k = \theta_k f_j + (1 -$
 210 $\theta_k)f_{j+1}$. From this formulation and the definition of X_N , it follows:

$$211 \quad x_k := X_N(\tilde{f}_k) = \theta_k X_j + (1 - \theta_k) X_{j+1} \quad , \quad (4)$$

212
 213 in other words, $X_N(\tilde{f}_k)$ is just a weighted average of the two consecutive values of X_j and X_{j+1} . In our implementation, the bins \tilde{f}_i
 214 are equidistant. In addition, instead of taking x_k as in (4), it was simplified with:

$$215 \quad x_k = \frac{1}{2}(X_j + X_{j+1}) \quad (5)$$

216
 217 as this did not affect the R^2 error much (as will be shown in the results of the test sets). In addition, a filter has been applied to X_n
 218 so that the integral is equal to 1 and that negative values are truncated to 0 similarly to the Fourier basis case.

219 4.3 WAVELET DECOMPOSITION

220 The Discrete Wavelet Transform is used to decompose the deposited energy into its approximation and detail components at
 221 different scales. This decomposition is achieved by applying low-pass and high-pass filters to the energy bins. Given a energy
 222 range divided in N bins $f = \{f_0, f_1, \dots, f_{N-1}\}$, the goal is to express it in terms of scaling functions and wavelet functions at different
 223 levels. This can be written mathematically as:

$$224 \quad f(t) = \sum_n c_A^1[n] \phi(t - n) + \sum_n c_D^1[n] \psi(t - n) \quad , \quad (6)$$

225 where:

- 226 1. $c_A^1[n]$ are the approximation coefficients,
- 227 2. $c_D^1[n]$ are the detail coefficients,
- 228 3. $\phi(t)$ is the scaling function (low-pass filter),
- 229 4. $\psi(t)$ is the wavelet function (high-pass filter).

230 For the Haar wavelets [26], the scaling function $\phi(t)$ and wavelet function $\psi(t)$ are defined as follows:

$$231 \quad \phi(t) = \begin{cases} 1 & \text{if } 0 \leq t < 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad \psi(t) = \begin{cases} 1 & \text{if } 0 \leq t < \frac{1}{2} \\ -1 & \text{if } \frac{1}{2} \leq t < 1 \\ 0 & \text{otherwise} \end{cases} . \quad (7)$$

232 The coefficients are given by:

$$233 \quad c_A^1[n] = \sum_{m=0}^1 h[m] f_{2n-m} \quad \text{and} \quad c_D^1[n] = \sum_{m=0}^1 g[m] f_{2n-m} \quad , \quad (8)$$

234 where:

$$235 \quad h = \left[\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right] \quad \text{and} \quad g = \left[-\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right] . \quad (9)$$

236 Finally, a Kernel Ridge Regression (KRR) model was applied between wavelet coefficients for a level equal to 1 and each bin of
 237 neutron spectrum ($y \in Y$). The hyperparameters were adjusted with a grid search only on:

- 238 1. the regularization parameter α of the RIDGE model.
- 239 2. the convolution kernel: linear, gaussian or cosine.

240 The best results were obtained with $\alpha=0.01$ and cosine kernel.

241 4.4 AGGLOMERATIVE CLUSTERING

242 This model was chosen to maintain consistency with the study on SiC. The algorithm is based on a regularized linear method (l^2
 243 regularization) associated with a convolution kernel [27]. The overall approach is structured as a two-stage pipeline, combining
 244 feature reduction via agglomerative clustering with supervised learning using Kernel Ridge Regression. The 600 energy bins of
 245 each deposited energy spectrum are reduced using agglomerative clustering with spatial connectivity constraints. Specifically, each
 246 energy bin is treated as a node in a graph, and edges are formed only between adjacent bins, enforcing locality-preserving clustering.
 247 This ensures that merged clusters correspond to contiguous energy regions. This hierarchical clustering minimizes a Ward-like
 248 objective function:

$$249 \quad \mathcal{L} = \sum_{c=1}^n \sum_{x_i \in C_c} \|x_i - \mu_c\|^2 \quad , \quad (10)$$

250

251

where C_c denotes the cluster c , μ_c its mean, and $n \in \{20,40,60,80,100\}$ is the number of clusters. The cluster centroids (means) are then used as the new features for downstream modeling, effectively reducing the input dimensionality from 600 to n . The KRR model is trained for each of the 25 neutron spectrum bins (i.e., a one-vs-all multi-output regression setup). The KRR model applies l^2 -regularization to linear regression in a reproducing kernel Hilbert space, using a convolution kernel for nonlinearity. The prediction function for a single output bin is:

$$\hat{y} = \sum_{i=1}^m \alpha_i K(x_i, x) \quad , \quad (11)$$

where $K(\cdot, \cdot)$ is the kernel function and α_i are learned weights. Two types of kernels were explored, the Radial Basis Function (RBF)

$$K(x, x') = \exp\left(-\gamma \|x - x'\|^2\right) \quad (12)$$

and the polynomial kernel

$$K(x, x') = (\gamma \langle x, x' \rangle + r)^d \quad , \quad (13)$$

with degree $d \in \{2, 3, 4\}$, and γ, r as tunable parameters.

The hyperparameters adjusted include:

- The number of clusters of explanatory variables;
- The regularization parameter α of the RIDGE model;
- The convolution kernel: polynomial or RBF;
- The γ parameter used in distance calculation via RBF and polynomial kernels;
- The degree of the polynomial kernel: 2, 3, or 4.

An identical model is created for each bin, and then the hyperparameters of each of these models are optimized using a random search grid for the best set of hyperparameters.

A summary of the main characteristics of the proposed approaches is provided in Table 1.

