

2 **A density-based clustering algorithm for the  
3 CYGNO data analysis**

---

4 **E. Baracchini,<sup>a,b</sup> L. Benussi,<sup>c</sup> S. Bianco,<sup>c</sup> C. Capoccia,<sup>c</sup> M. Caponero,<sup>c,d</sup> G. Cavoto,<sup>e,f</sup> A.  
5 Cortez,<sup>a,b</sup> I. A. Costa,<sup>g</sup> E. Di Marco,<sup>e</sup> G. D'Imperio,<sup>e</sup> G. Dho,<sup>a,b</sup> F. Iacoangeli,<sup>e</sup> G.  
6 Maccarrone,<sup>c</sup> M. Marafini,<sup>e,h</sup> G. Mazzitelli,<sup>c</sup> A. Messina,<sup>e,f</sup> R. A. Nobrega,<sup>g</sup> A. Orlandi,<sup>c</sup> E.  
7 Paoletti,<sup>c</sup> L. Passamonti,<sup>c</sup> F. Petrucci,<sup>i,j</sup> D. Piccolo,<sup>c</sup> D. Pierluigi,<sup>c</sup> D. Pinci<sup>e</sup> F. Renga,<sup>e</sup> F.  
8 Rosatelli,<sup>c</sup> A. Russo,<sup>c</sup> G. Saviano,<sup>c,k</sup> and S. Tomassini<sup>c</sup>**

9 *<sup>a</sup>Gran Sasso Science Institute,  
10 L'Aquila, I-67100, Italy*

11 *<sup>b</sup>Istituto Nazionale di Fisica Nucleare,  
12 Laboratori Nazionali del Gran Sasso, Assergi, Italy*

13 *<sup>c</sup>Istituto Nazionale di Fisica Nucleare,  
14 Laboratori Nazionali di Frascati, I-00044, Italy*

15 *<sup>d</sup>ENEA Centro Ricerche Frascati, Frascati, Italy*

16 *<sup>e</sup>Istituto Nazionale di Fisica Nucleare,  
17 Sezione di Roma, I-00185, Italy*

18 *<sup>f</sup>Dipartimento di Fisica Sapienza Università di Roma, I-00185, Italy*

19 *<sup>g</sup>Universidade Federal de Juiz de Fora, Juiz de Fora, Brasil*

20 *<sup>h</sup>Museo Storico della Fisica e Centro Studi e Ricerche "Enrico Fermi",  
21 Piazza del Viminale 1, Roma, I-00184, Italy*

22 *<sup>i</sup>Dipartimento di Matematica e Fisica, Università Roma TRE, Roma, Italy*

23 *<sup>j</sup>Istituto Nazionale di Fisica Nucleare, Sezione di Roma TRE, Roma, Italy*

24 *<sup>k</sup>Dipartimento di Ingegneria Chimica, Materiali e Ambiente, Sapienza Università di Roma, Roma, Italy*

25 *E-mail: [igor.abritta@engenharia.ufjf.br](mailto:igor.abritta@engenharia.ufjf.br)*

26 **ABSTRACT:** Time Projection Chambers (TPCs) working in combination with Gas Electron Multi-  
27 pliers (GEMs) produce a very sensitive detector capable of observing low energy events. This is  
28 achieved by capturing photons generated during the GEM electron multiplication process by means  
29 of a high-resolution camera. The CYGNO experiment has recently developed a TPC Triple GEM  
30 detector coupled to a low noise and high spatial resolution CMOS sensor. For the image analysis,  
31 an algorithm based on an adapted version of the well-known DBSCAN was implemented, called  
32 iDBSCAN. In this paper a description of the iDBSCAN algorithm will be given, including test  
33 and validation of its parameters, and a comparison with a widely used algorithm known as Nearest  
34 Neighbor Clustering (NNC). The results will show that the adapted version of DBSCAN is capable  
35 of providing full signal detection efficiency and very good energy resolution while improving the  
36 detector background rejection.

---

37	<b>Contents</b>	
38	<b>1 Experimental setup</b>	<b>2</b>
39	1.1 LEMON detector	2
40	1.2 Acquisition runs	2
41	1.3 Detector expected signals	3
42	<b>2 Data analysis flow</b>	<b>3</b>
43	2.1 Data structure	3
44	2.2 Overview of the event reconstruction procedure	4
45	2.3 The CYGNO intensity-based clustering algorithm	5
46	2.3.1 iDBSCAN	5
47	2.3.2 Validation of the iDBSCAN parameters	6
48	<b>3 iDBSCAN and NNC comparison</b>	<b>7</b>
49	3.1 Electronic noise, natural radioactivity and $^{55}\text{Fe}$ energy spectra	7
50	3.2 Slimness selection optimization	9
51	3.3 Light Yield Resolution	11
52	<b>4 Summary</b>	<b>13</b>

---

## 53 **Introduction**

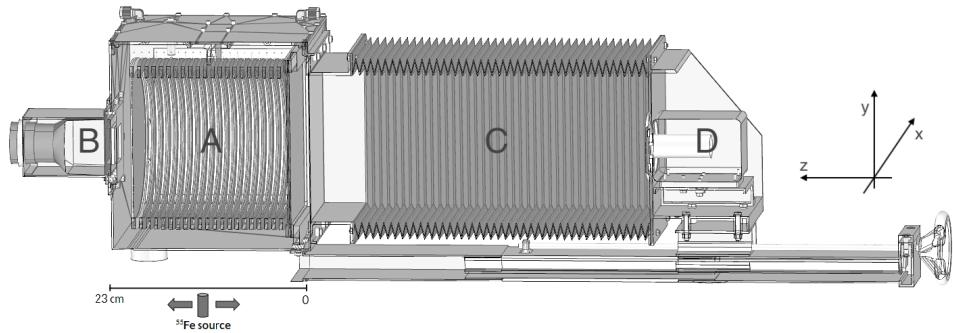
54 Clustering analysis is a widely used unsupervised technique to organize datasets into groups based  
55 on their similarities. One of the most known algorithms is the so-called Density-Based Spatial  
56 Clustering of Applications with Noise (DBSCAN) [1]. Given a set of elements distributed over a  
57 hyper-plane, DBSCAN seeks for areas of high density to form clusters. Such density is calculated  
58 considering the number of elements within a pre-defined hyper-sphere. The generalization power  
59 of DBSCAN and its simplicity, which make it a very attractive algorithm, can be understood in  
60 terms of its two parameters: the radius of the hyper-sphere ( $\epsilon$ ), which is applied over each element  
61 to count the number of neighboring elements around it, and the minimum number of points inside  
62 each hyper-sphere ( $N_{min}$ ), used to decide if those elements should make up a cluster. To fulfill  
63 the needs of the CYGNO experiment, a detector-specific algorithm, based on DBSCAN, has been  
64 developed. Within the context of the experiment, a detection apparatus composed of an optical  
65 readout system based on a high-resolution and low noise CMOS sensor capable of providing track  
66 images produced by interacting particles with release energies in the range of a few keV has been  
67 developed [2–7]. This modified version of DBSCAN, called intensity-based DBSCAN or simply  
68 iDBSCAN, has shown to be able to improve detector performance when compared to the previously  
69 used algorithm based on the Nearest Neighbor Clustering (NNC) technique [8]. This paper proposes  
70 a comparative study on the impact of NNC and iDBSCAN on two crucial detector's parameters,

71 background rejection and energy resolution, measured in the energy range of a few keV. For such,  
72 low energy particles (5.9 keV photons) produced by a  $^{55}\text{Fe}$  radioactive source, background from  
73 natural radioactivity and data with electronics noise only were employed.

