

# Multiwavelength classification of X-ray selected galaxy cluster candidates using convolutional neural networks

Matej Kosiba,<sup>1,2\*</sup> Maggie Lieu,<sup>2,3</sup> Bruno Altieri,<sup>2</sup> Nicolas Clerc,<sup>4</sup> Lorenzo Faccioli,<sup>5</sup> Sarah Kendrew,<sup>6</sup> Ivan Valtchanov,<sup>7</sup> Tatyana Sadibekova,<sup>5,8</sup> Marguerite Pierre,<sup>5</sup> Filip Hroch,<sup>1</sup> Norbert Werner,<sup>9,1,10</sup> Lukáš Burget,<sup>11</sup> Christian Garrel,<sup>5</sup> Elias Koulouridis,<sup>12,5</sup> Evelina Gaynullina,<sup>8</sup> Mona Molham,<sup>13</sup> Miriam E. Ramos-Ceja<sup>14</sup> and Alina Khalikova<sup>8</sup>

<sup>1</sup>*Department of Theoretical Physics and Astrophysics, Faculty of Science, Masaryk University, Kotlářská 2, Brno, 611 37, Czech Republic*

<sup>2</sup>*European Space Astronomy Centre, ESA, Villanueva de la Cañada, E-28691 Madrid, Spain*

<sup>3</sup>*Centre for Astronomy and Particle Theory, University of Nottingham, UK*

<sup>4</sup>*IRAP, Université de Toulouse, CNRS, CNES, UPS, (Toulouse), France*

<sup>5</sup>*AIM, CEA, CNRS, Université Paris-Saclay, Université Paris Diderot, Sorbonne Paris Cité, F-91191 Gif-sur-Yvette, France*

<sup>6</sup>*European Space Agency, Space Telescope Science Institute, 3700 San Martin Drive, Baltimore MD 21218, USA*

<sup>7</sup>*Telespazio Vega UK for ESA, European Space Astronomy Centre, Operations Department, 28691 Villanueva de la Cañada, Spain*

<sup>8</sup>*Ulugh Beg Astronomical Institute of Uzbekistan Academy of Science, 33 Astronomicheskaya str., Tashkent, UZ-100052, Uzbekistan*

<sup>9</sup>*MTA-Eötvös University Lendület Hot Universe Research Group, Pázmány Péter sétány 1/A, Budapest, 1117, Hungary*

<sup>10</sup>*School of Science, Hiroshima University, 1-3-1 Kagamiyama, Higashi-Hiroshima 739-8526, Japan*

<sup>11</sup>*Faculty of Information Technology, Brno University of Technology, Božetěchova 2, Brno, 612 00, Czech Republic*

<sup>12</sup>*Institute for Astronomy & Astrophysics, Space Applications & Remote Sensing, National Observatory of Athens, GR-15236 Palaia Penteli, Greece*

<sup>13</sup>*National Research Institute of Astronomy and Geophysics (NRIAG), 11421 Helwan, Egypt*

<sup>14</sup>*Max-Planck Institut für extraterrestrische Physik, Postfach 1312, 85741 Garching bei München, Germany*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

Galaxy clusters appear as extended sources in *XMM-Newton* images, but not all extended sources are clusters. So, their proper classification requires visual inspection with optical images, which is a slow process with biases that are almost impossible to model. We tackle this problem with a novel approach, using convolutional neural networks (CNNs), a state-of-the-art image classification tool, for automatic classification of galaxy cluster candidates. We train the networks on combined *XMM-Newton* X-ray observations with their optical counterparts from the all-sky Digitized Sky Survey. Our data set originates from the X-CLASS survey sample of galaxy cluster candidates, selected by a specially developed pipeline, the **X**Amin, tailored for extended source detection and characterisation. Our data set contains 1 707 galaxy cluster candidates classified by experts. Additionally, we create an official Zooniverse citizen science project, *The Hunt for Galaxy Clusters*, to probe whether citizen volunteers could help in a challenging task of galaxy cluster visual confirmation. The project contained 1 600 galaxy cluster candidates in total of which 404 overlap with the expert's sample. The networks were trained on expert and Zooniverse data separately. The CNN test sample contains 85 spectroscopically confirmed clusters and 85 non-clusters that appear in both data sets. Our custom network achieved the best performance in the binary classification of clusters and non-clusters, acquiring accuracy of 90 %, averaged after 10 runs. The results of using CNNs on combined X-ray and optical data for galaxy cluster candidate classification are encouraging and there is a lot of potential for future usage and improvements.

**Key words:** galaxies: clusters: general – methods: data analysis – techniques: image processing

## 1 INTRODUCTION

Galaxy clusters are massive systems at the peaks of the cosmic web. Their composition, rich in dark matter and hot baryonic gas makes them a potentially powerful tool to constrain cosmological parameters, growth of structure, neutrino mass and sterile neutrinos through cluster number counts, the cluster mass function and the baryon fraction (Allen et al. 2011; Mantz et al. 2015; Böhringer & Chon 2016).

In recent years, large cluster surveys such as XXL (Pierre et al. 2016; Pacaud et al. 2016), XCS (Mehrtens et al. 2012), X-CLASS (Clerc et al. 2012a; Ridl et al. 2017), *Planck* (Bartlett et al. 2008), redMaPPer (Rykoff et al. 2014), or the SPT-SZ survey (Bleem et al. 2015) have made it possible to statistically improve constraints on cosmology. However one of the challenges in using galaxy clusters for cosmology is understanding and modelling of the cluster selection function (e.g. Pacaud et al. 2006). The selection function has to be modelled in terms of observable parameters (like flux and apparent size), which can then be converted into galaxy cluster mass for a given cosmology and galaxy cluster physics evolution. The selection function of galaxy clusters is not trivial to model and often oversimplified. A selection function should not only take into account the volume and redshift of the survey but also the choice of clusters, which is often more complicated than a cut in flux. In X-ray wavelengths, whilst extended emission is generally a robust indicator of a galaxy cluster, the emission can also be attributed to nearby galaxies, saturated AGN and unresolved double point-sources. For this reason, galaxy cluster candidates are still visually examined together with optical data, prior to any spectroscopic confirmation (Adami et al. 2018). This process is tedious and out-dated with uncertainties impossible to model. With large X-ray sky surveys such as *e-ROSITA* (Merloni et al. 2012) expecting to discover tens of thousands of new galaxy clusters, combined with large optical surveys including LSST (Ivezic et al. 2008) and *EUCLID* (Racca et al. 2016), the old techniques will become obsolete. We need to prepare for the future with new methods that are able to deal with big data and improved accuracy.

Citizen science projects proved to be a great asset for scientific problems where human classifications are required for large amounts of data (e.g. Lintott et al. 2008; Willett et al. 2013). In the first version of the most well known of all citizen science projects, the Galaxy Zoo (Lintott et al. 2008), citizen volunteers managed to achieve more than 90% agreement with experts in a task of morphological classification of galaxies. While citizen projects are intended to provide huge manpower in the assessment of large astronomical data sets, the question whether this is an advantage over a limited number of evaluations by experts in the case of the confirmation of galaxy cluster candidates remains to be addressed. This paper scrutinizes this issue by evaluating the citizen volunteers success rate.

Machine learning offers a more constructive approach to the problem. The power of Machine learning has been demonstrated in astronomy for more than two decades, with applications including star-galaxy discrimination (Odewahn et al. 1992; Bertin 1993), classification of galaxy spectra (Folkes et al. 1996), photometric redshift estimation (Collister & Lahav 2004) or anomaly detection in X-ray spec-

tra (Ichinohe & Yamada 2019), to name a few. With the introduction of Convolutional Neural Networks (CNNs, Le-Cun et al. 1999) and deep learning (E Hinton 2007), it has been possible to automate human vision tasks such as image recognition (see e.g. Goodfellow et al. 2014; Schawinski et al. 2017; Ackermann et al. 2018; Lieu et al. 2018).

Supervised learning with convolutional neural networks (CNNs) was designed specifically for image classification tasks. If the true labels (classification classes) of the images are known, they can be used to train CNNs. The current way galaxy clusters are classified are liable to false positives and false negatives. Galaxy cluster candidates picked by an automated pipeline are visually analysed by several experts to create an initial catalogue of galaxy clusters, that are later verified with a spectroscopic confirmation. This process will not scale with large data volumes. Citizen science allows us to harness a large number of opinions on each object classification on a short timescale, speeding up the process significantly yet having a reasonable agreement with experts (see e.g. Willett et al. 2013; Dieleman et al. 2015). CNNs can be then trained on classifications made by either experts or citizen volunteers or both, to automate the final classification of galaxy cluster candidates, or even skipping the first step of the pipeline picking the candidate clusters. Applying CNN selection on simulations will enable modelling the selection function.

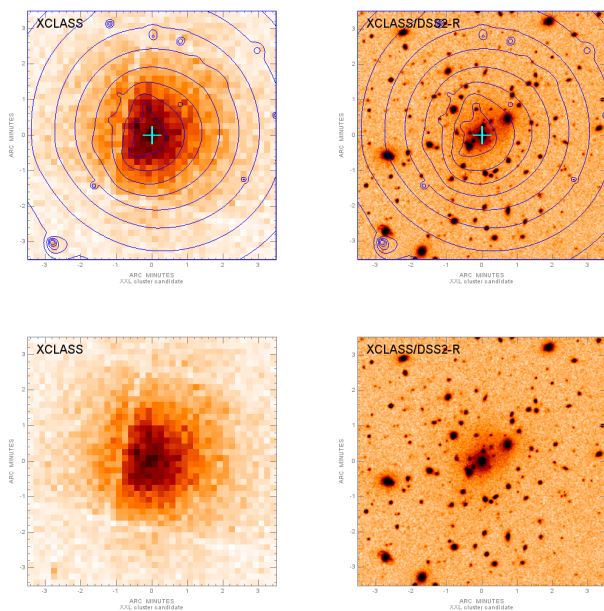
In this paper, we introduce a citizen science project we created to obtain large numbers of classified objects. We compare the performance of citizen volunteers with experts. We train CNNs on classifications of citizen volunteers and experts and compare their results. CNNs are tested on spectroscopically confirmed galaxy clusters and objects classified as non-clusters by experts.

