

UC Berkeley

UC Berkeley Previously Published Works

Title

Identifying Transient Candidates in the Dark Energy Survey Using Convolutional Neural Networks

Permalink

<https://escholarship.org/uc/item/2gv579r7>

Journal

Publications of the Astronomical Society of the Pacific, 134(1039)

ISSN

1538-3873

Authors

Ayyar, Venkitesh

Knop, Robert

Awbrey, Autumn

et al.

Publication Date

2022-09-01

DOI

10.1088/1538-3873/ac8375



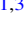

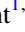
Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

**OPEN ACCESS**

Identifying Transient Candidates in the Dark Energy Survey Using Convolutional Neural Networks

Venkitesh Ayyar^{1,2} , Robert Knop, Jr.¹ , Autumn Awbrey^{1,3} , Alexis Andersen^{1,3} , and Peter Nugent^{1,3} ¹ Lawrence Berkeley National Laboratory, 1 Cyclotron Rd, Berkeley, CA, 94720, USA; vayyar@bu.edu² Hariri Institute for Computing and Computational Science and Engineering, Boston University, Boston, MA, 02215, USA³ Department of Astronomy, University of California, Berkeley, Berkeley, CA, 94720, USA

Received 2022 March 21; accepted 2022 July 20; published 2022 September 7

Abstract

The ability to discover new transient candidates via image differencing without direct human intervention is an important task in observational astronomy. For these kind of image classification problems, machine learning techniques such as Convolutional Neural Networks (CNNs) have shown remarkable success. In this work, we present the results of an automated transient candidate identification on images with CNNs for an extant data set from the Dark Energy Survey Supernova program, whose main focus was on using Type Ia supernovae for cosmology. By performing an architecture search of CNNs, we identify networks that efficiently select non-artifacts (e.g., supernovae, variable stars, AGN, etc.) from artifacts (image defects, mis-subtractions, etc.), achieving the efficiency of previous work performed with random Forests, without the need to expend any effort in feature identification. The CNNs also help us identify a subset of mislabeled images. Performing a relabeling of the images in this subset, the resulting classification with CNNs is significantly better than previous results, lowering the false positive rate by 27% at a fixed missed detection rate of 0.05.

Unified Astronomy Thesaurus concepts: [Type Ia supernovae \(1728\)](#); [Convolutional neural networks \(1938\)](#); [Random Forests \(1935\)](#)

1. Introduction

A major aspect of observational astronomy is the “survey” which involves the wholesale mapping of various regions of the sky to create catalogs which are subsequently mined for scientifically important astronomical objects. We refer to a *transient candidate* as the detection on a single image of a new or varying source with respect to a previously taken reference image, regardless of its astrophysical nature since at this stage its classification is unknown and will remain so until further data is taken (spectroscopy and/or additional photometry). Some examples of such transient candidates are solar system objects, supernovae, active galactic nuclei, variable stars, and neutron star mergers, etc. Since some of these events are quite rare and will fade rapidly, it is often important to trigger follow-up observations immediately to glean their underlying nature and discover new physics. Hence, identifying transient candidates in images quickly and efficiently is very important so as not to waste precious, and expensive, follow-up resources. For many years this process was conducted by manual inspection of images by humans. However, given the magnitude of image data

generated by modern telescopes, it became imperative to automate this process via machine learning techniques. This was first done by the SNFactory (Bailey et al. 2007) where boosted decision trees were employed to greatly reduce the number of candidates. Subsequently, Bloom et al. (2012) advanced this work using random forests to explore the Palomar Transient Factory data for new transients. Surveys such as the Dark Energy Survey (DES) (Flaugher 2005) map the sky both on a large scale and deeply, producing up to 170 GB of raw imaging data every night.

In this work, we describe our efforts to perform transient candidate identification in DES. This work builds on previous work with random forest (Goldstein et al. 2015; Wright et al. 2015; Mahabal et al. 2019) to classify Type Ia supernova from other artifacts of processing and instrumentation for the Dark Energy Survey Supernova program (DES-SN) data. Machine learning techniques such as Convolutional Neural networks (CNNs) (Lecun et al. 1998) have shown remarkable success in image classification problems. Here we apply these to identify transient supernova images.