Dimension reduction method	ML model	Input/output dimension	Reduced dimension	Regression details
Fourier basis	Linear Regression	600/25	60	$\alpha = 10^{-6}$
Nodal approach (linear interpolation)	Linear Regression	600/25	25	$\alpha = 10^{-6}$
Wavelet decomposition (Haar DWT)	Kernel Ridge Regression	600/25	300	Cosine kernel, $\alpha = 0.01$
Agglomerative clustering	Kernel Ridge Regression	600/25	20-100	RBF/ polynomial kernel

Table 1. Summary of reduction methods and characteristics.

5. RESULTS

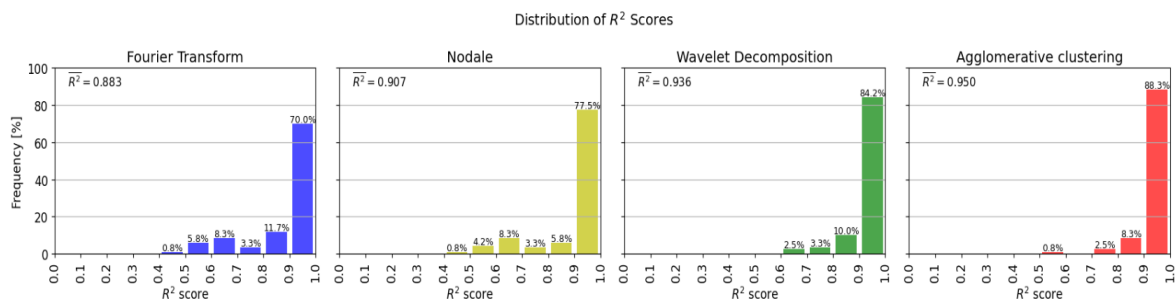
5.1 PREDICTION OF NEUTRON SPECTRA FROM IAEA COMPENDIUM DATABASE

To assess the accuracy in predicting the neutron spectrum, the R^2 metric was selected [28]. It was calculated for each sample and then averaged over the entire training or test subset, as formulated by:

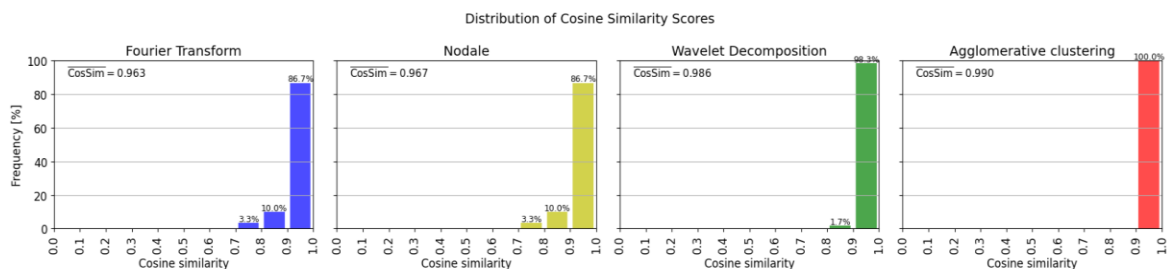
$$R_n^2 = \max\left\{0, \left(1 - \frac{\sum_{i=1}^m (t_{n,i} - \hat{t}_{n,i})}{\sum_{i=1}^m (t_{n,i} + \hat{t}_{n,i})}\right)\right\} \text{ and } \bar{R}^2 = \frac{1}{N} \sum_{n=1}^N R_n^2 \quad , \quad (14)$$

with $t_{n,i}$ and $\hat{t}_{n,i}$ the target and predicted fluence value respectively in the energy bin i for the spectrum n of the training or test subset and N is the total number of spectra analysed in the datasets. In addition to R^2 , other metrics have been analysed to better capture the similarity between predicted and simulated spectra. In particular, cosine similarity was evaluated as it directly compares the shape of two spectra independently of their absolute magnitude. This metric is well suited to spectral data as it highlights whether the two vectors point in the same direction. As illustrated in the following figure, the trends observed using cosine similarity closely mirrored those found with R^2 : samples with high R^2 also exhibited high cosine similarity, confirming the consistency of the two metrics in assessing prediction accuracy. Therefore, for the sake of clarity, interpretability, and consistency

302 with previous studies (e.g., the SiC case), R^2 is retained as the primary evaluation metric in the subsequent analysis of the collected
 303 results.



304
305 (a)

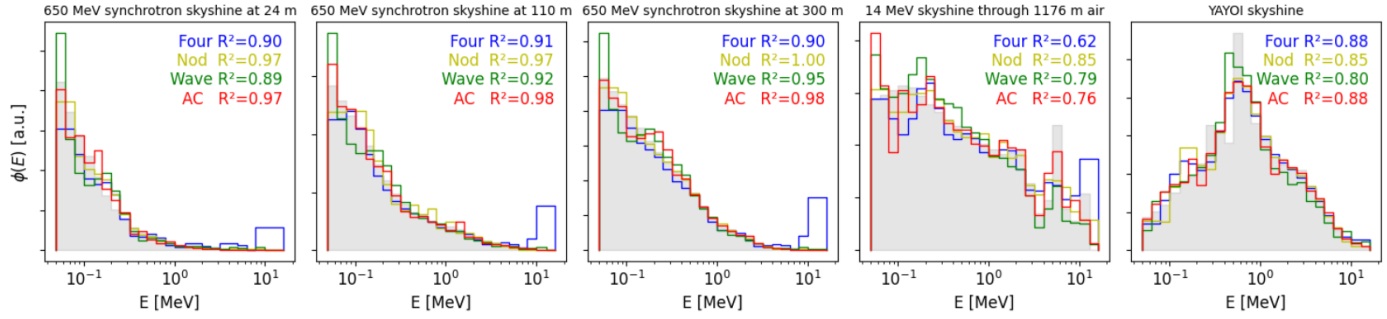


306
307 (b)