The structure of the paper is as follows: in Section 2 we present our citizen science project and its development together with a description of the observations and the construction of their classifications by the experts, in Section 3 we introduce the machine learning methods we use, Section 4 presents measurements used to evaluate classification or detection performance, Section 5 presents the results of the citizen science campaign as well as the results and discussion of neural networks analysis. Finally, we conclude in Section 6.

## 2 THE HUNT FOR GALAXY CLUSTERS

Our citizen science project, *The Hunt for Galaxy Clusters*<sup>1</sup>, was launched online as an official Zooniverse project on the 24th of October 2018. There were 1 600 galaxy cluster candidates in the project that have been detected as extended X-ray sources by the *XAmin* wavelet-based pipeline (Pacaud et al. 2006). Each object was classified by at least 30 different volunteers and this was completed by the 29th of April 2019. 1 227 volunteers participated in the project. Classifications of not logged in volunteers, as well as classifications which have been done on each object multiple times by the same volunteer, were not considered.

<sup>1</sup> <https://www.zooniverse.org/projects/matej-dot-kosiba/the-hunt-for-galaxy-clusters>



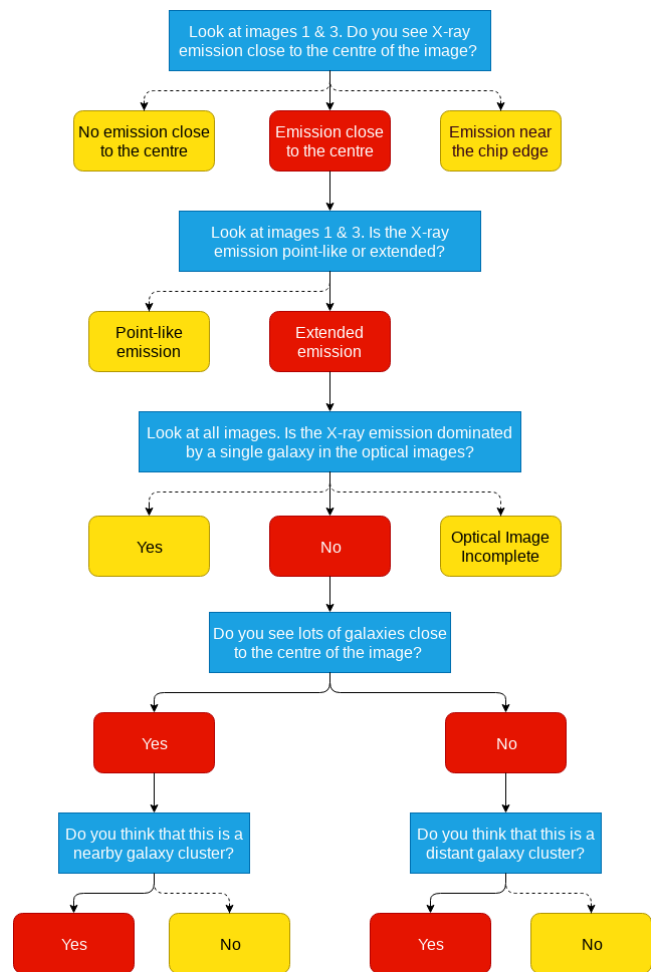
**Figure 1.** *Top left:* raw X-ray image with contours showing the areas of constant X-ray brightness and a cyan cross marking the object selected for classification. *Bottom left:* raw X-ray image without contours and markings. *Right:* corresponding optical images.

The project starts with a short tutorial briefly explaining how to navigate in the project’s page and how to classify candidate clusters. Each object comes with four images, covering the exact same area of the sky ( $7 \times 7$  arcmin<sup>2</sup>): two X-ray and two optical images. Figure 1 shows all four images of a galaxy cluster candidate as shown to the volunteers in *The Hunt for Galaxy Clusters*.

Our project uses six questions to help determine the class of a galaxy cluster candidate. Each question has two or three possible answers, and due to the structure of the decision tree (Figure 2), only a subset of the questions are answered. Those questions come with help notes, example images, as well as descriptions to each answer. We selected example images very carefully to cover a broad range of objects and/or instrument effects, in order to avoid biases. The Zooniverse volunteer’s answers were then used to create a binary classification scheme of *cluster* and *non-cluster*.

## 2.1 Data

The data in this work originates from the XMM CLuster Archive Super Survey (X-CLASS) (Clerc et al. 2012b), an X-ray galaxy cluster search in the archival data of the *European Space Agency’s* X-ray observatory *XMM-Newton*, combined with corresponding optical counterparts from the Digitized Sky Survey POSS-II (DSS2). We used *XMM-Newton* data obtained between 2000 and 2015, employing selection criteria described in (Clerc et al. 2012b), and excluding the data used by the XXL survey (Pierre et al. 2016).



**Figure 2.** The decision tree of *The Hunt for Galaxy Clusters* Zooniverse citizen science project. Blue cells represent questions, red are answers leading to the *cluster* class and yellow are answers leading to the *non-cluster* class.

## 2.2 X-ray pipeline

Our sample of galaxy cluster candidates has been constructed using the intermediate XAmin 3.5 version (new source models added: double point-source and point + extended source). This version, after the processing of the X-CLASS survey, appeared to suffer from a miss-centering problem randomly affecting a tiny fraction of the point-source population, that led to classify them as extended. In order to remove miss-classified sources, experts then performed an in-depth screening of the putative cluster candidate lists. The screening dealt as well with usual nearby galaxies and saturated AGNs, that both appear extended in the X-ray images

The pipeline is briefly described below. Firstly, a combined MOS1+MOS2+PN image of an *XMM-Newton* (Jansen 1999) observation is smoothed with a dedicated wavelet smoothing program called *mr.filter*, described by Starck et al. (1998) and shown in Starck & Pierre (1998) to effectively recover structures in X-ray images characterised by low numbers of photons.

Secondly, the wavelet smoothed image is analysed

by the source extraction software **SExtractor** (Bertin & Arnouts 1996). It creates a list of candidate sources for further analysis, returning an estimate of their position and their flux.

Note that, since **SExtractor** was developed for optical images which contain many more photons than the X-ray ones, smoothing the X-ray image is a necessity as **SExtractor** would not be able to work with raw data. This smoothing can be performed in several ways; the wavelet smoothing used by **XAmin** is one of the possible ways of smoothing the image and was shown by Valtchanov et al. (2001) to give the best results for X-ray images of diffuse sources like galaxy clusters.

Finally, we characterise the candidate sources found by **SExtractor**. This is done by fitting both a point source model given the *XMM-Newton* PSF computed at the source position and an extended  $\beta$  model (Cavaliere & Fusco-Femiano 1976) which better describes galaxy clusters. A source is declared to be a point source (AGN or an extended source too faint to be characterised as extended) or an extended source (galaxy cluster) depending on which of these two models best fits the candidates source. The details, including the relevant formulas and the selection criteria for defining an (almost) pure sample of galaxy clusters, are given in Pacaud et al. (2006).

Coordinates of the galaxy cluster candidates picked by **XAmin** are then used to produce normalised images (Appendix A) with and without X-ray contours to show lines of constant X-ray brightness. These contours are superimposed onto the optical counterpart image, together with a cyan cross mark and are used only for human screening to help visualise the X-ray emission.

### 2.3 Weighting volunteers classifications

Since each object is classified by 30 volunteers, we may end up with different classifications for the same galaxy cluster candidate. Each person's classification ability may vary according to the class and the question asked, and there may even be volunteers who purposely create malicious classifications. To mitigate those effects, we weight classifications of each user question-wise. Weighting is done according to the agreement of the majority, so each user has an accuracy determining a portion of his/her classifications being in agreement with the majority of votes, which is done question-wise,

$$G_i = \frac{C_i}{Q_i}, \quad i \in 1, \dots, 6 \quad (1)$$

where  $G_i$  is the weight applied for an individual on question  $i$ ,  $C_i$  is the number of answers to question  $i$  given by the individual that were in agreement with the majority and  $Q_i$  is the total number of answers the individual has made for question  $i$ .  $G_i$  essentially describes the ability of an individual to classify as the majority of volunteers would. Every classification in the project is then weighted according to the classifying volunteer's accuracy for the specific question. The bottom red leaves of the decision tree (Figure 2) are classification ending answers corresponds to the final answers stating that the classified object is a galaxy cluster. Similarly, all yellow leaves corresponds to the final answers stating that the object is not a galaxy cluster. Each galaxy cluster candidate gets 30 votes, each vote is an accuracy of the voting user

for the question of his/her classification ending answer (one of bottom red leafs or any yellow leaf). Those 30 weighted scores are summed to galaxy cluster (bottom red leafs) and non-galaxy cluster (yellow leafs) categories. The higher score determines the final Zooniverse weighted classification for the galaxy cluster candidate.