## 74 1 Experimental setup

### 75 1.1 LEMON detector

76 LEMOn (Large Elliptical MODule) is the most recent CYGNO experiment's prototype. Its core  
77 consists of a 7 liter active drift volume surrounded by an elliptical field cage ( $20 \times 20 \times 24 \text{ cm}^3$ )  
78 and a  $20 \times 24 \text{ cm}^2$  Triple GEM structure whose produced photons are readout by an Orca Flash 4  
79 CMOS-based camera [9] placed at a distance of 52.5 cm (i.e. 21 Focal Length, FL). More details  
80 are given in Ref. [8, 10, 11]. The drift chamber was filled with a He/CF<sub>4</sub> gas mixture in the  
81 proportion of 60/40 and a  $^{55}\text{Fe}$  source with an activity of about 740 MBq was used. For operation,  
82 electric fields are applied to the TPC drift volume and between the GEMs. They are called drift  
83 field ( $E_d$ ) and transfer field ( $E_t$ ) respectively. The typical operating conditions of the detector, as  
84 used in this work, are:  $E_d = 500 \text{ V/cm}$ ,  $E_t = 2.5 \text{ kV/cm}$ , and a voltage difference across the GEM  
85 sides ( $V_{\text{GEM}}$ ) of 460 V.



**Figure 1.** Drawing of the experimental setup. In particular, the elliptical field cage close on one side by the triple-GEM structure and on the other side by the semitransparent cathode (A), the PMT (B), the adaptable bellow (C) and the CMOS camera with its lens (D) are visible. The sliding external  $^{55}\text{Fe}$  source, positioned close to the TPC is also drawn.

### 86 1.2 Acquisition runs

87 Data were acquired using auto-trigger mode. For the proposed study presented in this document,  
88 three different acquisition datasets were used, as listed below:

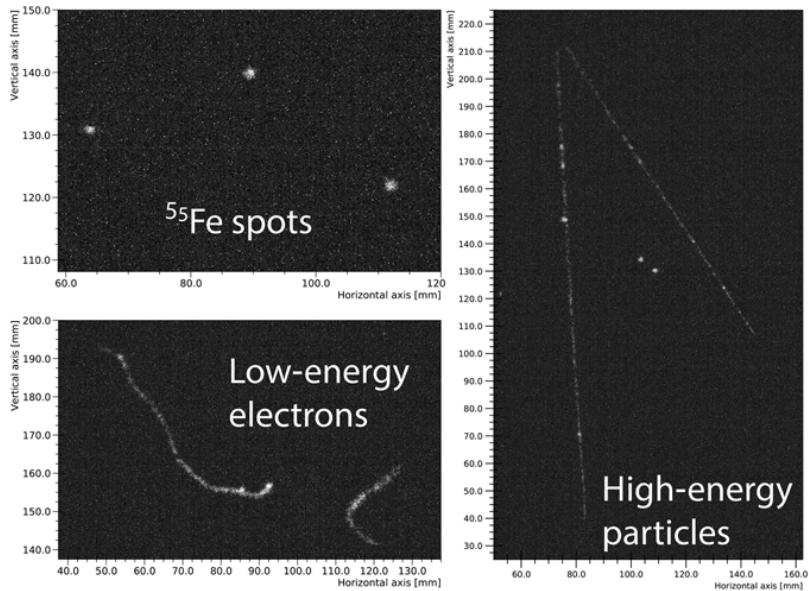
- 89 • Electronic noise (EN) dataset: produced by lowering down  $V_{\text{GEM}}$  to a value where the  
90 multiplication process is forbidden (6478 images recorded);
- 91 • Natural radioactivity (NRAD) dataset (composed of cosmic rays and environmental radioac-  
92 tivity): produced by exposing the camera lens and turning on the detector power supplies

93 and raising  $V_{\text{GEM}}$  to the nominal value of 460 V to allow charge multiplication and secondary  
 94 light emission during this process (864 images recorded);

- 95 • Electron Recoils (ER) dataset: the same as the previous item but placing a  $^{55}\text{Fe}$  source near  
 96 to the detector drift volume (864 images recorded).

### 97 1.3 Detector expected signals

98 Based on the acquisition datasets defined in section 1.2, particles interacting with the detector gas  
 99 can have two distinct origins:  $^{55}\text{Fe}$  source and natural radioactivity. The former releases 5.9 keV  
 100 photons which produce round spots on the image while the latter can be composed of few different  
 101 particles as photons, electrons and muons. Typical signals are shown in Fig. 2: three interactions  
 102 of  $^{55}\text{Fe}$  photons in the left top image; two low-energy electrons in the left bottom image; and two  
 103 high-energy particles and, between them, two interactions of  $^{55}\text{Fe}$  photons in the right image.



104 **Figure 2.** Examples of signals that can occur using the described configuration.

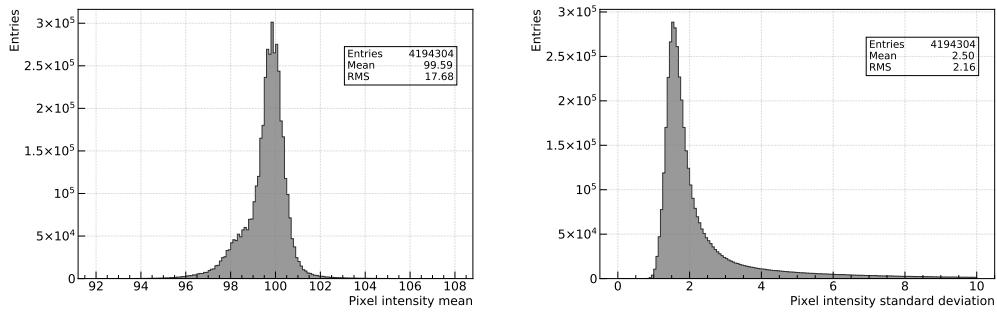
105 In this work, the signals of interest are those generated by the 5.9 keV photons, which are used  
 106 to assess the impact of the proposed clustering algorithms on the detector characteristics, focusing  
 107 mainly on its energy resolution and background-events rejection performance in the energy range  
 of few keV.