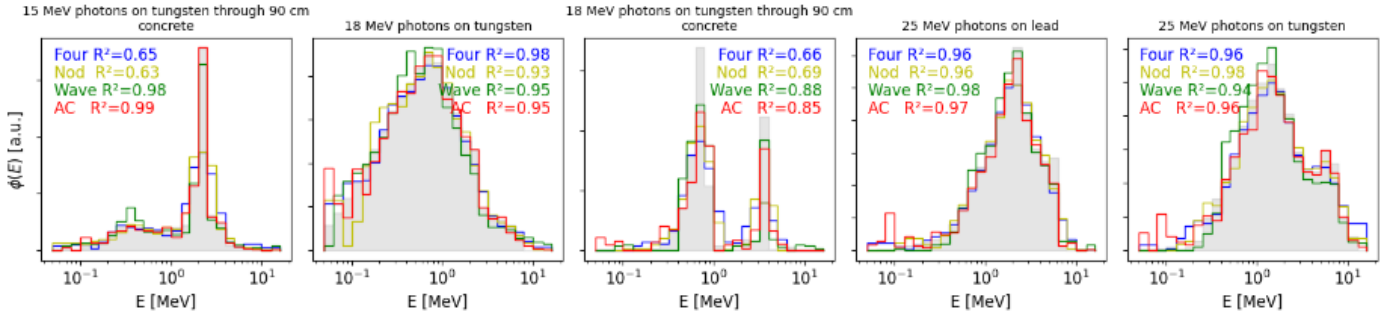
308 **Fig. 3.** R^2 score (a) and cosine similarity score (b) achieved for the test subset for the four described algorithms.

309 In Figure 3, the average R^2 obtained in predicting the neutron spectra from the test dataset is reported for the four feature extraction
 310 methods under study. An R^2 score greater than 0.9 guarantees a high-quality of neutron spectrum reconstruction. For the four
 311 algorithms, Fourier presents 70% of predictions in agreement with this criterion, Nodal 78%, Wavelet 84% and the Agglomerative
 312 clustering predicts in 88% of the cases a result with the desired R^2 value. Additionally, none of the spectra appear to be predicted
 313 with an R^2 lower than 0.4 for any of the four approaches. The relatively lower performance observed with the Fourier transform
 314 can be attributed to its intrinsic assumption of periodicity. This limitation arises from the fact that the Fourier series expansion is
 315 based on the principle of periodicity, which causes it to fail in representing the last energy bin if it does not align with the first one.
 316 These behaviours are visible in Figure 4. The Nodal basis approach yields intermediate performance: better than Fourier but inferior
 317 to Wavelet and Agglomerative clustering. This outcome is partially due to the use of a Linear Regression model, which was deemed
 318 appropriate in this case because the nodal projection significantly reduces the dimensionality of the input and tends to linearize the
 319 relationships between features and target variables. However, linear regression lacks the flexibility to model more complex,
 320 nonlinear dependencies, which can limit reconstruction accuracy when such interactions are present in the data. For both Wavelet
 321 decomposition and Agglomerative clustering, KRR was adopted as the learning model. These feature extraction methods preserve
 322 localized and potentially nonlinear structures in the input data, which justifies the use of KRR—capable of capturing such
 323 nonlinearities through the kernel trick. The kernel transforms the input space into a higher-dimensional feature space where linear
 324 regression can approximate complex relationships. This combination yields significantly better predictive performance,
 325 particularly in cases where the input-output relationship is highly nonlinear. Conversely, the Agglomerative clustering
 326 demonstrates poor predictive performance when attempting to reproduce a spectrum characterized by a 3-bump shape. This
 327 limitation arises from the absence of spectra with such a form in the training dataset. Despite this omission, the Fourier, Wavelet
 328 and Nodal models exhibit superior adaptability to representing the 3-bump spectra, achieving a higher average R^2 score compared
 329 to the Agglomerative clustering method for these cases.