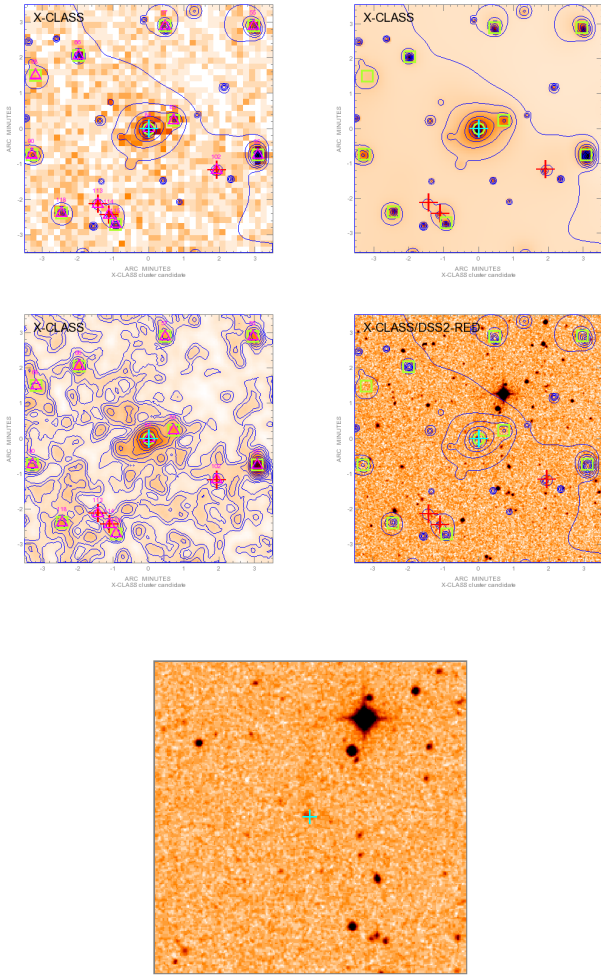
### 2.4 Classifications of experts

The galaxy cluster candidates generated by the **XAmin** pipeline are manually classified by the X-CLASS collaboration. Each galaxy cluster candidate is classified by two experts and three moderators make the final classification on conflicting decisions. Figure 3 shows how a galaxy cluster candidate is presented to the experts. The images are provided without redshift or sky coordinate information, and the experts make decisions without consulting with each other to avoid any bias. The experts were given the opportunity to classify objects as a low redshift cluster ( $0 < z < 0.3$ ), high redshift cluster ( $z > 0.3$ ), nearby galaxy, point source, star or AGN, double source, artefact, edge, fossil group, high background image, no optical image or dubious source. We create a binary classification scheme where the last four categories in the list are not used, low and high redshift clusters are collectively referred to as clusters and the remaining classes are collectively referred to as non-clusters.

## 3 MACHINE LEARNING APPROACH

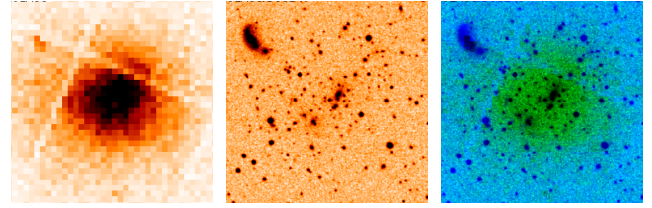
Now, we turn our attention to a machine learning approach, which allows us to automatically process astronomical data on much larger scales than what is possible to achieve by human annotations. We use neural networks – a parametric model, that is able to learn to approximate a complex function from training examples of inputs and the corresponding outputs. In our case, each training example consists of combined X-ray and optical image as the input and the corresponding output class label obtained from a human annotator. In our experiments, we consider binary classification, where the class labels are *galaxy cluster* and *non-cluster*, but also multi-class classification with subcategories that will be discussed in Section 5.5. From the training examples, our neural networks learn to predict posterior probabilities of all classes given an input image. In our experiments, we evaluate the performance of the neural networks using measures discussed in Section 4. For some of the measures, we need to make a hard classification decision for each input image from our evaluation set. In such a case, we simply select the most probable class.

In this work, we use Convolutional Neural Networks (CNN), which is currently the most popular and very effective neural network architecture for image processing (Lecun et al. 1989; Ciresan et al. 2012; Krizhevsky et al. 2012). A deeper knowledge of CNNs is not necessary for interpreting our results and understanding the presented analyses. It is only necessary for understanding some of the technical details. This paper also can not give a complete tutorial to CNNs, therefore, we do not provide a further introduction to CNNs and we kindly refer the interested reader to the relevant textbooks (Goodfellow et al. 2016; Bishop 2006) or the numerous tutorials available online. We use two



**Figure 3.** Images of a galaxy cluster candidate classified by experts. Top left: an X-ray raw image overplotted with contours showing areas of constant X-ray brightness, and marks produced by the *XAmin* pipeline. Top right and bottom left images are smoothed versions of the X-ray images, wavelet and Gaussian smoothing produced by the *XAmin* pipeline, respectively. The Gaussian smoothed image is overplotted with Gaussian contours, the sigma is chosen to be 3 pixels (with a pixel size of 2.5 arc seconds so the sigma is 7.5 arc seconds). Bottom right: the optical counterpart of the X-ray image with superimposed marks and wavelet contours. All images cover the exact same area of the sky,  $7 \times 7$  arcmin<sup>2</sup>, except for the bottom panel, where we focus in the central region ( $4 \times 4$  arcmin<sup>2</sup>) of the optical image, because with the contours and the symbols it is not easy to see the central cluster brightest galaxy and overdensity of faint galaxies.

CNNs architectures for our experiments: Using the Keras toolkit (Chollet et al. 2015), we build and train our *custom network*, which uses a conventional CNN architecture with interleaving convolutional and pooling layers and final dense layers. The second architecture is MobileNet (Howard et al. 2017). We take these networks as provided by its authors pre-trained on the ImageNet (Deng et al. 2009) data, which is a large data set of millions of real-world images categorised into thousands of classes. We assume that such pre-training



**Figure 4.** Left is the  $356 \times 356$  pixel X-ray rgb .PNG image, middle is its  $356 \times 356$  pixel optical .PNG counterpart and right is an rgb .PNG image made by stacking grayscale optical image as blue channel, grayscale X-ray image as green channel and the red channel was filled with zeros.

can serve as a good initialisation of the CNN parameters, which are further retrained on our training data for galaxy cluster classification.

### 3.1 Data preprocessing

For training neural networks, we use images without contours and marks. For each candidate cluster, a pair of X-ray and optical PNG images were merged into a single PNG image. As well as our custom network, we use existing architectures, that were designed to take input images with 3 colour channels. In order to achieve this, we grayscale the X-ray and optical images and stack them together as individual channels, leaving one channel empty (zero-filled) to create a single RGB image. Although training of our custom network can be done with any number of input channels, we use the same 3-channel images as the input to the network unless stated otherwise. By default, we construct the input images as follows: the blue channel contains the grayscale optical image, the green channel contains the grayscale X-ray image and the red is filled with a matrix of zeros (Figure 4).

### 3.2 Data augmentation

With smaller data sets, the risk of over-fitting increases, resulting in poor generalisation to data outside of the training set. To prevent overfitting, we use data augmentation to reduce the probability that the network will see exactly the same image twice and to essentially increase our training sample size. At each training step, the input image is randomly scaled to a uniform value between  $1/1.3$  and  $1.3$ , rotated by a random uniform angle between  $0$  and  $360^\circ$  and translated in x and y directions by a random uniform value between  $-4$  and  $4$  pixels.

## 4 PERFORMANCE MEASUREMENTS

This section describes the measurement methods we chose to evaluate our neural networks compared to a baseline.

Accuracy is the most intuitive performance measurement. It is the ratio of correct predictions to all predictions and is defined as

$$A = \frac{TP + TN}{TP + TN + FP + FN}, \quad (2)$$

where  $TP$  refers to the number of true positives, in our case the number of clusters correctly classified as clusters,  $TN$  is

a number of true negatives (number of non-clusters correctly classified as non-clusters),  $FP$  is a number of false positives (number of non-cluster incorrectly classified as clusters) and  $FN$  states for a number of false negatives (number of clusters incorrectly classified as non-clusters).

Precision is the ratio of the correctly classified positives (i.e. clusters) and all objects classified as positives. This is defined as

$$P = \frac{TP}{TP + FP}. \quad (3)$$

Recall is the ratio of the correctly classified positives and all positives examples in the test data. It is defined as

$$R = \frac{TP}{TP + FN}. \quad (4)$$

The receiver operating characteristic (ROC) is a performance measurement of detection problems plotted as a true positive rate (recall) against the false positive rate, defined as

$$FPR = \frac{TN}{TN + FP} \quad (5)$$

at various thresholds. The area under the curve (AUC) describes the model's capability to distinguish between two classification classes and is independent of the choice of the threshold. When reporting detection performance for a class (from the CNN output) in terms of ROC curve, we compare the posterior probability of the class to a varying detection threshold.

## 5 RESULTS AND DISCUSSION

### 5.1 The Hunt for Galaxy Clusters results

The data set of 1600 galaxy cluster candidates in *The Hunt for Galaxy Clusters* contained 404 objects previously classified by experts.

Table 1 displays a comparison of the unweighted and weighted classifications of the Zooniverse volunteers (subsection 2.3) based on the agreement with the experts. Figure 5 shows ROC curves computed for the whole crossmatch sample of 404 objects classified by both the Zooniverse volunteers and experts and the ROC computed on a subsample of 170 objects, 85 spectroscopically confirmed galaxy clusters and 85 objects classified as non-clusters by experts. This subsample is also used for the testing of the CNNs. The Zooniverse volunteers performed better on the subsample of 170 objects than on the whole crossmatch sample of 404 objects. This could be an indication of a bias towards correctly classifying easier objects since spectroscopically confirmed galaxy clusters tend to be larger.

Figure 6 shows the fraction of the Zooniverse volunteer's individual answers in agreement with experts to all Zooniverse answers for classification ending answers, except for *not a nearby cluster* and *not a distant cluster*, which do not have a direct counterpart in the classification of experts. Assuming that the expert classifications are the ground truth, the biggest difficulty for the volunteers seems to be distinguishing extended from point-like X-ray emission. Also,

the volunteers inconsistently classified a large fraction of *no emission* classes, suggesting that they struggled to interpret the X-ray images. The huge discrepancy between volunteer's individual classifications and classifications of experts were in the *edge* category, used for galaxy cluster candidates close to the edge of *XMM-Newton's* chips and its field of view. Based upon discussions within the online forum, we assume that this bias could emerge from *XMM-Newton's* grid-like pattern created by small gaps between its individual detectors, which volunteers often mistaken for the edge of the chips. The *nearby galaxy* category was also a difficult question for the volunteers. Again based on the forum discussion we find that volunteers often classified nearby galaxy clusters with a prominent brightest central galaxy as a *nearby galaxy* class, which could lead to many nearby galaxy clusters missed. In general, the Zooniverse volunteers preferentially classified objects as *non-clusters*.