## 108 2 Data analysis flow

### 109 2.1 Data structure

110 The acquisition system provides images with  $2048 \times 2048$  pixels captured by the ORCAD Flash  
 111 CMOS sensor. The photo sensor has an sensitive area of  $13312 \mu\text{m}^2$  and each pixel has a size of  
 112  $6.5\mu\text{m} \times 6.5\mu\text{m}$ . The camera's exposure time was set to 40 ms and it covers an area of  $26 \times 26 \text{ cm}^2$

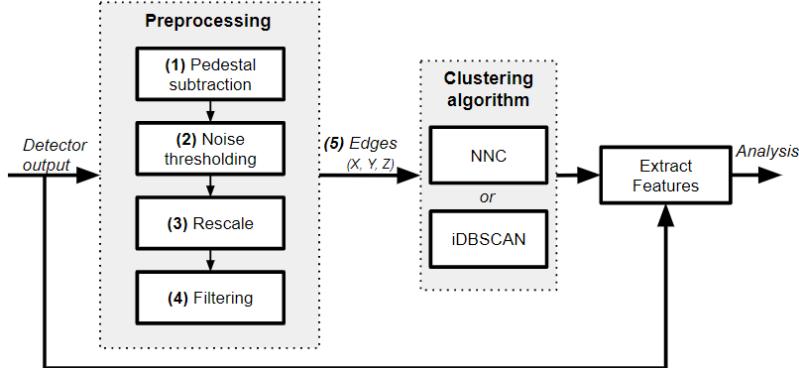
113 in relation to the plane of the last layer of the GEM detector. Each pixel provides a response, here  
 114 called intensity, proportional to the number of collected photons [6] added to a baseline, also known  
 115 as pedestal, which can be defined as the intensity value corresponding to zero photons. Specifically,  
 116 the pedestal average value of the sensor is about 99 counts, however it can vary from pixel to pixel.  
 117 Additionally, the noise level is another important parameter that can vary from pixel to pixel. Those  
 118 effects can be seen in Fig. 3 which shows the mean and standard deviation distributions of the  
 119 noise as computed for each pixel, produced with the EN dataset. To account for such variations,  
 120 both the pixel baseline ( $\mu_i$ ) and its average noise ( $\sigma_i$ ), calculated as the standard deviation of the  
 121 pedestal distribution, are estimated for every single pixel  $i$  before running the event reconstruction  
 122 procedure.



**Figure 3.** Mean and standard deviation distributions of the sensor's pixels noise.

## 123 2.2 Overview of the event reconstruction procedure

124 The current CYGNO's event-reconstruction algorithm is represented in the flowchart shown in  
 125 Fig. 4 and it is described below.



**Figure 4.** Flowchart of the CYGNO's event-reconstruction algorithm.

- 126 1. Pedestal subtraction is carried out pixel by pixel by subtracting  $\mu_i$  from their original intensity  
 127 values, generating new intensity values defined as  $I_i$ .

- 128 2. Lower and upper thresholds are applied to  $I_i$ . While the upper limit is set to 100 counts, the  
 129 lower limit is set to 1.3 times  $\sigma_i$ . The upper limit allows to remove pixels with a too large  
 130 intensity, very likely not due to ionization in gas, while the lower limit was optimized and  
 131 set to be just above noise level to ensure a good detection efficiency, but not too low in order  
 132 not to overload the event-reconstruction algorithm with pixels dominated by noise. Pixels  
 133 outside those limits have their intensities reset to zero.
- 134 3. Images are then rescaled to  $512 \times 512$  pixels, for CPU reasons, so that each  $4 \times 4$  matrix, called  
 135 macro-pixel, is assigned an intensity value corresponding to the average of the intensities  $I_i$   
 136 of the 16 pixels occupying the same area of the sensor.
- 137 4. The rescaled image goes then through a filtering stage based on a  $4 \times 4$  median filter that  
 138 replaces a given macro-pixel intensity by the median of all macro-pixels in its neighborhood  
 139  $w$ , as given by Equation 2.1 [12], where  $f(x, y)$  is the intensity of the macro-pixel  $(x, y)$ .

$$g(x, y) = \text{median}\{f(x, y), (x, y) \in w\} \quad (2.1)$$

140 Such filter is widely used in many applications due to its effective noise suppression capability  
 141 and high computational efficiency [13]. Tests performed on the EN dataset (see section 1.2)  
 142 showed that this filter is able to reduce the number of noise pixels sent to the clustering  
 143 algorithm by a factor of  $3.07 \pm 0.02$ .

- 144 5. Finally, the coordinates (X, Y) and respective intensities (Z) of the pixels with non-zero  $I_i$   
 145 values are sent to the clustering algorithm whose output is used to extract clusters' features  
 146 such as integrated light, length and width, computed over the full-resolution image. Those  
 147 features are then used to select events of interest.

148 In this work three features, extracted from the clusters, are used:

- 149 • Length and width: the full length of the major and minor axes along the two eigenvectors of  
 150 the (X,Y) pixel matrix in the context of Principal Component Analysis [14] are assigned as  
 151 the length and width of the cluster, respectively.
- 152 • Cluster light: calculated as the sum of all the pixel  $I_i$  intensities belonging to the cluster.

153 As mentioned before, prior to iDBSCAN, the CYGNO clustering algorithm was based on the  
 154 widely employed NNC method. Basically it groups neighboring pixels that went through a selection  
 155 similar to the one in step 3. A detector performance study using such method was presented in [8].