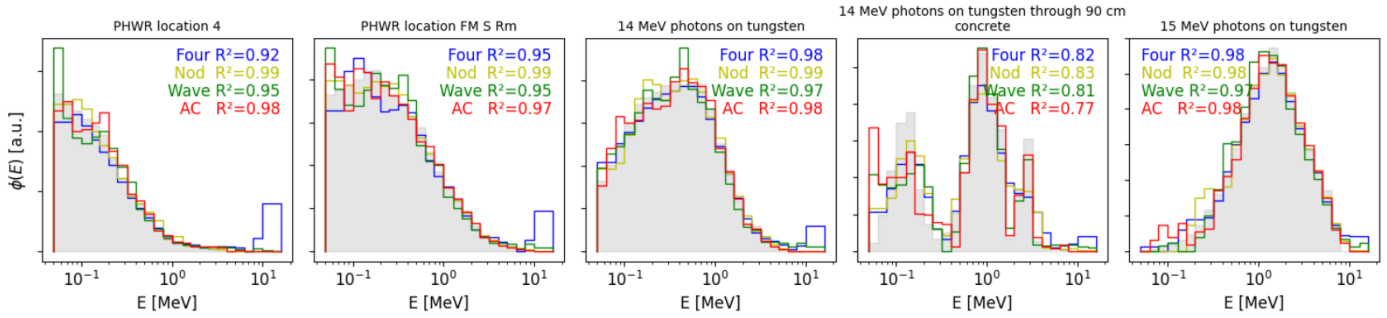
330



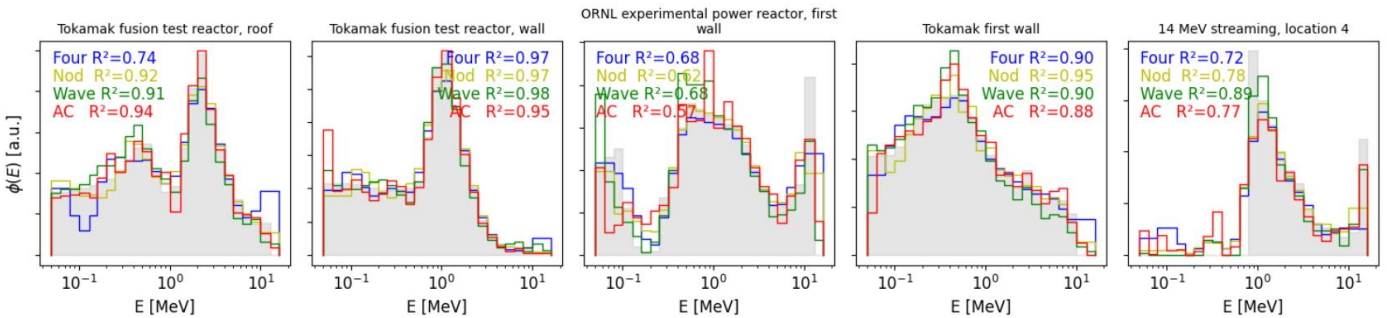
331



332



333

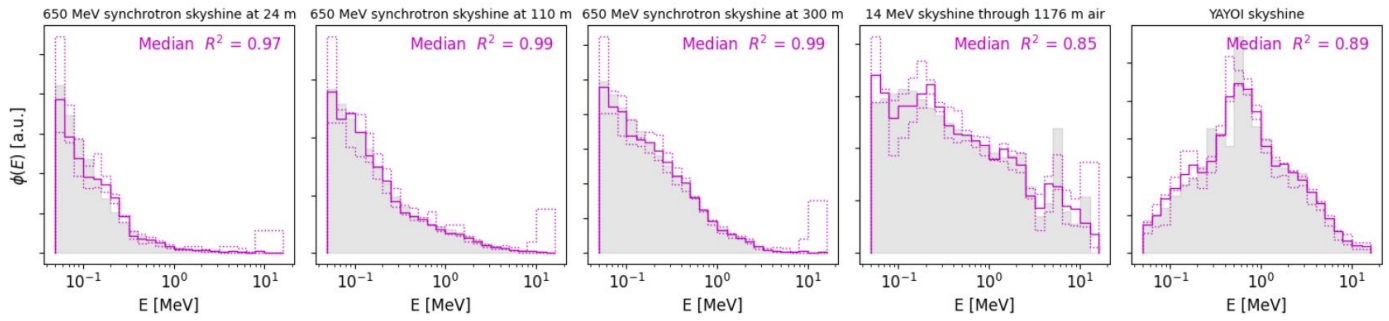


334

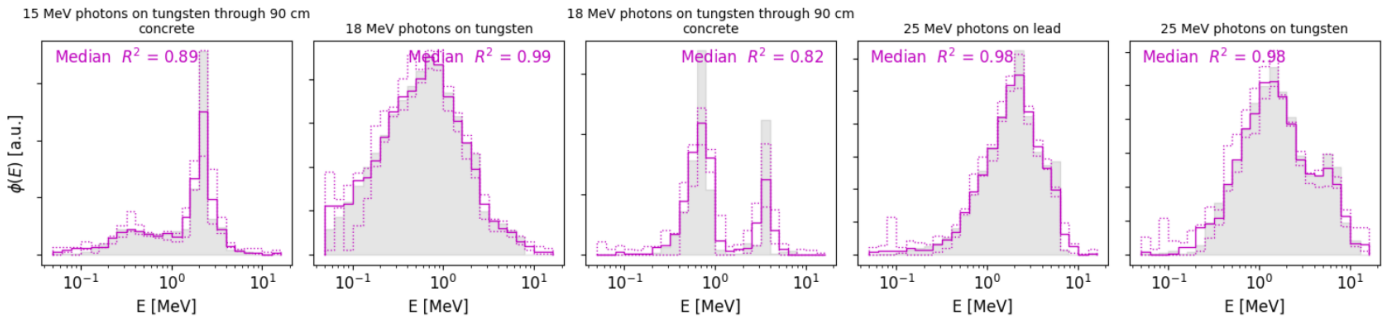
335 **Fig. 4.** Comparison of some real spectra taken from the IAEA Compendium database and predicted ones with the four described
 336 ML models (abbreviations; Four : Fourier Transform, Nod : Nodale; Wave : Wavelet Decomposition; AC : Agglomerative
 337 clustering).

338 Each model exhibits distinct strengths and limitations in reconstructing the neutron spectra, which opens the possibility of
 339 combining their predictions to enhance overall robustness and accuracy. One of the simplest methods is to calculate the median of
 340 predictions by bins. The advantage of the median is not to be penalized by extreme values as observed at the ends of the spectra
 341 predicted with the extracted features with the Fourier model.