Some of the biases could be mitigated in possible future versions of the project if explanations were clearer and more focus was put on example images in the help notes. Possibly the most important biases were often a classification of an X-ray emission as *no emission* and misclassification of an extended X-ray emission as a point-like X-ray emission. This are the main reasons why clusters were missed by the Zooniverse volunteers. We tried to keep in mind the possibility of low scientific knowledge of the volunteers and not to overwhelm the volunteers with huge amounts of information, which could discourage them, but we were still able to provide a detailed explanation of the X-ray emission in the tutorial and the help notes, with nice example images and diagrams to help with the X-ray contours. Small interviews with our beta testers revealed that around 20% of them did not read the supporting texts. It might be possible that classifications with a lot of disagreement in the interpretation of the X-ray emission preferentially came from volunteers who did not adequately read the supporting material. A questionnaire would be needed to further probe this possibility. These biases could be cut down with simpler and shorter explanations of the X-ray properties, so it would be easier to understand and less information to digest. Another common tendency was the misclassification of nearby clusters that contain prominent BCGs (brightest cluster galaxies), with that of nearby galaxies. This could be reduced with a dedicated pair of images for the two situations in the help notes.

We have to note that even the classifications of experts could be biased towards *low-z* clusters, since we use DSS optical images, which are limited to  $z \sim 0.3$ .

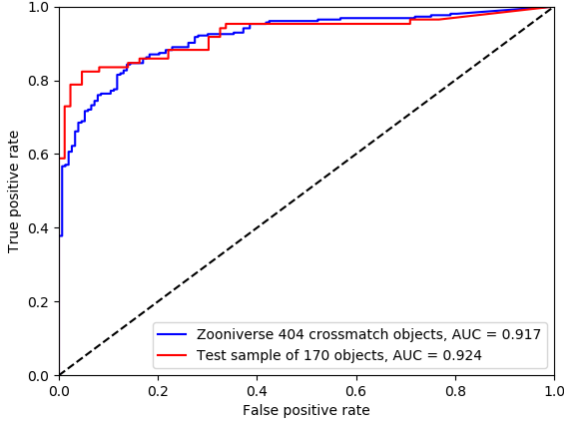
Another possible bias may come from the fact that spectroscopically confirmed clusters are biased to big clusters, which might affect our interpretation.

To explore if the Zooniverse volunteers were biased finding preferentially most prominent galaxy clusters, we made extent – extension likelihood plane plots (see Appendix B). We found that the galaxy clusters found by the Zooniverse volunteers populate all of the space, not showing bias and their sample of galaxy clusters also can not be recreated by a simple cut in this space.

Even though the Zooniverse volunteers did not show a high accuracy compared to experts, misclassifying many galaxy clusters as other options, the sample of galaxy clusters they selected is pure. This makes us conclude that, via

**Table 1.** The results of cluster classification by Zooniverse volunteers on two data sets, 404 objects are those classified by both, scientists and Zooniverse volunteers, the 170 objects data set is a subsample of the 404 objects, where 85 objects are spectroscopically confirmed clusters and 85 are objects classified as non-clusters by experts.

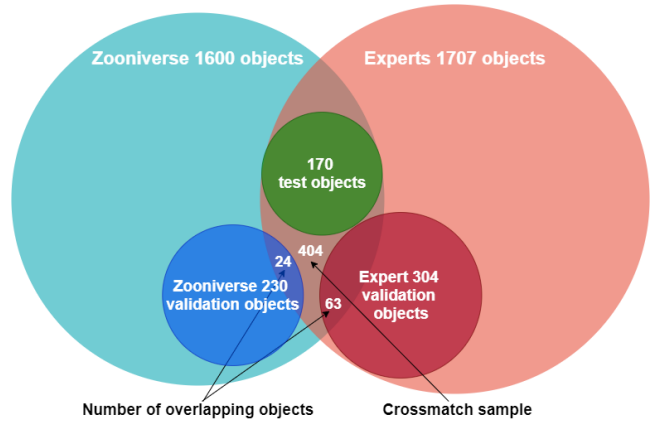
| Data set    | Zooniverse Classifications | TP  | TN  | FP | FN  | accuracy | precision | recall |
|-------------|----------------------------|-----|-----|----|-----|----------|-----------|--------|
| 404 objects | unweighted                 | 69  | 150 | 0  | 185 | 0.542    | 1.000     | 0.272  |
| 404 objects | weighted                   | 102 | 149 | 1  | 152 | 0.621    | 0.990     | 0.401  |
| 170 objects | weighted                   | 55  | 84  | 1  | 30  | 0.818    | 0.982     | 0.647  |



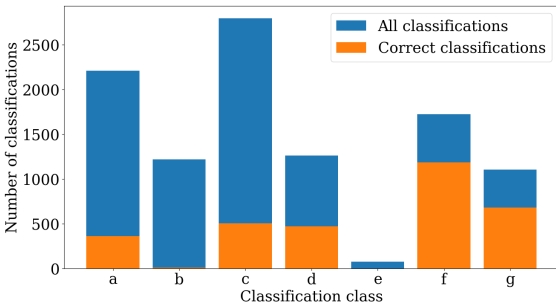
**Figure 5.** The receiver operating characteristic (ROC) curves for the classifications by Zooniverse volunteers, taking the classifications of experts as the ground truth. Closer the curve copies the left vertical and top horizontal axis, better the classifier. The dashed line shows how would the results be if the people guessed totally randomly.

**Table 2.** The number of objects in the training, validation and test data sets classified by Zooniverse and experts.

| Class       | Zooniverse |          | Experts |          |      |
|-------------|------------|----------|---------|----------|------|
|             | Train      | Validate | Train   | Validate | Test |
| cluster     | 320        | 130      | 845     | 200      | 85   |
| non-cluster | 880        | 100      | 388     | 104      | 85   |
| total       | 1200       | 230      | 1233    | 304      | 170  |



**Figure 7.** A Venn diagram presenting the data sets.



**Figure 6.** A quantification of the Zooniverse classifications for a) no emission, b) edge, c) point, d) nearby galaxy, e) no optical image, f) nearby galaxy cluster, g) distant galaxy cluster, assuming the ground truth is the expert classification.

the Zooniverse project, the general public can help scientific research where a very pure sample of galaxy clusters is required, but it did not prove to be helpful in a case where a sample of galaxy clusters should be complete.

## 5.2 CNN training

We use two different data sets, one classified by experts and one by the Zooniverse volunteers. We use balanced training

batches, containing the same number of classification classes, randomly sampled from the training data. This is to prevent the network from being biased towards the class that occurs most frequently in the training sample.

Regardless of the training data, all the networks were tested on the same data set of 85 spectroscopically confirmed galaxy clusters and 85 objects classified as non-clusters by the experts, the 170 test objects. Table 2 and Figure 7 describe the numbers of objects used in the training, validation and test data sets, classified by experts and the Zooniverse volunteers for testing on the 170 object test sample. All the networks were trained on grayscale and combined X-ray and optical images as described in subsection 3.1 if not stated otherwise.

We experimented with both a custom network (Table 3) and using 3 different state of the art CNN architectures: VGG19 (Simonyan & Zisserman 2014), InceptionV3 (Szegedy et al. 2015) and MobileNet (Howard et al. 2017). We used those networks with their pre-trained weights, using a large learning rate and unfreezing all the

**Table 3.** The architecture of our custom network which achieved the best performance. Each of the convolutional and dense layers is followed by a ReLU non-linearity with the exception of the final output dense layer which has the softmax for classification.

| Layer | Layer type | filter shape / stride | input shape |
|-------|------------|-----------------------|-------------|
| 1     | conv       | 3×3×64/(1, 1)         | 356×356×3   |
| 2     | max pool   | 2×2/(2, 2)            | 356×356×64  |
| 3     | conv       | 3×3×32/(1, 1)         | 178×178×64  |
| 4     | max pool   | 2×2/(2, 2)            | 178×178×32  |
| 5     | conv       | 3×3×32/(1, 1)         | 89×89×32    |
| 6     | max pool   | 2×2/(2, 2)            | 89×89×32    |
| 7     | conv       | 3×3×32/(1, 1)         | 45×45×32    |
| 8     | max pool   | 2×2/(2, 2)            | 45×45×32    |
| 9     | conv       | 3×3×32/(1, 1)         | 23×23×32    |
| 10    | max pool   | 2×2/(2, 2)            | 23×23×32    |
| 11    | conv       | 3×3×32/(1, 1)         | 12×12×32    |
| 12    | max pool   | 2×2/(2, 2)            | 12×12×32    |
| 13    | flatten    | -                     | 6×6×32      |
| 14    | dense      | 256                   | 1152        |
| 15    | dense      | 2                     | 256         |

**Table 4.** Hyperparameters of our custom network and the MobileNet network. The number of iterations, batches yielded during training, is shown for training on the data set classified by experts.

| Hyperparameters   | Custom net        | MobileNet         |
|-------------------|-------------------|-------------------|
| Batch size        | 10                | 20                |
| Iterations        | 153 000           | 3 825             |
| Optimizer         | SGD               | Adadelta          |
| Nest. Momentum    | 0.90              | -                 |
| Rho               | -                 | 0.95              |
| Initial lr.       | 0.0001            | 1.0               |
| lr. decay         | 10 <sup>-6</sup>  | 0.95              |
| Minimal lr.       | 10 <sup>-4</sup>  | 0.01              |
| lr. red. patience | 14                | 4                 |
| lr. red. factor   | 0.75              | 0.85              |
| Dense dropout     | 0.65              | 0.65              |
| Output activation | softmax           | softmax           |
| Loss function     | cat. crossentropy | cat. crossentropy |
| Input image size  | 356×356           | 224×224           |

layers. Of the 3 models, MobileNet, pre-trained on the ImageNet (Deng et al. 2009), achieved the best performance and therefore we only discuss this architecture. Similarly, Lieu et al. (2018) found MobileNet to be the superior architecture for classifying solar system objects. The hyperparameters for our custom network and the MobileNet network are given in Table 4. We used Keras (Chollet et al. 2015) with TensorFlow (Abadi et al. 2015) backend. The `lr. red. patience` and `lr. red. factor` are parameters of the `ReduceLRonPlateau` Keras callback. The parameter `lr. red. patience` defines how many epochs without improvement of the validation accuracy (different proxy can be chosen to monitor) have to pass to change the current learning rate by multiplying it with the `lr. red. factor`.