## 156 2.3 The CYGNO intensity-based clustering algorithm

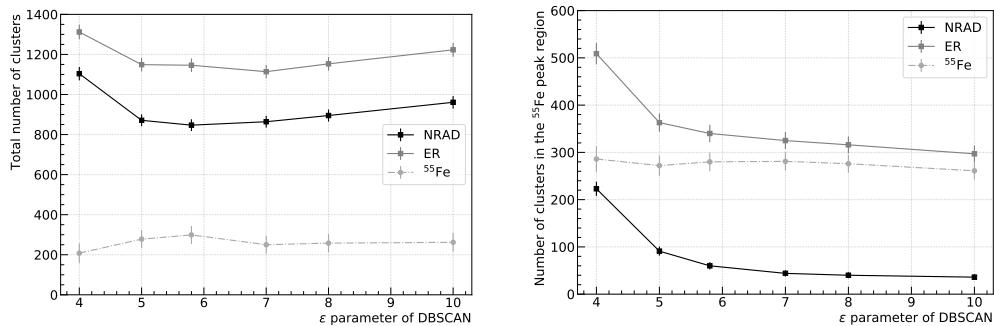
### 157 2.3.1 iDBSCAN

158 As in many areas, in particle physics it is possible to insert a priori knowledge about the detection  
 159 system and its data to improve the performance of the clustering task [15]. In this sense, a modifi-  
 160 cation of DBSCAN [16] clustering algorithm was implemented, to better match the experimental

161 conditions and data of the LEMOn detector. As mentioned before, DBSCAN has only two parameters:  
 162  $\epsilon$  and  $N_{min}$ . Whenever the number of neighboring elements inside a hyper-sphere reaches the  
 163  $N_{min}$  value, the center element and all its neighbors are activated to start the formation of a cluster.  
 164 Then, the same process happens to all the neighboring elements in order to expand the starting  
 165 cluster, to form a final cluster. This process is repeated to all the data elements. To be applied to  
 166 CYGNO, instead of using the number of elements as a parameter to decide if the elements inside  
 167 a hyper-sphere make part of a cluster, the sum of their intensity values is used. Consequently,  
 168 the  $N_{min}$  become a parameter related to the total intensity within a hyper-sphere instead of to the  
 169 number of elements. Therefore, rather than having each pixel counted as a unit when computing  
 170 the number of pixels inside a given hyper-sphere, each pixel counts  $I_i$  times. If the total intensity is  
 171 equal or greater than a certain value ( $N_{min}$ ), they are considered as making part of a cluster.

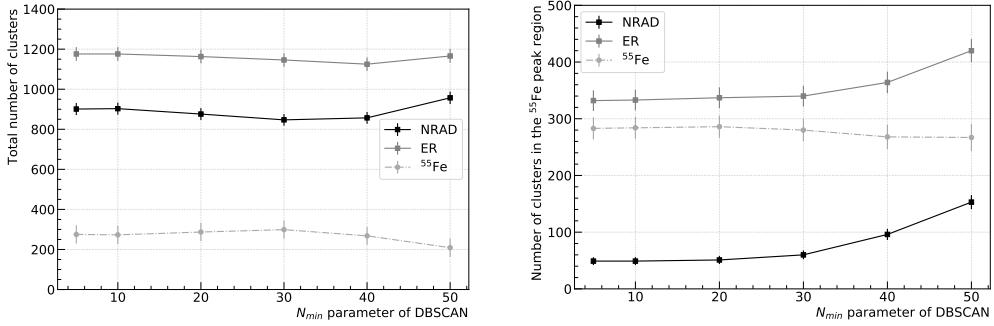
### 172 2.3.2 Validation of the iDBSCAN parameters

173 The CYGNO Collaboration is currently using iDBSCAN for the clustering method in its event-  
 174 reconstruction. The iDBSCAN performance for signals produced by the interactions of photons  
 175 from  $^{55}\text{Fe}$  has been studied as a function of different values of the its parameters:  $\epsilon$  and  $N_{min}$ . In  
 176 order to evaluate those values, a test on the detector efficiency and background rejection was carried  
 177 out: a scan over the two iDBSCAN parameters has been performed. While the  $\epsilon$  ( $N_{min}$ ) parameter  
 178 will be fixed to a value of 5.8 (30), the other parameter's value will be swept from 5 to 50 (4 to 10).  
 179 Figure 5 (left) shows the total number of clusters found as a function of  $\epsilon$  for two distinct datasets:  
 180 ER and NRAD. For low  $\epsilon$  values the number of NRAD clusters tends to increase, indicating an  
 181 increase of background contamination. However, for  $\epsilon$  values between 5 and 7, this contamination  
 182 rate stabilizes around a minimum value. Figure 5 (right) shows the same trend, while counting  
 183 only clusters with an integral in the range 2000–4000 photons, characteristic of  $^{55}\text{Fe}$  deposits. This  
 184 region refers to the energy region of the  $^{55}\text{Fe}$  produced electron recoils (see Fig. 9).



**Figure 5.** Total number of reconstructed clusters (left) and Spot number (right) as a function of  $\epsilon$  for runs with NRAD events only and with  $^{55}\text{Fe}$  + NRAD events.

185 Similarly, a scan over the  $N_{min}$  parameter has been performed as shown in Fig. 6. Applying  
 186 the same logic as for the  $\epsilon$  parameter, the plot on the left indicates a low contamination region for  
 187  $N_{min}$  values between 20 and 40, and the right plot to a region for  $N_{min} \leq 30$ . In both cases, when  
 188 stable, the difference between the results indicate a number of  $^{55}\text{Fe}$  clusters of about 280.



**Figure 6.** Total number of reconstructed clusters (left) and Spot number (right) as a function of  $N_{min}$  for runs with NRAD events only and with  $^{55}\text{Fe} + \text{NRAD}$  events.

Finally, energy resolution has also been measured in function of the iDBSCAN parameters. Values around 12.2% have been measured for all the  $\epsilon$  and  $N_{min}$  considered values, with negligible variation. Section 3.3 provides details about the energy resolution measurement.

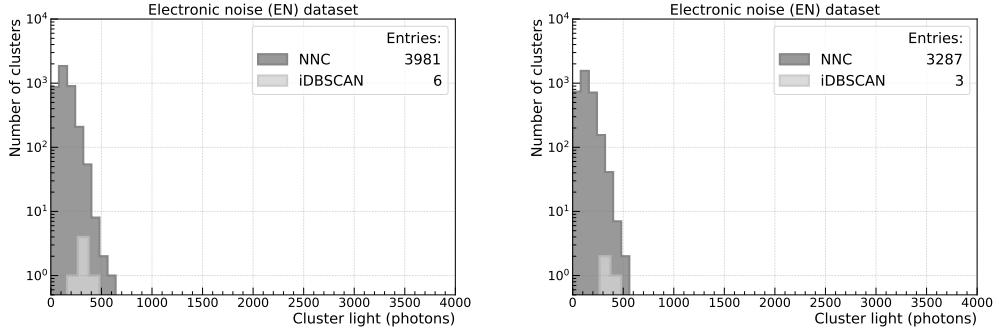
### 3 iDBSCAN and NNC comparison

#### 3.1 Electronic noise, natural radioactivity and $^{55}\text{Fe}$ energy spectra

The well-known energy deposition signature of 5.9 keV photons coming out from the  $^{55}\text{Fe}$  source is exploited in order to evaluate the detection efficiency and background rejection of both methods. While the ER dataset will be used for signal characterization, EN and NRAD datasets will be deployed for background rejection measurements. The EN acquired data produces low energy clusters with a distribution squeezed in the region below 500 photons as shown in Fig. 7, NRAD produces an energy distribution widely spread by a heavy tail component as shown in Fig. 8 while ER forms an additional narrow distribution centered at around 3000 photons as shown in Fig. 9. In this last case, the energy spectrum is composed of background and  $^{55}\text{Fe}$  induced deposits and, therefore, to reconstruct the  $^{55}\text{Fe}$  energy distribution, the background distribution should be subtracted. All the distributions were generated with the same amount of images, 864 of them, except for the iDBSCAN distributions of Fig. 7 which used 6478 images, in order to collect enough EN-clusters, which occur at a low rate. Additionally, the signal purity is enhanced accounting for the cluster aspect ratio, called slimness, defined as the ratio between the minor axis (width) and major axis (length) of each cluster.