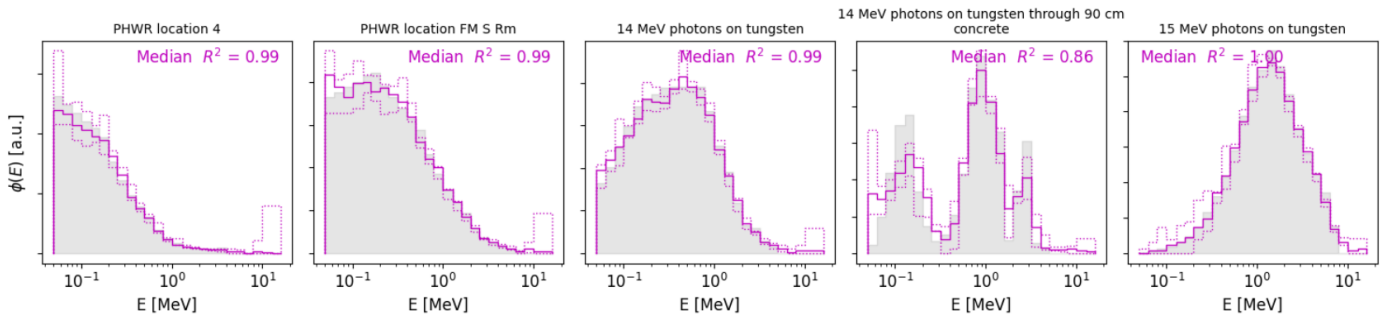
342



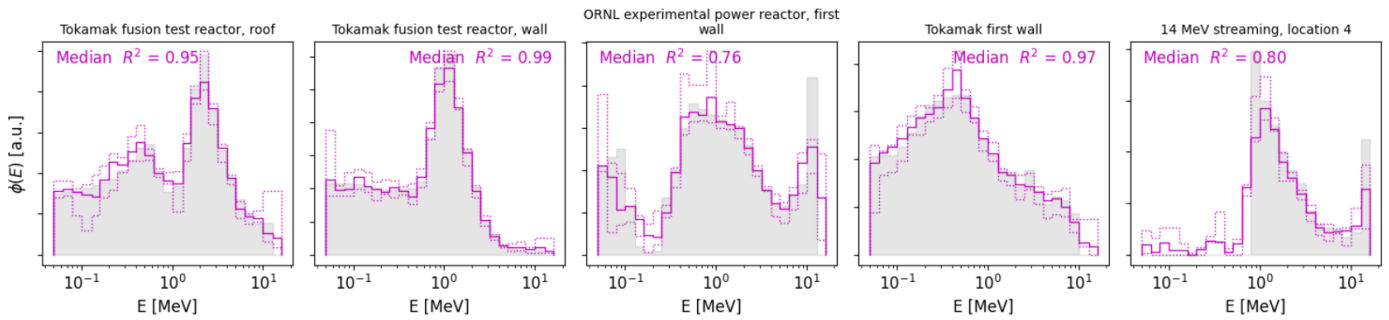
343



344



345



346

347

Fig. 5. Comparison of some real spectra taken from the IAEA Compendium database and predicted ones with the median prediction of the four described ML models.

348

349

350

351

352

353

354

355

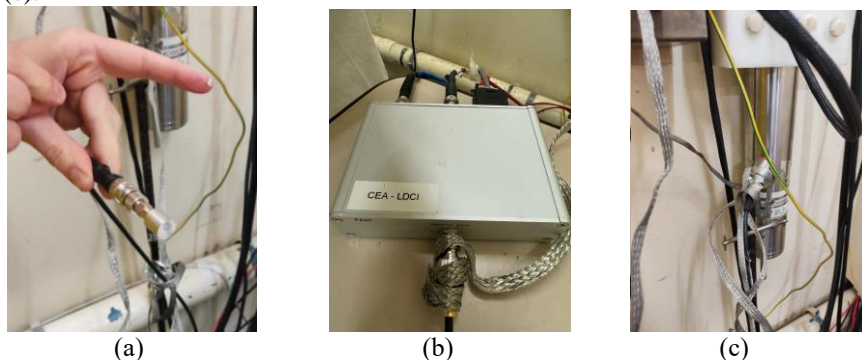
Using the median of the predictions of the four algorithms produces spectra that are mostly better predicted. On the other hand, these predictions are more robust with regard to the individual behavioural differences of each model. It therefore works like a meta-model whose combination function is the median instead of the mean. In Figure 5, the maximum and minimum spectra predicted by the four models for each energy bin are shown with dashed lines. The comparison between the spectra predicted by the median model and the actual spectra from the IAEA Compendium database yields an \bar{R}^2 of 0.953, which is slightly higher than that of the Kernel Ridge model with Agglomerative clustering. However, in some cases, the spectra predicted by the median model have a lower R^2 compared to those predicted by Kernel Ridge with Agglomerative clustering. This occurs because the median model, by construction, selects the median prediction from the four models for each energy bin and it does not always correspond

356 to the best possible prediction. If one of the individual models systematically produces more accurate spectra in a particular region,
 357 the median model does not take full advantage of this superior performance but instead blends predictions from all models.
 358 Moreover, if multiple models exhibit a similar systematic error—such as a consistent underestimation or overestimation in certain
 359 energy bins—the median model will inherit this bias. In contrast, an individual model that happens to be less affected by this bias
 360 may produce a higher R^2 for specific spectra. Therefore, it is advisable to compare the estimation of the performance with an
 361 appropriate metric of individual models before selecting the median one.

362 5.2 PREDICTION OF 14 MEV AND AMBe SOURCE SPECTRA

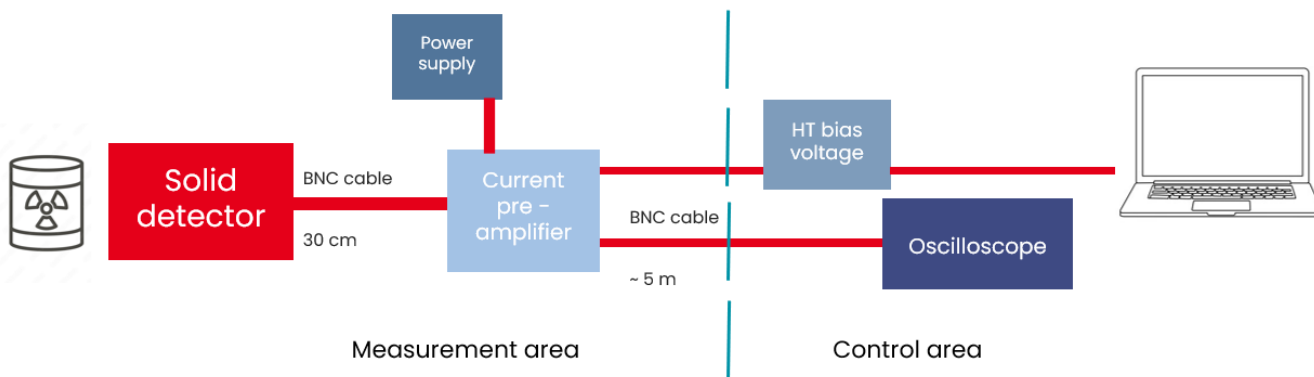
363 To evaluate the practical applicability of the proposed methodology, we tested the performance of the Median model on real
 364 measurements collected during an experimental campaign at the DANAIDES irradiation casemate within the TOTEM facility at
 365 CEA Cadarache [29]. In this facility, the CVD diamond detector was exposed to a 14.1 MeV neutron flux and AmBe neutron
 366 source. The D-T neutron generator GENIE16 model from SODERN of this facility was used for this experiment in continuous
 367 mode to obtain the 14.1 MeV, with the neutron emission equals to $2 \times 10^8 \text{ ns}^{-1}$ and the Am-Be source with neutron emission
 368 $2 \times 10^5 \text{ ns}^{-1}$.