The batches used to train the networks were randomly generated during training, always from the whole training sample. Validation started once a satisfying number of generated batches was presented to the network, this is the training data set size divided by the batch size. This was done to maximise the use of our data while keeping balanced numbers of classes in the yielded training batches, in order to avoid biasing the network.

### 5.3 CNN results

We demonstrate that convolutional neural networks are capable of high accuracy, automated galaxy cluster candidate classification. We trained each of our networks 10 times with the exact same hyperparameters, differing only in the seed for generation of random numbers during network’s initialisation, the order of random image selection into balanced mini-batches during training and the random sampling of augmentation values applied during training but keeping the same objects in the training, validation and test data sets. The results of individual runs are averaged and presented together with their standard deviations in Table 5 and Figure 8, helping to compare various networks.

To report accuracy (A), precision (P) and recall (R) in Table 5), we need to make hard classification decision for each example image from our test set. Our neural networks are trained to output the probability that the input image is a galaxy cluster. Therefore, we classify input images as galaxy cluster if this probability is higher than 0.5.

Our best-performing custom network (CN-E), trained on the expert classified data set, achieved an average accuracy of  $(90 \pm 3) \%$ . We also explored training on concatenated PNG images, without the grayscale, so having six channels instead of three, but this did not change the performance significantly.

The MobileNet architecture trained on the data classified by experts achieved an average accuracy of  $(88 \pm 2) \%$ . Perhaps MobileNet has slightly different sensitivity for individual colour channels due to the potential bias in its original training sample. We explored this possibility by training it on two additional channel configurations, X-ray green, optical red, empty blue and X-ray red, optical green, empty blue, but its performance did not change significantly.

Training using the labels obtained in the Zooniverse project resulted in lower performance for both, our custom network (CN-Z) and the MobileNet (MN-Z), achieving average accuracies  $(82 \pm 1) \%$  and  $(79 \pm 2) \%$ , respectively.

Lastly, we also explored the training of neural networks on single wavelength PNG images. Our custom network using expert labels trained only on the X-ray images without their optical counterparts (CN-E solo X-ray) achieved an average accuracy of  $(81 \pm 1) \%$ . Our custom network using expert labels trained only on the optical images (CN-E solo optical) performed the worse, achieving an accuracy of only  $68 \pm 2) \%$ . This is rather easily understandable knowing that the *XMM-Newton* data are much deeper than the POSS-II images used for the current analysis: while *XMM-Newton* can detect galaxy clusters as extended sources out to  $z = 1$  at least, the POSS sensitivity strongly drops beyond  $z \sim 0.3$  rendering galaxies are hardly identifiable.

Using augmentation (subsection 3.2) was critical to achieving good performance, the accuracy of the network CN grayscale would drop from  $(90 \pm 3) \%$  to  $(75 \pm 2) \%$  without the augmentation and from  $(88 \pm 2) \%$  to  $(81 \pm 1) \%$  for MobileNet.

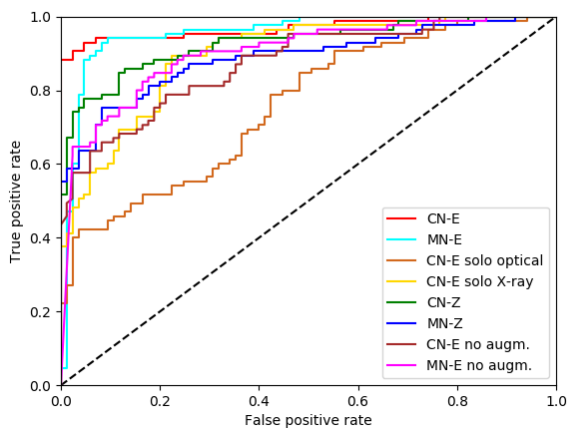
### 5.4 Interpreting the results

We further investigate the results of the best training run of our custom network (CN-E), which can classify even faint clusters and those close to the edge of *XMM-Newton*’s field



**Table 5.** Averaged galaxy cluster candidate classification results of the networks each trained 10 times with the exact same hyperparameters, only with a different seed for generation of random numbers during its initialisation.

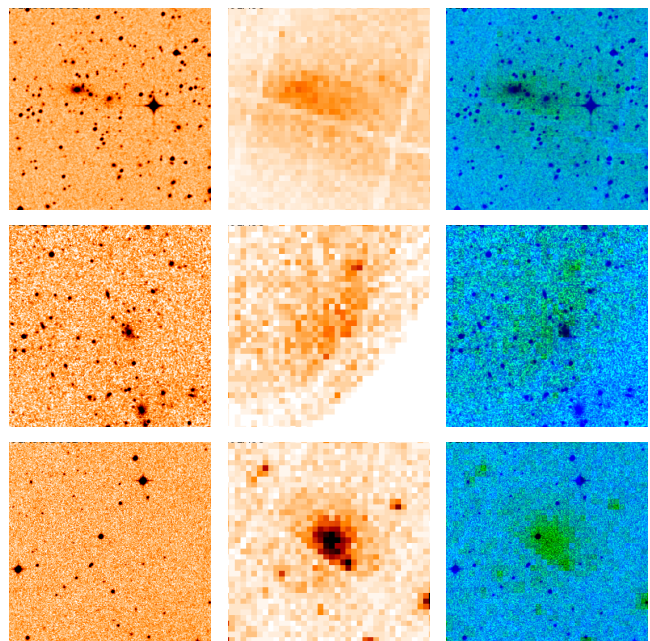
| network           | A $\pm$ std     | P $\pm$ std     | R $\pm$ std     | AUC $\pm$ std   |
|-------------------|-----------------|-----------------|-----------------|-----------------|
| CN-E              | 0.90 $\pm$ 0.03 | 0.89 $\pm$ 0.05 | 0.91 $\pm$ 0.03 | 0.96 $\pm$ 0.01 |
| MN-E              | 0.88 $\pm$ 0.02 | 0.87 $\pm$ 0.03 | 0.91 $\pm$ 0.03 | 0.94 $\pm$ 0.01 |
| CN-E solo optical | 0.68 $\pm$ 0.02 | 0.64 $\pm$ 0.02 | 0.85 $\pm$ 0.04 | 0.77 $\pm$ 0.02 |
| CN-E solo x-ray   | 0.81 $\pm$ 0.01 | 0.78 $\pm$ 0.03 | 0.86 $\pm$ 0.04 | 0.89 $\pm$ 0.01 |
| CN-Z              | 0.82 $\pm$ 0.01 | 0.96 $\pm$ 0.01 | 0.67 $\pm$ 0.02 | 0.91 $\pm$ 0.01 |
| MN-Z              | 0.79 $\pm$ 0.02 | 0.96 $\pm$ 0.03 | 0.62 $\pm$ 0.03 | 0.86 $\pm$ 0.02 |
| CN-E no augm.     | 0.75 $\pm$ 0.02 | 0.70 $\pm$ 0.02 | 0.87 $\pm$ 0.03 | 0.87 $\pm$ 0.01 |
| MN-E no augm.     | 0.81 $\pm$ 0.01 | 0.75 $\pm$ 0.02 | 0.91 $\pm$ 0.01 | 0.90 $\pm$ 0.02 |



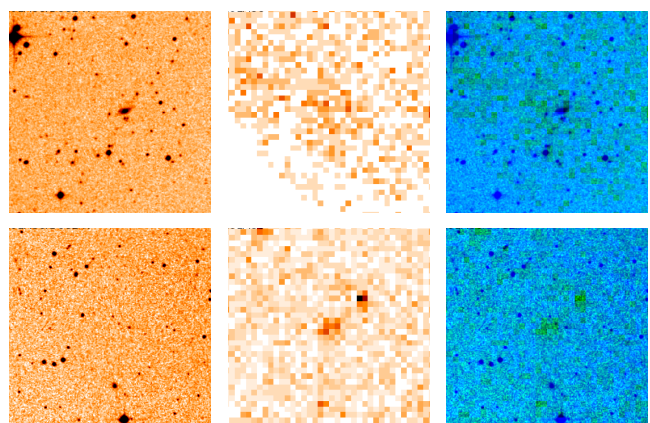
**Figure 8.** ROC curves for the best-performing networks when trained on different data formats. Closer the curve copies the left vertical and top horizontal axis, better the classifier. The dashed line represents how would an untrained, randomly guessing classifier score. Training on optical data only ended up with the poorest results, using only X-ray data achieved much better results, however, the combination of optical and X-ray data resulted in the best performance. CN refers to our custom network, MN to the MobileNet architecture, E to the data set classified by experts, Z to the data set classified by the Zooniverse volunteers.

of view. Figure 9 shows some of these randomly selected correctly classified galaxy clusters.