CNNs have been used for transient candidate identification by multiple groups. The work of Cabrera-Vives et al. (2016) and Cabrera-Vives et al. (2017) were focused on *u*-band imaging from DECam. In Gieseke et al. (2017), they used data from the SkyMapper Survey and their true-positives were



Original content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](#). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

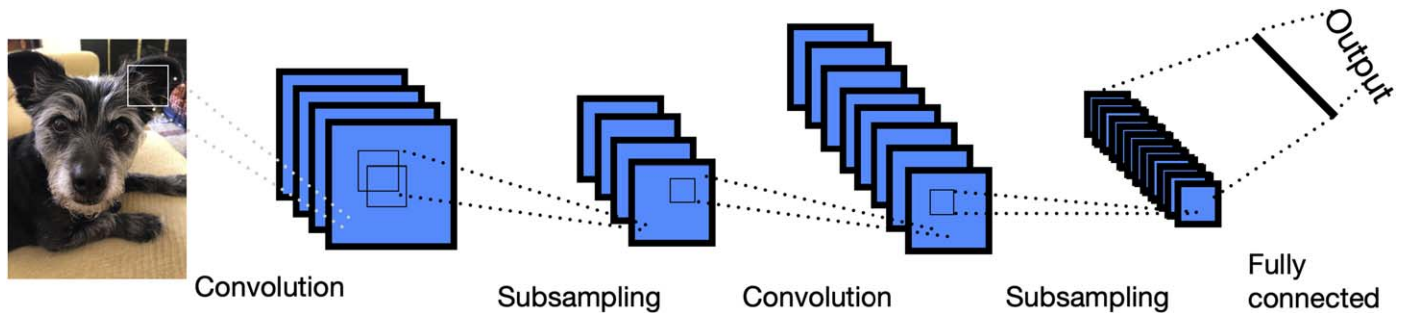


Figure 1. The general structure of a CNN with *convolution*, *subsampling* and *fully connected* layers.

solely drawn from discovered supernovae. Several groups have undertaken the challenge of transient candidate discovery in very wide-field, under-sampled imaging for surveys like TESS, GOTO and GWAC (Jayaraman et al. 2021; Killestein et al. 2021; Turpin et al. 2020).

In Gómez et al. (2020) they use both the spatio and *temporal* data (from a higher cadence survey) to train their CNN’s. This extra data is invaluable for classifying many transient candidates, but lacks the flexibility to work with a single image. We are also aware of another unpublished work with CNNs with a DES data set.⁴ In Acero-Cuellar et al. (2022) they train CNNs for transient detection while focusing on avoiding the use of a difference image in the training process. In Duev et al. (2019), the authors include galactic transients while training their CNNs.

Here, we use the data set found in Goldstein et al. (2015) where stamps were placed on galaxies, drawn from an appropriate redshift range for the survey, and within the galaxy at a location proportional to its surface brightness. This is different than the approach in Cabrera-Vives et al. (2016) and Cabrera-Vives et al. (2017), where the true-positives in their training and validation data were generated by selecting stamps of real PSF-like sources and placing them at a different location at the same epoch and in the same CCD they were observed. This approach works well for *u*-band transients, as potential backgrounds are quite faint (i.e., a host galaxy). However, for a number of transient candidates in an optical search (AGN, supernovae, etc.) this is troublesome as their brightness is comparable to or fainter than their associated galaxies.

After giving a brief description of CNNs in Section 2, we describe our data set in Section 3. In Section 4, we discuss our procedure and present our results. Finally, we summarize our findings in Section 5. Our code is available at the [github repository](#).⁵

⁴ Previous work using CNN’s on DES data can be found [here](#)

⁵ We provide the best saved models and notebooks for plotting and visualization [here](#).

2. Convolutional Neural Networks

Here we give a brief introduction to convolutional neural networks. *Neural networks* are machine learning computing systems that are very efficient at learning patterns in input data. These are generic functions consisting of weight parameters organized in layers. Acting on the input data, after periodic application of nonlinear *activation* functions, they produce outputs which can be either numbers (for regression) or class labels (for classification). By minimizing the deviation between computed output and the expected output, one arrives at the optimal weight parameters. The procedure to compute the weights of a network is called *training the network*. A properly trained network learns the generic function and can correctly predict the outputs for an unseen data set. Essentially, they are universal function approximators.