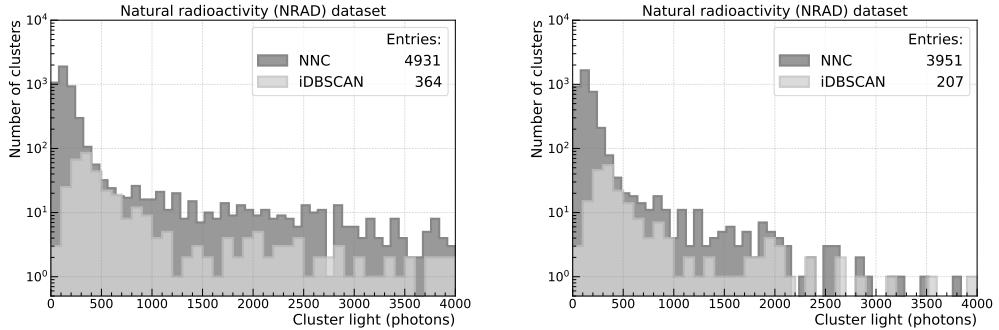
Figure 7 compares the energy spectrum of clusters generated by NNC with those generated by iDBSCAN for EN events without and with a selection based on the slimness parameter, considering only clusters with slimness greater than 0.4 for the later case. The computed numbers of EN-clusters per image for NNC and iDBSCAN were  $4.61 \pm 0.07$  and  $(9 \pm 4) \times 10^{-4}$ , respectively. Regarding NNC, EN-clusters dominate the background rate for energies below 500 photons which can be noticed by comparing the EN energy distribution of Fig. 7 with that of the NRAD shown in Fig. 8. Slimness selection decreases the number of clusters per image to  $3.80 \pm 0.07$  and  $(5 \pm 3) \times 10^{-4}$

215 for NNC and iDBSCAN, respectively. Therefore, when compared to NNC, iDBSCAN is able to  
 216 reduce the number of EN-clusters per image by a factor of  $(5 \div 7) \times 10^3$ .



**Figure 7.** Clusters energy distribution for NNC and iDSBSCAN applied to the EN dataset, without (left) and with (right) a selection on the slimness.

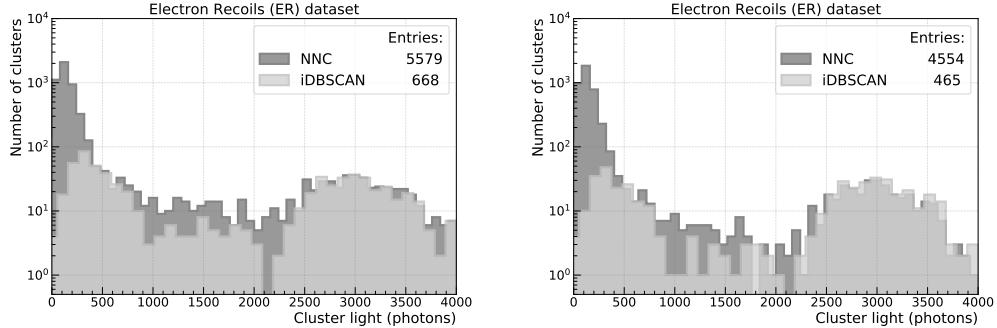
217 Figure 8 presents the energy distributions for the NNC and iDBSCAN clusters using the NRAD  
 218 dataset without (left) and with (right) a selection on slimness. iDBSCAN presents a clear peak  
 219 evolution around 300 photons while NNC accumulates clusters with lower energies. As mentioned  
 220 before, for NNC this region is highly populated by EN-clusters. iDBSCAN reduces the number of  
 221 background events in the region between 2000 and 4000 photons, which is the region where the  
 222  $^{55}\text{Fe}$  events are expected to be, as mentioned before, providing better background rejection for low  
 223 energy events as for the 5.9 keV photons. On the right of Fig. 8, the distribution of light, only  
 224 considering clusters with slimness greater than 0.4 is shown. As can be observed, this operation  
 225 has reduced even more the number of background events in the  $^{55}\text{Fe}$  region. However, for the lower  
 226 energy region, the number of background clusters has not changed much, causing iDBSCAN to  
 227 maintain a better background rejection efficiency when compared to NNC.



**Figure 8.** Clusters energy distribution for NNC and iDSBSCAN applied to the NRAD dataset, without (left) and with (right) a selection based on the slimness.

228 Figure 9 shows the results of the same analysis performed on the ER dataset. In this case,  
 229 the sum of the distribution obtained in the NRAD sample and the one from  $^{55}\text{Fe}$  interactions is  
 230 expected. As shown, both clustering algorithms are sensitive to the 5.9 keV photon events. However,

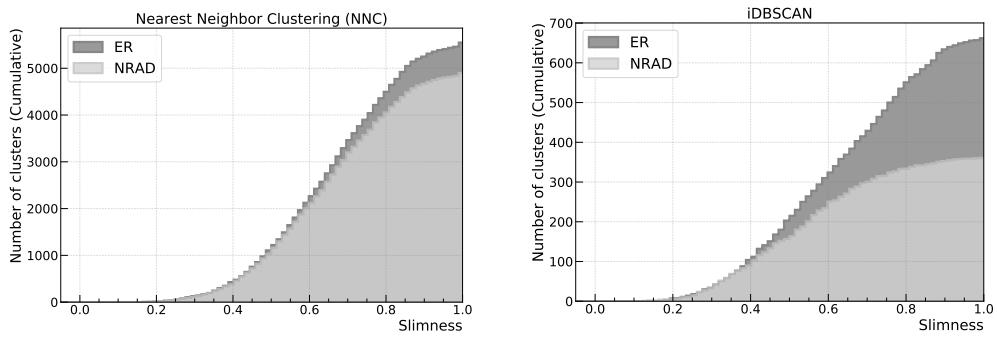
as commented previously, a higher purity level is achieved using iDBSCAN. After applying the slimness threshold, as shown in the right plot of Fig. 9, the NNC and iDBSCAN  $^{55}\text{Fe}$  peaks get closer indicating that both methods have similar detection efficiency considering that the number of  $^{55}\text{Fe}$  spots found by each method is practically the same.