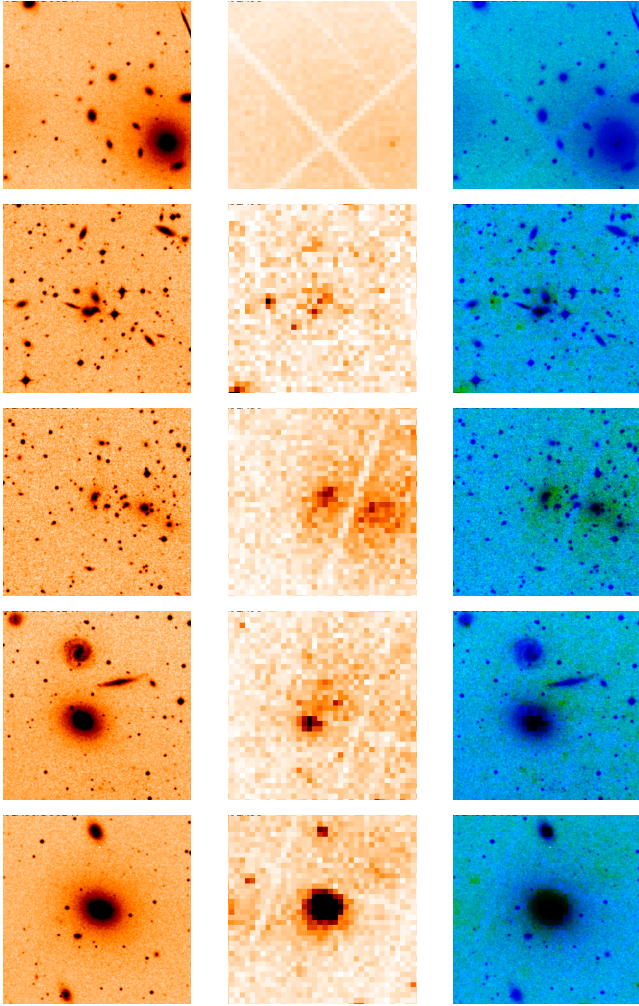
Figure 10 shows two objects classified as non-clusters by the experts, but as clusters by our custom network. The top object raised a concern that it was actually a galaxy cluster. We assume that it was classified as a galaxy cluster by our custom network because of the presence of the faint X-ray emission in the centre and that it is a promising candidate for further investigation and spectroscopic redshift confirmation. Figure 11 displays images of spectroscopically confirmed galaxy clusters which have been incorrectly classified by our custom network as a *non-cluster* class. The first object from the top is a non-centered galaxy cluster. The second contains a group of nearby galaxies with faint extended X-ray emission, which might have fooled our network. The third is a cluster that falls on a chip gap. The fourth is a galaxy cluster with three prominent nearby galaxies along the line of sight which is probably what fooled our network, and the last object appears like a nearby galaxy, which can be hard to classify even for the experts.



**Figure 9.** Spectroscopically confirmed galaxy clusters correctly classified by our custom network randomly selected from the test sample (TP). *Left:* optical, *middle:* X-ray, *right:* combined.



**Figure 10.** Non-galaxy clusters incorrectly classified as galaxy clusters (FP) by our custom network. *Left:* optical, *middle:* X-ray, *right:* combined.

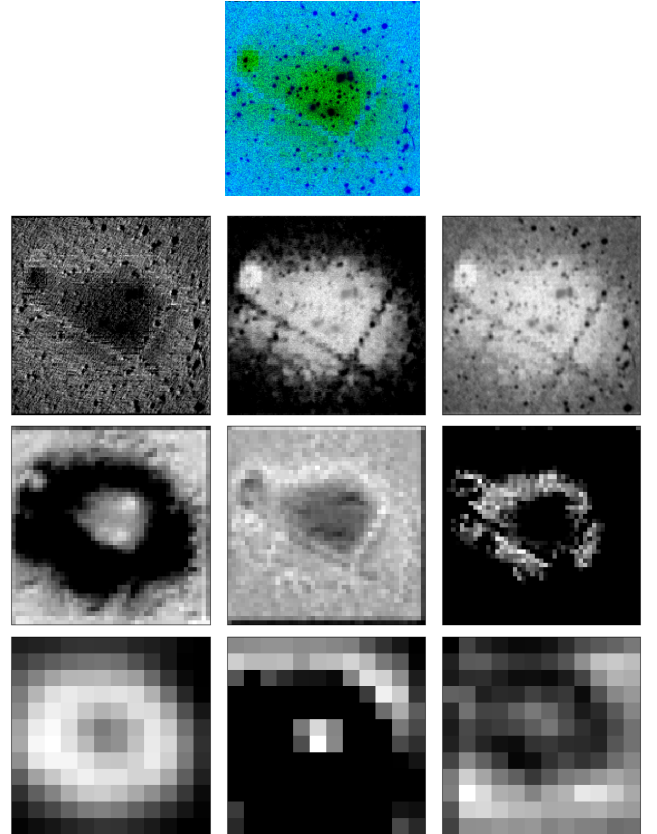


**Figure 11.** Galaxy clusters incorrectly classified as non-galaxy clusters (FN) by our custom network. *Left:* optical, *middle:* X-ray, *right:* combined.

Figure 12 shows outputs of the selected filters of our custom network for a spectroscopically confirmed nearby galaxy cluster. We can see how the network learned to search for edges and colour patches of X-ray or optical light. Some filters learned to search primarily for X-ray emission and others for optical emission. Most of the filters detected both of the emission components simultaneously. Multiple filters in the same layer usually learned to search for X-ray emission, but their sensitivity is different. There are filters which get activated only by stronger emission, while other filters are more sensitive to X-ray emission. The network uses the filters to probe the presence and extent of the X-ray emission in the input image. Note that the filter output size decreases deeper within the network because of the max-pooling operation applied in the pooling layer after each convolutional layer.

### 5.5 Multi-class classification

We also trained neural networks for multi-class classification using the labels of the experts. We segregated objects into 5 classification classes - *low z cluster*, *high z cluster*,



**Figure 12.** Top: Input image to the trained network. Each row from second to last shows outputs (activation maps) of 3 selected filters from 2nd, 4th and 6th convolutional layer of our custom network, respectively.

**Table 6.** Results from the multi-class classification networks.

| class          | A    | P    | R    | AUC  |
|----------------|------|------|------|------|
| MN grayscale   |      |      |      |      |
| Low-z cluster  | 0.77 | 0.62 | 0.94 | 0.93 |
| High-z cluster | 0.87 | 0.56 | 0.22 | 0.91 |
| Point source   | 0.87 | 0.88 | 0.36 | 0.89 |
| Nearby galaxy  | 0.90 | 0.70 | 0.73 | 0.92 |
| Other          | 0.91 | 0.65 | 0.68 | 0.92 |
| CN grayscale   |      |      |      |      |
| Low-z cluster  | 0.79 | 0.68 | 0.81 | 0.89 |
| High-z cluster | 0.84 | 0.44 | 0.65 | 0.89 |
| Point source   | 0.84 | 0.75 | 0.27 | 0.88 |
| Nearby galaxy  | 0.89 | 0.74 | 0.57 | 0.85 |
| Other          | 0.87 | 0.52 | 0.64 | 0.88 |

*nearby galaxy*, *point source* (point, star or AGN, double source) and *other* (artefact, edge). The ROC curves and performance measurements were calculated as one versus all problem.

In this regime, the MobileNet architecture and our custom network achieved an AUC and accuracy, averaged over all classes, within 1 sigma. The MobileNet achieved an AUC score of  $(91 \pm 2)\%$  and accuracy of  $(86 \pm 6)\%$ , and our custom network obtained an AUC of  $(88 \pm 2)\%$  and  $(85 \pm 4)\%$  accuracy (Table 6).

In the case of multi-class classification problems, ROC and AUC are plotted for each of the classes separately as one

**Table 7.** The number of objects in the training, validation and test data sets in a single fold of the 10 fold cross-validation.

| Class       | Experts |          |      |
|-------------|---------|----------|------|
|             | Train   | Validate | Test |
| cluster     | 904     | 113      | 113  |
| non-cluster | 399     | 57       | 114  |

**Table 8.** Classification results of our custom networks for a 10 fold cross-validation on classifications done by experts.

| Fold | A    | P    | R    |
|------|------|------|------|
| 1    | 0.89 | 0.89 | 0.88 |
| 2    | 0.92 | 0.91 | 0.93 |
| 3    | 0.90 | 0.90 | 0.91 |
| 4    | 0.88 | 0.91 | 0.83 |
| 5    | 0.87 | 0.88 | 0.86 |
| 6    | 0.87 | 0.87 | 0.88 |
| 7    | 0.88 | 0.90 | 0.86 |
| 8    | 0.92 | 0.89 | 0.95 |
| 9    | 0.88 | 0.84 | 0.94 |
| 10   | 0.89 | 0.92 | 0.87 |

versus all, reducing the problem to the binary case. From the ROC curves (Figure 13), we see that the *point source* and *high-z galaxy cluster* were the hardest classes to detect, and in the custom network, the *nearby galaxy* class was the easiest to distinguish. We interpret this as a consequence of nearby galaxies being very distinct from the other classes in the optical. Interestingly, this category did not achieve the best performance for the MobileNet network, however, it still placed among the top-performing classes.

We note that since we have trained the neural networks on a sample of galaxy cluster candidates picked by the XAmin pipeline, our sample of point sources is biased towards objects with some spatially extended emission. Thus we can not consider the networks trained for multi-class classification as a reliable point source classifiers since they are not representative of the population and do not reflect the typical appearance of an X-ray point source. If one would like to use our neural networks for point source detection, re-training or fine-tuning of our models on a representative sample of X-ray point sources would be required.

## 5.6 Cross-validation

We perform 10-fold cross-validation of CN-E to explore, if the test data set, having all of its galaxy clusters spectroscopically confirmed, shows significant bias compared to the galaxy cluster sample in the training data set. Table 7 contains the number of example images in each data set for a single fold of the cross-validation. The cross-validation accuracy scores between 87% and 92% (Table 8, Figure 14) and our CN-E achieved accuracy 90% on average (Table 5, Figure 8). Those results are consistent and the test sample we used does not seem to have any significant bias on the network’s performance.