369 The sensor was linked to a current amplifier, supplied by Metraware, ensuring optimized integration with the subsequent
 370 components of the counting system and it was located 1cm far from the neutron emission point. The bias voltage was applied by a
 371 CAEN bias voltage supply controlled via laptop and the detector was operated at -300 V. Finally, the signals were acquired through
 372 a Lecroy Waverunner 9254 oscilloscope. A photograph of the experimental apparatus is shown in Figure 6, which includes the
 373 CVD diamond detector encapsulated in a protective housing (a), the current preamplifier (b) and the positioning of the sensor close
 374 to the neutron generator (c).



375
 376
 377 **Fig. 6.** (a) CVD diamond detector encapsulated in a protective housing, (b) the current preamplifier and (c) the positioning of the
 378 sensor close to the neutron generator.
 379

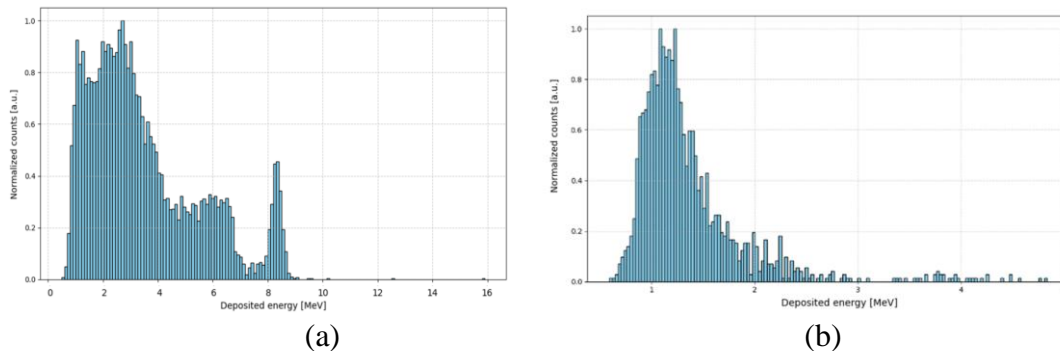
380 To better illustrate the overall measurement configuration, Figure 7 provides a schematic layout of the experimental setup,
 381 including the neutron source, detector positioning, electronics chain (amplifier, HV supply), and acquisition system.
 382



383
 384
 385 **Fig. 7.** Experimental setup scheme.
 386

387 The collected pulses with the measurements with the CVD diamond detector were processed in order to obtain the deposited energy
 388 in the detectors, in accordance to [30]. It is important to note that for an accurate prediction of the neutron spectrum using this
 389 machine learning approach, a high count rate is not required; rather, a precise characterization of the characteristic peaks in the
 390 energy deposition spectrum—as a function of the incident neutron energy—is essential. This is because the counts in the energy
 391 deposition spectrum are normalized to their maximum value before being used as input for the ML model. In this experimental
 392 campaign, it was observed that 10^4 acquisitions per source (corresponding to approximately 90 minutes for the 14.1 MeV source

393 and six hours for the AmBe source) were sufficient to reproduce the characteristic energy deposition peaks, shown in Figure 8.
 394 These experimental results are consistent with values reported in the literature. In particular, Figure 8a shows a peak at 8.3 MeV,
 395 typically used to assess the resolution of this sensor, arising from the $^{12}\text{C}(n,\alpha)^9\text{Be}$ reaction in which beryllium is produced in its
 396 ground state. The bump observed between 4 and 7 MeV is attributed to the production of three α particles, while the higher peak
 397 at lower energy is primarily due to scattering. For the AmBe source (Figure 8b), the emitted neutrons have lower energies and do
 398 not activate all threshold reactions with carbon atoms; hence, the main peak is predominantly due to scattering reactions, with a
 399 smaller contribution from alpha particle production, in agreement with the simulations.

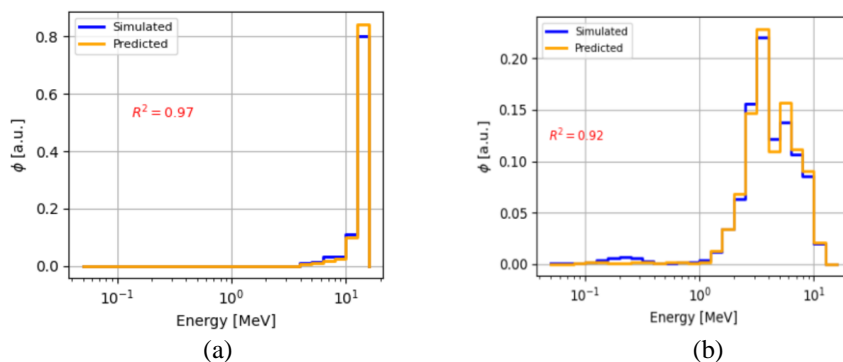


400
401

402 **Fig. 8.** Deposited energy in CVD diamond detector from 14.1 MeV (a) and AmBe neutron source (b).