**Figure 9.** Clusters energy distribution for NNC and iDSBSCAN applied to the ER dataset, without (left) and with (right) a selection based on the slimness.

### 3.2 Slimness selection optimization

Figure 10 shows the slimness cumulative distribution of clusters for an interval between zero and one applied to the NRAD and ER datasets. NNC and iDBSCAN cases are shown on the left and right plots, respectively. As it is possible to see, in both cases,  $^{55}\text{Fe}$  spots tend to have slimness higher than about 0.4. This variable can be used in conjunction with energy measurement to discriminate  $^{55}\text{Fe}$  spots from background clusters. In this section the value of slimness will be swept so that it is possible to choose the most suitable value for its use as an event selection parameter as well as to evaluate its impact when applied together with the energy measurement.

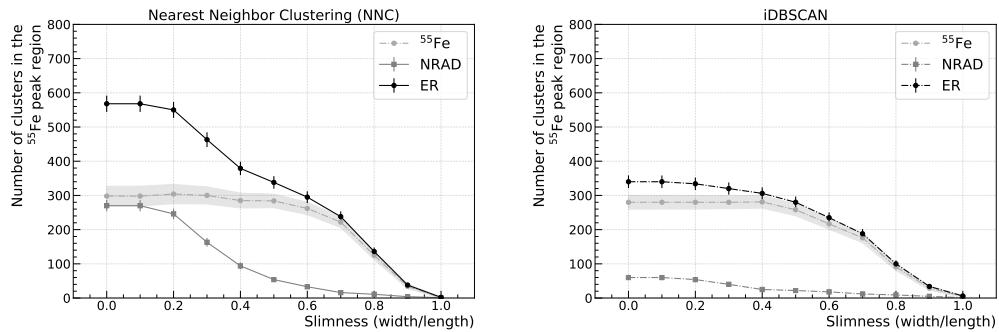


**Figure 10.** Cumulative distribution of the slimness for NRAD and  $^{55}\text{Fe}+\text{NRAD}$  data, for NNC (left) and iDSBSCAN (right).

In order to evaluate the signal efficiency and purity as a function of the slimness selection for the two algorithms, which should return a purer cluster sample given the differences between

the  $^{55}\text{Fe}$  and natural radioactive tracks topologies, the number of clusters within the selected  $^{55}\text{Fe}$  energy region (from 1500 to 4500 photons) was measured for various slimness threshold values ( $X \geq x$ ) as shown in Fig. 11 for the NNC (left) and iDBSCAN (right) algorithms. According to these curves, iDBSCAN finds a similar number of clusters in the  $^{55}\text{Fe}$  region when compared to NNC for slimness below 0.4, given by the difference between the black and gray curves, but with lesser contamination (gray curve).

Considering that the  $^{55}\text{Fe}$  clusters produce an intensity that follows a Gaussian distribution with an average value of about 3000 photons and standard deviations of 550 and 371, for NNC and iDBSCAN respectively (see Fig. 12), then more than 99% of the  $^{55}\text{Fe}$  clusters are selected between 1500 and 4500 photons. On the other hand, for the same region, the subtraction of the natural radioactivity events between the  $^{55}\text{Fe}$  and NRAD acquisition runs has a mean value equal to zero but a fluctuation of about 23 (14) and 11 (7) clusters for slimness equal to 0.0 (0.4), for NNC and iDBSCAN respectively. Therefore, the dashed line of Fig. 11 is composed mainly of  $^{55}\text{Fe}$  events plus few background events produced by the statistical fluctuation that occurs in the process of subtracting natural radioactivity. As can be noticed by observing Fig. 8, iDBSCAN tends to have less background contamination than NNC, reducing the statistical uncertainty related to the background subtraction. This effect is also shown by the shaded band drawn around the dashed lines of Fig. 11.



**Figure 11.** Scan in the number of clusters on the  $^{55}\text{Fe}$  peak region (between 1500 and 4500 photons) when changing the threshold on the slimness for NRAD and  $^{55}\text{Fe}$ +NRAD data, for NNC (left) and iDSBSCAN (right).

Based on the measurements of Fig. 11, the impact of the slimness parameter can be assessed by measuring selection efficiency ( $\varepsilon_{\text{sel}}$ ) and fake events ( $F_{\text{evts}}$ ), as defined below:

- $\varepsilon_{\text{sel}}$ : number of clusters found in the ER dataset ( $n_{\text{Fe}}$ ) subtracted by the number of clusters found in the NRAD dataset ( $n_{\text{Rd}}$ ) divided by the maximum value of the  $n_{\text{Fe}} - n_{\text{Rd}}$  subtraction among all slimness values (see Equation 3.1);

$$\varepsilon_{\text{sel}} = \left( \frac{n_{\text{Fe}} - n_{\text{Rd}}}{\max(n_{\text{Fe}} - n_{\text{Rd}})} \right) \quad (3.1)$$

- $F_{\text{evts}}$ : ratio between the number of clusters found in the NRAD dataset ( $n_{\text{Rd}}$ ) and the number

of clusters found in the ER dataset (nFe) (see Equation 3.2a). This measure can also be understood in terms of background rejection ( $B_{rj}$ ) as shown by Equation 3.2b;

$$F_{\text{evts}} = \left( \frac{nRd}{nFe} \right) \quad (a), \quad B_{rj} = 1 - F_{\text{evts}} \quad (b) \quad (3.2)$$

Figure 10 shows that for slimness below 0.4 the efficiency for background events is very small, while most of the  $^{55}\text{Fe}$  events are retained. Table 1 shows the computed  $\varepsilon_{\text{sel}}$  and  $F_{\text{evts}}$  for both clustering methods and different thresholds on the slimness variable ranging from 0.0 to 0.8. For the high efficiency region ( $\geq 0.94$ ), occurring for slimness values from 0.0 to 0.4, iDBSCAN achieved a lower fake event probability, always about 3 times less than NNC. For slimness greater than or equal to 0.6 both methods begin to lose efficiency. For slimness greater than 0.4, the signal is still almost 100% efficient, while the background is reduced by a factor 1/4.