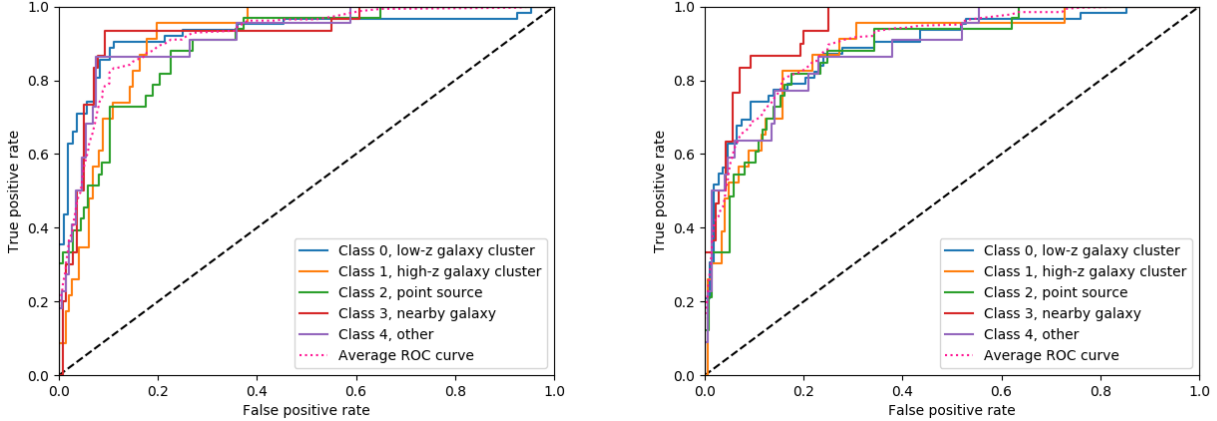
## 6 SUMMARY

In this paper, we have presented convolutional neural networks to classify extended X-ray sources detected by the XAmin pipeline. This automated method can be used to replace the traditional manual screening confirmation task of the XAmin galaxy cluster candidates, which is often tedious and slow.

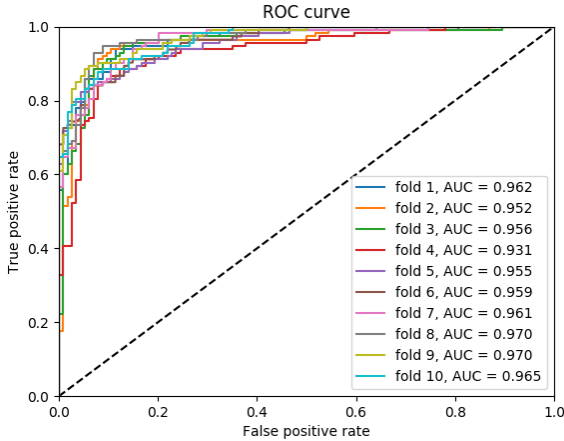
Firstly, we built a crowd-sourcing Zooniverse project - *The Hunt for Galaxy Clusters*, to obtain a classification of a large number (1600) of galaxy cluster candidates in a short time frame (6 months). Our volunteers obtained 62% agreement with experts for identifying clusters and non-clusters in an overlapping sample of 404 objects. We found that the volunteers were often incorrectly classifying objects as point sources or no emission. Out of 254 objects classified as galaxy clusters by experts in the overlapping sample, volunteers agreed on 104 of those (66/146 *low-z* and 38/108 *high-z* galaxy clusters), which is only about 40%, but they inconsistently classified only 1 non-cluster as a galaxy cluster. In total, the volunteers found 506 clusters from 1600 candidates. We suspect the reason behind this low performance of the Zooniverse volunteers in *The Hunt for Galaxy Cluster*, if compared to e.g. Galaxy Zoo, to be the complexity of combined X-ray and optical data of galaxy cluster candidates, burdened by multiple projection and instrumental effects (see subsection 5.1 for discussion of biases the Zooniverse volunteers exhibited). We also tested a hypothesis, that the Zooniverse volunteers would preferentially find prominent galaxy clusters and that their sample could be easily recreated by a cut in the extent – extension likelihood plane (Pacaud et al. 2006), however, the Zooniverse volunteers found galaxy clusters across the entire extent – extension likelihood space (Appendix B), pointing out that their help could be used for a galaxy cluster science.

Next, we trained CNNs on *XMM-Newton* X-ray images combined with their optical counterparts from DSS2, to distinguish galaxy clusters from non-clusters. The cross-validation of our custom network shows consistent results (Table 8, Figure 14) with accuracy scoring between 87% and 92%. We further developed networks on a fixed training, validation and test samples, the networks trained on Zooniverse classified data having a different training and validation samples than those trained on data classified by experts, but both having the same test sample. Our best network (CN-E) obtained an average accuracy of 90% (subsection 5.3). This network used our custom architecture and was trained on labels made by experts. The test sample of 170 objects is composed of 85 spectroscopically confirmed galaxy clusters (62 *low-z* and 23 *high-z*), and 85 galaxy cluster candidates classified as non-clusters by experts. For comparison, a similar network using the MobileNet architecture (MN-E) obtained an average accuracy of 88% and using the custom architecture with the Zooniverse classifications (CN-Z) gave an average accuracy of 82% at best.

In this work, we show that CNNs trained using either X-ray only or optical only images had significantly lower performance in reliably identifying galaxy clusters in comparison to using the combined data. While in the X-rays *XMM-Newton* detects galaxy clusters as extended sources to  $z = 1$  at least, the optical POSS-II data sensitivity strongly drops beyond  $z \sim 0.3$ , making galaxies hardly identifiable.



**Figure 13.** ROC curves for multi-class classification performed by the MobileNet architecture (*left*) and our custom network (*right*).



**Figure 14.** ROC curves for 10 fold cross-validation of our custom network trained on expert classifications.

This is evident from the high number of false-positive detections of galaxy clusters (low precision) using the optical only data. The X-ray only network achieved higher accuracy (81%) than the optical only network (68%).

Additionally we train our networks for multi-class classification using expert classified labels: *low-z galaxy cluster*, *high-z galaxy cluster*, *point source*, *nearby galaxy* and *other*. In this case, the MobileNet architecture performed slightly, but not significantly, better than our custom network (Table 5).

This project is a pilot study to determine the potential of CNNs for the detection of galaxy clusters. In the future, we intend to apply our methods to large sky surveys such as the new *eROSITA* or LSST and *Euclid*. Their enormous data sets are expected to contain tens of thousands of new galaxy clusters, which will require automated, fast and reliable methods to identify, as human screening of such large data volumes will be impossible. Our methods can also be applied to simulated data. Our custom network can be easily fine-tuned to, e.g., *eROSITA* simulations and deliver an au-

tomated search tool for galaxy clusters from X-ray images. Applying our CNN on simulations will also enable modelling of the cluster selection function, important for cosmological studies, which cannot be done with clusters selected by human inspection due to their inconsistent biases.

## ACKNOWLEDGEMENTS

We would like to thank the referee Dr Florence Durret for the valuable comments that helped to improve the paper.

We would like to acknowledge the scientists from the X-CLASS collaboration who manually classified the XAmin galaxy cluster candidates, mainly Jean-Baptiste Melin, one of the moderators overseeing all classifications and Edoardo Cucchetti. We are also very thankful to all of our citizen volunteers who participated in *The Hunt for Galaxy Clusters*. We use data generated via the <http://zooniverse.org> platform, development of which is funded by generous support, including a Global Impact Award from Google, and by a grant from the Alfred P. Sloan Foundation. The Digitised Sky Surveys were produced at the Space Telescope Science Institute under U.S. Government grant NAG W-2166. The images of these surveys are based on photographic data obtained using the Oschin Schmidt Telescope on Palomar Mountain and the UK Schmidt Telescope. The plates were processed into the present compressed digital form with the permission of these institutions. The Second Palomar Observatory Sky Survey (POSS-II) was made by the California Institute of Technology with funds from the National Science Foundation, the National Geographic Society, the Sloan Foundation, the Samuel Oschin Foundation, and the Eastman Kodak Corporation.

We implement our machine learning codes using Keras (Chollet et al. 2015) with a TensorFlow (Abadi et al. 2015) backend, and data augmentation using scikit-learn (Pedregosa et al. 2011). We also used Numpy (Oliphant 2006), Matplotlib (Hunter 2007) and Astropy (Astropy Collaboration et al. (2013), Astropy Collaboration et al. (2018))

Python3 (Van Rossum & Drake 2009) packages. Our codes are open source<sup>2</sup>.

Matej Kosiba is supported by the *European Space Agency* traineeship and the ERASMUS program and travel funding from the ESAC science faculty, Maggie Lieu was supported by *European Space Agency* research fellowship at the European Space Astronomy Centre and a research fellowship at the University of Nottingham.