Convolutional neural networks are a class of neural networks that specialize in recognizing patterns in image data. Using blocks of kernels that scan through the images, they extract features at different scales. Figure 1 gives the general layout of a CNN. A typical CNN is made up of the following basic layers:

1. *Convolution layers*: These perform convolutional operations on the images to extract feature maps.
2. *Subsampling layers*: Operating on feature maps, the subsampling layers compress the dimensionality to reduce the number of parameters.
3. *Fully connected layers*: These layers combine different features of a single layer together.

CNNs typically have very large number of parameters and hence are prone to overfitting. One way to mitigate this is by using dropout layers that help suppress the unimportant weights by setting weight parameters to zero during training.

CNNs have been used extensively for image recognition, classification for images obtained both in the real world and in scientific experiments (Lecun et al. 1998; Ciregan et al. 2012). Large multi-layered Neural networks, despite having large number of free parameters have shown remarkable success in image classification (Szegedy et al. 2015; Krizhevsky et al.

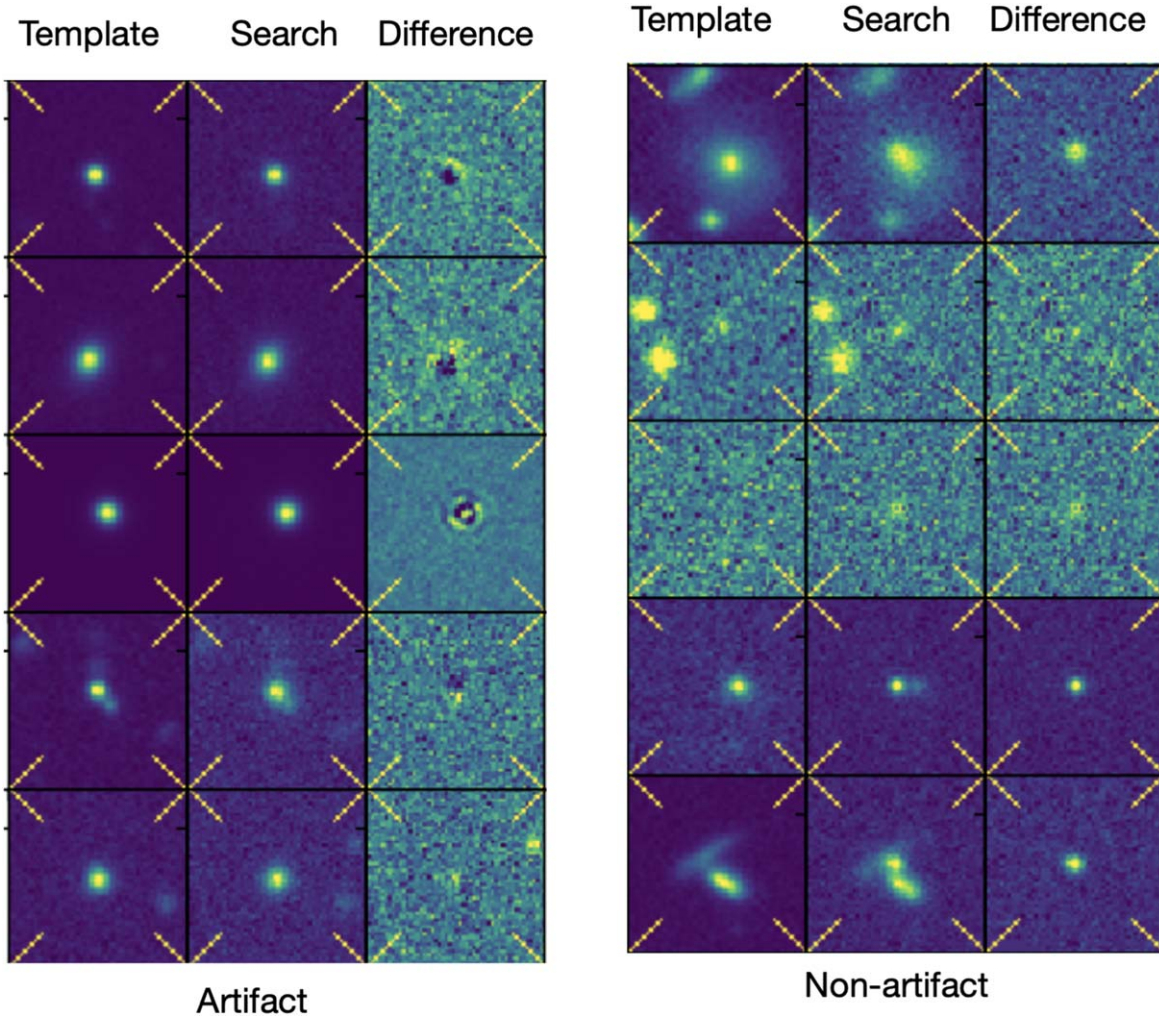


Figure 2. The figure shows the three channels: template, search and difference for five different artifact and non-artifact samples.

2012; Bhimji et al. 2018). Nevertheless, many studies have used specialized CNNs that have connections between non-adjacent layers such as *Resnet* (He et al. 2015) and *Unet* (Ronneberger et al. 2015). In a previous work (Ayyar et al. 2020), it was found that instead of using CNNs with a few but long layers, CNNs with a reasonably high number of layers with fewer parameters were remarkably successful in classifying signal from background for data sets in high energy physics experiments. This prompted us to explore the potential of such layered *deep* CNNs for classification problems in astronomy.

3. Dataset

3.1. Dataset

We used the same data set used in Goldstein et al. (2015). The data collected is from the DES science operations from 2013 August to 2014 February and consists of 898,963 independent samples. Each sample in turn consists of three

types of 2D images of dimension 51×51 . The three types are labeled: *Template*, *Search*, and *Difference*. To incorporate the information from all three images, we use them as channels. In other words, each input sample is a 3-channel image of dimensions 51×51 , having an expected label: Artifact (1) or non-artifact (0). Due to the original timing of the data collection, it lacked non-artifact sources. Hence, the original authors used the method of artificial source construction, and thus injected these non-artifacts into the images. This method has been used extensively before (Bailey et al. 2007; Bloom et al. 2012). For this data set, all non-artifacts were artificially generated.