403 Normalized deposited energy counts were used as the input for the aforementioned Median model, with the output being the
 404 neutron spectra shown in Figure 9. The predicted spectra were compared with reference spectra obtained from PHITS simulations
 405 of the DANAIDES irradiation casemate. The comparison reveals excellent agreement, with R^2 values exceeding 0.9 for both
 406 spectra. For these experiments, the energy grid discretization adopted for the training and test databases was not optimal, as the
 407 selected neutron sources have a much narrower energy range than the grid used. Consequently, adjusting the energy binning
 408 according to the specific source to which the detector is exposed could improve the energy discretization, while still respecting the
 409 intrinsic energy resolution of CVD diamond sensor [31]. Nonetheless, these results demonstrate a proof of concept regarding the
 410 feasibility of this neutron spectrum reconstruction methodology.
 411

411



412
413

414 **Fig. 9.** Predicted and simulated neutron energy spectra of 14.1 MeV (a) and AmBe source (b).

415 **6. CONCLUSION**

416 In this study, a comparative machine learning framework was developed to reconstruct neutron energy spectra from energy
 417 deposition data acquired in a CVD diamond detector. The input data were generated through detailed Monte Carlo simulations
 418 using the PHITS transport code, modelling the deposited energy in the detector to a set of 544 synthetic neutron spectra and 120
 419 real spectra taken from the IAEA Compendium database. Four distinct feature extraction methods were investigated to reduce the
 420 high dimensionality of the raw deposited energy data (600 bins) prior to regression: Fourier Transform, Nodal Basis Projection,
 421 Wavelet Decomposition, and Agglomerative Clustering approach. These transformations were combined with two regression
 422 models: Linear Regression was used for the Fourier and Nodal projections due to their compact and smooth representation of the
 423 signal, while Kernel Ridge Regression was applied to the Wavelet and Agglomerative features to better capture complex nonlinear
 424 relationships in the data. The performance of each model was evaluated using the coefficient of determination R^2 , revealing that
 425 all approaches consistently achieved R^2 scores above 0.7 across the test set. The Kernel Ridge model combined with Agglomerative

426 Clustering yielded the highest average R^2 , with 88% of the spectra reconstructed above the 0.9 threshold. Nevertheless, specific
427 limitations were observed: the Fourier-based model struggled with non-symmetric spectra due to its periodicity assumptions, while
428 the Kernel Ridge model exhibited reduced accuracy in reconstructing spectra with three-peaked ("3-bump") structures, which were
429 underrepresented in the training data. To enhance prediction robustness and generalization, a median-based ensemble model was
430 introduced by computing the bin-wise median across the four individual model outputs. This meta-model showed improved overall
431 R^2 performance ($\bar{R}^2 = 0.953$) and mitigated individual model shortcomings, such as boundary artefacts and localized overfitting.
432 However, the median approach does not always outperform the best individual model on specific spectra, particularly when one
433 model dominates in accuracy for a certain class of inputs. The proposed methodology was further tested through an experimental
434 campaign conducted at the TOTEM facility of CEA Cadarache, where the CVD diamond detector was irradiated with 14.1 MeV
435 and AmBe neutron sources. The reconstructed spectra showed good agreement with reference spectra generated via PHITS
436 simulations under the same experimental conditions, demonstrating the feasibility of applying the method in real-world scenarios.
437 This study lays a proof of concept for applying machine learning to neutron spectroscopy using solid-state detectors. It offers an
438 alternative to classical unfolding methods, avoiding common pitfalls such as instability in matrix inversion and artefacts due to
439 over-regularization. Future developments may include the expansion of the training dataset to cover a wider variety of spectral
440 shapes, refinement of the energy grid discretization, and the integration of advanced learning models such as neural networks or
441 probabilistic method like Gaussian processes, as well as the introduction of the uncertainty quantification. An additional key
442 perspective concerns the incorporation of the total neutron fluence in the reconstruction pipeline: currently, spectra are normalized
443 to unit area, but estimating the absolute fluence would enable a physically meaningful reweighting of the reconstructed spectrum,
444 providing not only spectral shape but also absolute flux values. Moreover, extending the methodology to other detector materials
445 and geometries could broaden its applicability across a wider range of neutron energy regimes.

446 ACKNOWLEDGMENTS

447 The authors would like to acknowledge the organizers and participants of the SEME (Semaines d'Études Mathématiques et
448 Entreprises) event for their valuable contributions and insightful discussions. Special thanks to AMIES (Agence pour les
449 Mathématiques en Interaction avec l'Entreprise et la Société) and Archimedes Institute from Aix-Marseille University for
450 facilitating this initiative.