## REFERENCES

- Abadi M., et al., 2015, TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, <http://tensorflow.org/>
- Ackermann S., Schawinski K., Zhang C., Weigel A. K., Turp M. D., 2018, *MNRAS*, **479**, 415
- Adami C., et al., 2018, *A&A*, **620**, A5
- Allen S. W., Evrard A. E., Mantz A. B., 2011, *ARA&A*, **49**, 409
- Astropy Collaboration et al., 2013, *A&A*, **558**, A33
- Astropy Collaboration et al., 2018, *AJ*, **156**, 123
- Bartlett J. G., Chamballu A., Melin J.-B., Arnaud M., Members of the Planck Working Group 5 2008, *Astronomische Nachrichten*, **329**, 147
- Bertin E., 1993, Science with Astronomical Near-Infrared Sky Surveys: Proceedings of the Les Houches School, Centre de Physique des Houches, Les Houches, France, 20–24 September, 1993. Springer Netherlands
- Bertin E., Arnouts S., 1996, *A&AS*, **117**, 393
- Bishop C. M., 2006, Pattern Recognition and Machine Learning (Information Science and Statistics). Springer-Verlag, Berlin, Heidelberg
- Bleem L. E., et al., 2015, *ApJS*, **216**, 27
- Böhringer H., Chon G., 2016, *Modern Physics Letters A*, **31**, 1640008
- Cavaliere A., Fusco-Femiano R., 1976, *A&A*, **49**, 137
- Chollet F., et al., 2015, Keras, <https://keras.io>
- Ciresan D. C., Meier U., Masci J., Schmidhuber J., 2012, *Neural Networks*, **32**, 333
- Clerc N., Sadibekova T., Pierre M., Pacaud F., Le Fèvre J.-P., Adami C., Altieri B., Valtchanov I., 2012a, *MNRAS*, **423**, 3561
- Clerc N., Sadibekova T., Pierre M., Pacaud F., Le Fèvre J.-P., Adami C., Altieri B., Valtchanov I., 2012b, *VizieR Online Data Catalog*, **742**
- Collister A. A., Lahav O., 2004, *PASP*, **116**, 345
- Deng J., Dong W., Socher R., Li L., Li K., Fei-Fei L., 2009, in 2009 IEEE Conference on Computer Vision and Pattern Recognition. pp 248–255, [doi:10.1109/CVPR.2009.5206848](https://doi.org/10.1109/CVPR.2009.5206848)
- Dieleman S., Willett K. W., Dambre J., 2015, *MNRAS*, **450**, 1441
- E Hinton G., 2007, **11**, 428
- Folkes et al 1996, Highlights of Spanish Astrophysics II. Springer
- Goodfellow I. J., Pouget-Abadie J., Mirza M., Xu B., Warde-Farley D., Ozair S., Courville A., Bengio Y., 2014, preprint, ([arXiv:1406.2661](https://arxiv.org/abs/1406.2661))
- Goodfellow I., Bengio Y., Courville A., 2016, Deep Learning. MIT Press
- Howard A. G., Zhu M., Chen B., Kalenichenko D., Wang W., Weyand T., Andreetto M., Adam H., 2017, preprint, ([arXiv:1704.04861](https://arxiv.org/abs/1704.04861))
- Hunter J. D., 2007, Computing in science & engineering, **9**, 90
- Ichinohe Y., Yamada S., 2019, *MNRAS*, **487**, 2874
- Ivezic Z., Tyson J. A., Abel B., Acosta E., Allsman R., et al. A., 2008, preprint, ([arXiv:0805.2366](https://arxiv.org/abs/0805.2366))
- Jansen F. A., 1999, *ESA Bulletin*, **100**, 9

<sup>2</sup> <https://github.com/matej-kosiba/CNN-multiwavelength-classification-of-X-ray-selected-galaxy-cluster-candidates>

**Table A1.** Threshold values used by the XAmin pipeline, std and median are the standard deviation and the median of the image data.

|         | X-ray       | Optical          |
|---------|-------------|------------------|
| min cut | 0           | median - std     |
| max cut | median × 14 | median + 5 × std |

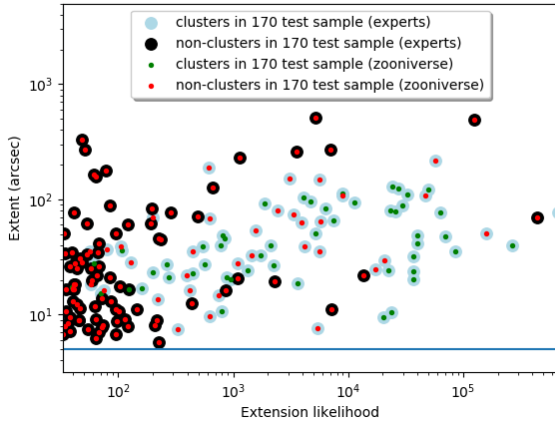
- Krizhevsky A., Sutskever I., Hinton G. E., 2012, in Pereira F., Burges C. J. C., Bottou L., Weinberger K. Q., eds., Advances in Neural Information Processing Systems 25. Curran Associates, Inc., pp 1097–1105, <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-network.pdf>
- LeCun Y., Haffner P., Bottou L., Bengio Y., 1999, in Shape, Contour and Grouping in Computer Vision. Springer-Verlag, London, UK, UK, pp 319–
- Lecun Y., Boser B., Denker J., Henderson D., Howard R., Hubbard W., Jackel L., 1989, *Neural Computation*, **1**, 541
- Lieu M., Conversi L., Altieri B., Carry B., 2018, preprint, ([arXiv:1807.10912](https://arxiv.org/abs/1807.10912))
- Lintott C. J., Schawinski K., Slosar A., Land K., et al. B., 2008, *MNRAS*, **389**, 1179
- Mantz A. B., et al., 2015, *MNRAS*, **446**, 2205
- Mehrtens N., et al., 2012, *MNRAS*, **423**, 1024
- Merloni A., et al., 2012, preprint, ([arXiv:1209.3114](https://arxiv.org/abs/1209.3114))
- Odewahn S. C., Stockwell E. B., Pennington R. L., Humphreys R. M., Zumach W. A., 1992, *AJ*, **103**, 318
- Oliphant T. E., 2006, A guide to NumPy. Vol. 1, Trelgol Publishing USA
- Pacaud F., et al., 2006, *MNRAS*, **372**, 578
- Pacaud F., et al., 2016, *A&A*, **592**, A2
- Pedregosa F., et al., 2011, *Journal of Machine Learning Research*, **12**, 2825
- Pierre M., et al., 2016, *A&A*, **592**, A1
- Racca G. D., et al., 2016, in Space Telescopes and Instrumentation 2016: Optical, Infrared, and Millimeter Wave. p. 990400 ([arXiv:1610.05508](https://arxiv.org/abs/1610.05508)), [doi:10.1117/12.2230762](https://doi.org/10.1117/12.2230762)
- Ridl J., et al., 2017, *MNRAS*, **468**, 662
- Rykoff E. S., et al., 2014, *ApJ*, **785**, 104
- Schawinski K., Zhang C., Zhang H., Fowler L., Santhanam G. K., 2017, *MNRAS*, **467**, L110
- Simonyan K., Zisserman A., 2014, preprint, ([arXiv:1409.1556](https://arxiv.org/abs/1409.1556))
- Starck J.-L., Pierre M., 1998, *A&AS*, **128**, 397
- Starck J.-L., Murtagh F. D., Bijaoui A., 1998, *Image Processing and Data Analysis*
- Szegedy C., Vanhoucke V., Ioffe S., Shlens J., Wojna Z., 2015, preprint, ([arXiv:1512.00567](https://arxiv.org/abs/1512.00567))
- Valtchanov I., Pierre M., Gastaud R., 2001, *A&A*, **370**, 689
- Van Rossum G., Drake F. L., 2009, Python 3 Reference Manual. CreateSpace, Scotts Valley, CA
- Willett K. W., Lintott C. J., Bamford S. P., Masters K. L., et al. S., 2013, *MNRAS*, **435**, 2835

## APPENDIX A: IMAGE PREPROCESSING

The output of the XAmin pipeline is an image with the following normalisation: if a pixel value is lower than min cut, it is attributed a value of 255; if a pixel is greater than max cut it is attributed a value of 0; and  $255 \times (1 - (\text{data} - \text{min cut}) / (\text{max cut} - \text{min cut}))$  otherwise Table A1. To produce the .png images used in the neural networks, XAmin applies the normalisation separately to each of the channels according to Table A2.

**Table A2.** PNG image channel values as constructed by the XAmin pipeline. *pix* refers to the pixel value after cutting.

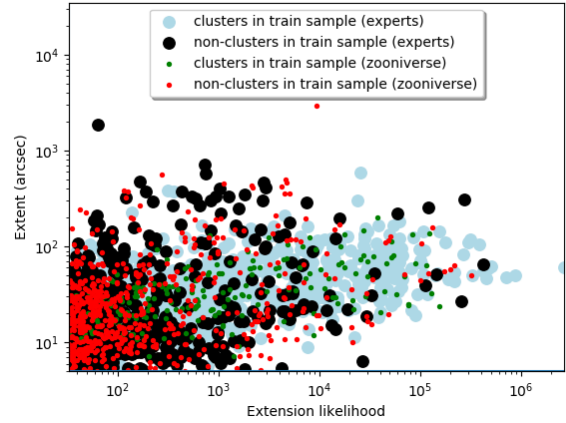
| Channel | Pixel value           | Normalised pixel value                        |
|---------|-----------------------|---|
| R       | $\text{pix} \geq 176$ | 255   |
|         | $\text{pix} < 176$    | $\text{pix} \times 255 / 176$                 |
| G       | $\text{pix} \geq 120$ | $(\text{pix} - 120) \times 255 / (255 - 120)$ |
|         | $\text{pix} < 120$    | 0   |
| B       | $\text{pix} \geq 190$ | $(\text{pix} - 190) \times 255 / (255 - 190)$ |
|         | $\text{pix} < 190$    | 0   |

**Figure B1.** Extent – extension likelihood plane for objects of the 170 test sample classified by experts and the Zooniverse volunteers.

## APPENDIX B: EXTENT – EXTENSION LIKELIHOOD PLANE PLOTS

The extent – extension likelihood plane plots (Figure B1, Figure B2) of our C1 sample of galaxy cluster candidates as described in (Pacaud et al. 2006), were used to analyse the Zooniverse sample of galaxy clusters and investigate our initial hypothesis, that the Zooniverse volunteers will preferentially find most prominent galaxy clusters. We find that the sample of the Zooniverse galaxy clusters span the entire extent – extension likelihood plane and can not be recreated by a simple cut in this space. Please note however that the XAmin v3.5 we used to make the C1 sample had an issue fitting the point source peak, resulting in many non-clusters in the C1 region on the plots and that it is not the same pipeline as the XXL collaboration used before.

This paper has been typeset from a  $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$  file prepared by the author.

**Figure B2.** Extent – extension likelihood plane for objects of the experts train sample and the Zooniverse train sample.