**Table 1.**  $\varepsilon_{\text{sel}}$  and  $F_{\text{evts}}$  comparison between NNC and iDBSCAN.

Slimness (width/length)	$\varepsilon_{\text{sel}}$		$F_{\text{evts}}$		$B_{rj}(\%)$	
	iDBSCAN	NNC	iDBSCAN	NNC	iDBSCAN improvement	
0.0	1.00 +0.00 -0.02	0.98 +0.01 -0.02	0.18 +0.04 -0.04	0.48 +0.04 -0.04	56.98 +5.15 -5.41	
0.2	1.00 +0.00 -0.02	1.00 +0.00 -0.01	0.16 +0.04 -0.04	0.45 +0.04 -0.04	51.67 +4.50 -4.73	
0.4	1.00 +0.00 -0.01	0.94 +0.02 -0.03	0.08 +0.02 -0.03	0.25 +0.05 -0.04	22.60 +1.37 -1.62	
0.6	0.77 +0.05 -0.05	0.86 +0.03 -0.04	0.08 +0.04 -0.03	0.11 +0.04 -0.03	03.96 +0.19 -0.26	
0.8	0.32 +0.06 -0.05	0.41 +0.05 -0.06	0.09 +0.07 -0.04	0.08 +0.06 -0.04	-01.00 +0.10 -0.06	

The last column of Table 1 shows the iDBSCAN background-rejection improvement compared to NNC. For slimness equal to 0.4, for example, iDBSCAN has 92% of background rejection efficiency while NNC has 75%, leading to a relative improvement of  $(92-75)/75 \approx 23\%$ .

### 3.3 Light Yield Resolution

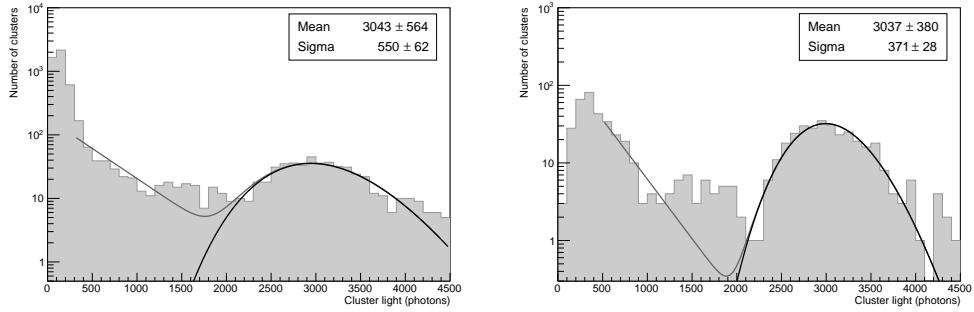
The detector energy resolution was estimated by a fit to the clusters energy distributions accounting for natural radioactivity and the  $^{55}\text{Fe}$  events. The former was modeled by an exponential function and the latter by a Polya function [17]:

$$P(n) = \frac{1}{b\bar{n}} \frac{1}{k!} \left( \frac{n}{b\bar{n}} \right)^k \cdot e^{-n/b\bar{n}} \quad (3.3)$$

where  $b$  is a free parameter and  $k = 1/b - 1$ . The distribution has  $\bar{n}$  as expected value, while the variance is governed by  $\bar{n}$  and the  $b$  parameter, as follows:  $\sigma^2 = \bar{n}(1 + b\bar{n})$ . The total likelihood is given by the sum of the two functions.

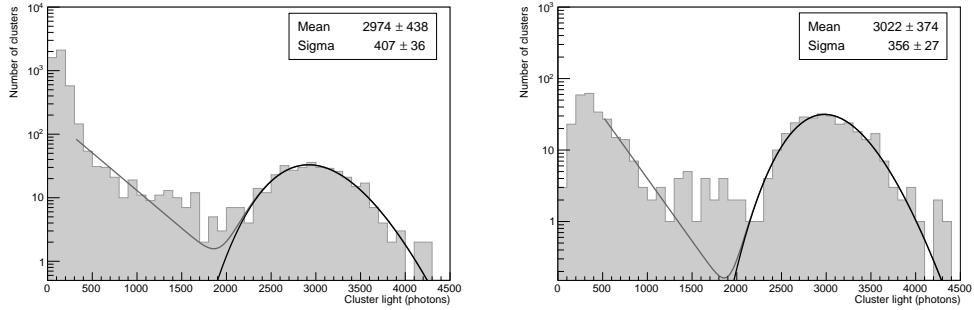
Figure 12 shows the fit results for NCC (left) and iDBSCAN (right) clusters without applying any selection on the slimness parameter. Based on the computed values, energy resolution were measured to be  $(18.1 \pm 3.9)\%$  and  $(12.2 \pm 1.8)\%$  for NNC and iDBSCAN respectively, and the energy conversion factor approximately 515 ADC units per keV for both. Conversion factor and Energy resolution are computed using the *mean* and *sigma* parameters shown in Fig. 12. The former is given by dividing the *sigma* by the *mean*, while for the latter the *mean* is divided by 5.9 keV (ER energy).

Figure 13 shows the fit results after applying a threshold on slimness: only clusters with slimness higher than 0.4 were considered. The estimated energy resolutions were  $13.7 \pm 2.4\%$  and



**Figure 12.** Results of the fit applied to the NNC (left) and iDBSCAN (right) energy distributions.

297 11.8  $\pm$  1.7% for NNC and iDBSCAN, respectively, with a conversion factor of about 510 ADC units  
 298 per keV. Finally, Table 2 shows the resulting energy resolution for NNC and iDBSCAN for different  
 299 values of slimness. Note that due to its higher background contamination, the energy resolution  
 300 obtained with NNC decreases while the slimness threshold increases, reaching eventually the energy  
 301 resolution obtained with iDBSCAN, which is already quite pure without any selection on slimness.



**Figure 13.** Results of the fit applied to the NNC (left) and iDBSCAN (right) energy distributions for clusters with slimness higher than 0.4.