451 REFERENCES

- 452 [1] Q. Huang, Z. Chen, X. Xu, and J. Zhang, "A review of the application of artificial intelligence to nuclear reactors: Where we
453 are and what's next," *Heliyon*, vol. 9, no. 3, p. e14273, 2023. DOI: 10.1016/j.heliyon.2023.e14273
- 454 [2] A. Lyoussi et al., "I_SMART a collaborative project on innovative sensor for material ageing and radiation testing," KIC
455 InnoEnergy, Grenoble, France, 2012. [Online]. Available: https://www.kic-innoenergy.com/projects/i_smart
- 456 [3] M. Birk, G. Goldring, and P. Hillman, "Fast neutron spectroscopy with solid state detectors," *Nucl. Instrum. Methods*, vol.
457 21, pp. 107–113, 1963. DOI: 10.1016/0029-554X(63)90061-6
- 458 [4] J. Gómez-Ros, "Solid state detectors for neutron radiation monitoring in fusion facilities," *Radiation Measurements*, vol. 71,
459 pp. 57–63, 2014. DOI: 10.1016/j.radmeas.2014.04.006
- 460 [5] E. Belfiore et al., "SiC Detector Thickness Optimization for Enhanced Variability Response," in *SNA + MC 2024 Conf.*
461 *Proc.*, Paris, 2024, pp. 45–50. (in press)
- 462 [6] E. Belfiore et al., "Proof of Concept of the Prediction of Neutron Spectra using Solid Detectors via Machine Learning
463 Approach," in *IEEE NSS MIC RTSD Conf.*, Tampa, 2024. DOI: 10.1109/NSSMICRTSD49039.2024.10001
- 464 [7] R. Antoni, P.-G. Allineï, and L. Bourgois, "Prediction of fast neutron spectra with a spherical TEPC using a machine-learning
465 algorithm," *Nucl. Instrum. Methods Phys. Res. A*, vol. 1050, p. 168139, 2023. DOI: 10.1016/j.nima.2023.168139
- 466 [8] E. Belfiore et al., "Theoretical analysis of neutron spectra measurement with SiC detectors using a machine learning
467 technique," *J. Instrum.*, vol. 19, no. 1, p. P01007, 2024. DOI: 10.1088/1748-0221/19/01/P01007
- 468 [9] Y. Liu et al., "Simultaneous measurement of energy spectrum and fluence of neutrons using a diamond detector," *Sci. Rep.*,
469 vol. 12, no. 1, p. 15325, 2022. DOI: 10.1038/s41598-022-19711-7
- 470 [10] E. Griesmayer et al., "Neutron cross section measurements with diamond detectors," *EPJ Web Conf.*, vol. 225, p. 05005,
471 2019. DOI: 10.1051/epjconf/201922505005
- 472 [11] R. Slavickas et al., "Applications of carbon-based diamond detectors: A critical review," *Sensors*, vol. 22, no. 15, p. 5736,
473 2022. DOI: 10.3390/s22155736
- 474 [12] O. Obraztsova et al., "Comparing the response of a SiC and a sCVD diamond detectors to 14-MeV neutron radiation," *IEEE*
475 *Trans. Nucl. Sci.*, vol. 65, no. 3, pp. 578–585, 2018. DOI: 10.1109/TNS.2018.2800440
- 476 [13] C. Fabjan, "Collider detectors for multi-TeV particles," in *Encyclopedia of Physical Science and Technology*, 3rd ed., vol.
477 3, pp. 253–268, 2003.
- 478 [14] C. E. Nebel, "CVD diamond: a review on options and reality," *Materials*, vol. 14, no. 22, p. 7081, 2021. DOI:
479 10.3390/ma14227081
- 480 [15] Z. Cheng et al., "High thermal conductivity in wafer scale cubic silicon carbide crystals," *arXiv preprint*, arXiv:2207.05292,
481 2022.

482 [16] Q. Wahab and M. Willander, “Silicon carbide and diamond for high temperature device applications,” Academia.edu, 2015.
483 [Online]. Available: <https://www.academia.edu/20207274>

484 [17] J. Achard, V. Jacques, and A. Tallaire, “CVD diamond single crystals with NV centres: a review of material synthesis and
485 technology for quantum sensing applications,” arXiv preprint, arXiv:1912.09749, 2019.

486 [18] Z. He et al., “Machine learning in nuclear physics at low and intermediate energies,” *Front. Phys.*, vol. 11, p. 1152892,
487 2023. DOI: 10.3389/fphy.2023.1152892

488 [19] M. Kamuda et al., “A comparison of machine learning methods for automated gamma-ray spectroscopy,” *J. Radioanal.
489 Nucl. Chem.*, vol. 324, pp. 1115–1124, 2020. DOI: 10.1007/s10967-020-07125-3

490 [20] T. Sato et al., “Features of Particle and Heavy Ion Transport code System (PHITS) version 3.02,” *J. Nucl. Sci. Technol.*, vol.
491 55, no. 6, pp. 684–690, 2018. DOI: 10.1080/00223131.2018.1432933

492 [21] L. Sobczak, “Dimensionnement d’un détecteur diamant pour la spectrométrie de neutrons rapides à partir d’une mesure
493 micro-dosimétrique,” Internal Technical Report, CEA Cadarache, 2022. (Available upon request)

494 [22] J. Kelley et al., “Energy levels of light nuclei A = 12,” *Nucl. Phys. A*, vol. 968, pp. 71–142, 2017. DOI:
495 10.1016/j.nuclphysa.2017.04.002

496 [23] K. Shibata et al., “JENDL-4.0: A new library for nuclear science and engineering,” *J. Nucl. Sci. Technol.*, vol. 48, no. 1, pp.
497 1–30, 2011. DOI: 10.1080/00223131.2011.555445

498 [24] R. V. Griffith et al., “Compendium of neutron spectra and detector responses for radiation protection purposes,” IAEA
499 Technical Report Series No. 403, 1990.

500 [25] F. Pedregosa et al., “Scikit-learn: Machine learning in Python,” *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011. DOI:
501 10.5555/1953048.207819

502 [26] C. K. Chui, *An Introduction to Wavelets*, San Diego: Academic Press, 1992, pp. 1–272. ISBN: 9780121745841

503 [27] P. Exterkate, “Model selection in kernel ridge regression,” *Comput. Stat. Data Anal.*, vol. 68, pp. 1–16, 2013. DOI:
504 10.1016/j.csda.2013.03.005

505 [28] D. Britto, J. A. Silva, and E. Rocha, “What is R² all about?,” *Leviathan – Cadernos de Pesquisa Política*, vol. 3, no. 1, pp.
506 60–68, 2011.

507 [29] J. Roult dit Rouaux et al., “Experimental gamma coincidence spectra recorded in prompt gamma neutron activation
508 analysis,” *J. Radioanal. Nucl. Chem.*, vol. 333, pp. 6577–6592, 2024. DOI: 10.1007/s10967-024-09822-x

509 [30] Q. Potiron et al., “Modelling of a SiC based detector for the interpretation of 14.1 MeV neutrons measurements,” *EPJ Web
510 Conf.*, vol. 288, p. 04011, 2023. DOI: 10.1051/epjconf/202328804011

511 [31] C. Weiss and E. Griesmayer, “Fusion neutron diagnostics with CVD diamond detectors,” *Fusion Eng. Des.*, vol. 203, p.
512 113605, 2024. DOI: 10.1016/j.fusengdes.2023.113605