**Table 2.** Detector resolution comparison between NNC and iDBSCAN as a function of slimness.

Slimness (width/length)	Resolution (%)	
	iDBSCAN	NNC
0.0	12.2 $\pm$ 1.8	18.1 $\pm$ 3.9
0.2	12.0 $\pm$ 1.7	17.3 $\pm$ 3.7
0.4	11.8 $\pm$ 1.7	13.7 $\pm$ 2.4
0.6	12.0 $\pm$ 2.0	11.8 $\pm$ 1.8
0.8	12.3 $\pm$ 3.8	11.1 $\pm$ 2.8

302 **4 Summary**

303 An adapted version of DBSCAN, named intensity-based DBSCAN, has recently been developed and  
304 tested on data acquired with a CYGNO TPC prototype. The impact of this algorithm on the detector  
305 performance has been studied using 5.9 keV photons from a  $^{55}\text{Fe}$  radioactive source and compared  
306 with results obtained with a simple NNC approach. The iDBSCAN parameters were optimized  
307 for the running conditions of LEMOn, which uses a 4M pixels sCMOS camera, and for signals  
308 from  $^{55}\text{Fe}$  photons. The obtained results showed that, with iDBSCAN, the clustering process of the  
309 CYGNO's event-reconstruction algorithm can achieve, without any other event-selection routine, a  
310 natural radioactivity background rejection in the energy region around 5.9 keV (from 3.0 keV to 8.8  
311 keV) of  $0.82_{-0.04}^{+0.04}$  and a number of electronic-noise clusters per image of  $(9 \pm 4) \times 10^{-4}$ , occurring  
312 predominantly in the region below 1 keV ( $\approx 500$  photons). Compared with NNC, these results  
313 represent an enhancement of 57% for the former and, for the latter, an improvement by a factor of  
314 a few thousand. Finally, the detector energy resolution using iDBSCAN was measured to be  $(12.2 \pm 1.8)\%$   
315 for 5.9 keV electron recoil events. By requiring spots with slimness larger than 0.4, a rate  
316 of electronic-noise clusters per image of  $(5 \pm 3) \times 10^{-4}$ , a natural radioactive background rejection  
317 of  $0.92_{-0.04}^{+0.03}$  and an energy resolution of  $(11.8 \pm 1.7)\%$  were achieved.

318 **Acknowledgments**

319 This work was supported by the European Research Council (ERC) under the European Union's  
320 Horizon 2020 research and innovation program (grant agreement No 818744) and also by the  
321 Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Finance Code  
322 001.

323 **References**

- 324 [1] M. Ester, H.-P. Kriegel, J. Sander and X. Xu, *A density-based algorithm for discovering clusters a*  
325 *density-based algorithm for discovering clusters in large spatial databases with noise*, in *Proceedings*  
326 *of the Second International Conference on Knowledge Discovery and Data Mining*, KDD'96,  
327 pp. 226–231, AAAI Press, 1996, <http://dl.acm.org/citation.cfm?id=3001460.3001507>.
- 328 [2] L. M. S. Margato, F. A. F. Fraga, S. T. G. Fetal, M. M. F. R. Fraga, E. F. S. Balau, A. Blanco et al.,  
329 *Performance of an optical readout GEM-based TPC*, *Nucl. Instrum. Meth. A* **535** (2004) 231.
- 330 [3] C. M. B. Monteiro, A. S. Conceicao, F. D. Amaro, J. M. Maia, A. C. S. S. M. Bento, L. F. R. Ferreira  
331 et al., *Secondary scintillation yield from gaseous micropattern electron multipliers in direct dark*  
332 *matter detection*, *Phys. Lett. B* **677** (2009) 133.
- 333 [4] C. M. B. Monteiro, L. M. P. Fernandes, J. F. C. A. Veloso, C. A. B. Oliveira and J. M. F. dos Santos,  
334 *Secondary scintillation yield from GEM and THGEM gaseous electron multipliers for direct dark*  
335 *matter search*, *Phys. Lett. B* **714** (2012) 18.
- 336 [5] A. Bondar, A. Buzulutskov, A. Grebenuk, A. Sokolov, D. Akimov, I. Alexandrov et al., *Direct*  
337 *observation of avalanche scintillations in a THGEM-based two-phase Ar avalanche detector using*  
338 *Geiger-mode APD*, *JINST* **5** (2010) P08002 [[1005.5216](#)].
- 339 [6] M. Marafini, V. Patera, D. Pinci, A. Sarti, A. Sciubba and E. Spiriti, *ORANGE: A high sensitivity*  
340 *particle tracker based on optically read out GEM*, *Nucl. Instrum. Meth. A* **845** (2017) 285.

- 341 [7] V. C. Antochi, E. Baracchini, G. Cavoto, E. D. Marco, M. Marafini, G. Mazzitelli et al., *Combined*  
 342 *readout of a triple-GEM detector*, *JINST* **13** (2018) P05001 [[1803.06860](https://arxiv.org/abs/1803.06860)].
- 343 [8] I. A. Costa, E. Baracchini, F. Bellini, L. Benussi, S. Bianco, M. Caponero et al., *Performance of*  
 344 *optically readout GEM-based TPC with a 55Fe source*, *Journal of Instrumentation* **14** (2019) P07011.
- 345 [9] Hamamatsu, *ORCA-Flash4.0 V3 Digital CMOS camera*, 2018.
- 346 [10] D. Pinci, E. Di Marco, F. Renga, C. Vona, E. Baracchini, G. Mazzitelli et al., *Cygnus: development*  
 347 *of a high resolution TPC for rare events*, *PoS EPS-HEP2017* (2017) 077.
- 348 [11] G. Mazzitelli, V. C. Antochi, E. Baracchini, G. Cavoto, A. De Stena, E. Di Marco et al., *A high*  
 349 *resolution tpc based on gem optical readout*, in *2017 IEEE Nuclear Science Symposium and Medical*  
 350 *Imaging Conference (NSS/MIC)*, pp. 1–4, Oct, 2017, [DOI](#).
- 351 [12] G. Lopes, E. Baracchini, F. Bellini, L. Benussi, S. Bianco, G. Cavoto et al., *Study of the impact of*  
 352 *pre-processing applied to images acquired by the cygno experiment*, in *Pattern Recognition and*  
 353 *Image Analysis*, pp. 520–530, Springer International Publishing, (09, 2019), [DOI](#).
- 354 [13] Y. Dong and S. Xu, *A new directional weighted median filter for removal of random-valued impulse*  
 355 *noise*, *IEEE Signal Processing Letters* **14** (2007) 193.
- 356 [14] I. T. Jolliffe, *Principal component analysis*, *Springer series in statistics* **29** (2002) .
- 357 [15] K. Wagstaff and C. Cardie, *Clustering with instance-level constraints*, in *Proceedings of the*  
 358 *Seventeenth International Conference on Machine Learning*, ICML '00, (San Francisco, CA, USA),  
 359 p. 1103–1110, Morgan Kaufmann Publishers Inc., 2000, [DOI](#).
- 360 [16] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel et al., *Scikit-learn:*  
 361 *Machine learning in python*, *Journal of Machine Learning Research* **12** (2011) 2825.
- 362 [17] W. Blum, L. Rolandi and W. Riegler, *Particle detection with drift chambers*, Particle Acceleration and  
 363 Detection, ISBN = 9783540766834. Springer Science & Business Media, 2008,  
 364 [10.1007/978-3-540-76684-1](https://doi.org/10.1007/978-3-540-76684-1).