Figure 2 depicts the three channels for five independent samples for both artifacts and non-artifacts. Distinguishing them visually requires some level of expertise. A more detailed explanation of the data set can be found in Goldstein et al. (2015).

		True labels	
		Non-artifact (=0)	Artifact (=1)
Predicted Labels	Non-artifact (=0)	T_p True positive	F_p False positive
	Artifact (=1)	F_n False negative	T_n True negative

Figure 3. The Confusion matrix for a binary classifier.

3.2. Classification and ROC Curves

In classification problems, the model provides a class prediction for each sample. The aim is to develop a model that categorizes most samples correctly. For a binary classification problem like this one, the performance can be summarized by a 2×2 confusion matrix shown in Figure 3. Some common classifier performance metrics are the True Positive Rate (TPR), False Positive Rate (FPR) and Missed Detection Rate (MDR). They are defined as:

$$\begin{aligned}
 \text{MDR} &= \frac{F_n}{T_p + F_n} \\
 \text{FPR} &= \frac{F_p}{F_p + T_n} \\
 \text{TPR} &= \frac{T_p}{T_p + F_n}
 \end{aligned} \tag{1}$$

where the quantities T_p , T_n , F_p and F_n are defined in Figure 3. Since the classifier prediction is a floating point number between 0 and 1, one uses a *threshold* parameter to determine a predicted class. The behavior of the classifier as the threshold is varied, can be seen through the Receiving Operator Characteristic (ROC) curve. The ROC curve is the most commonly used method to compare a set of classifier models. A useful quantity used for comparison of models is the Area Under the Curve (AUC), which is expected to approach 1.

4. Analysis and Results

4.1. CNN Architecture Search

The goal of this work was to develop optimized CNNs by exploring various CNN architectures. Starting with 4–5 different CNN architectures, we obtained new architectures by varying kernel sizes, dropout layer locations, dropout ratios, types of pooling, learning rates, etc., and compared the classification performance of these models. Picking the best performing models among these, we performed further

Table 1
Table Describing the Best Performing CNN Models

Model Name	Number of Parameters	Area Under Curve (AUC)	Training Time per Epoch on GPU
1	266 k	0.994	60 s
2	415 k	0.994	127 s
3	853 k	0.993	190 s
4	954 k	0.994	47 s

variations and compared their performance. After a few such iterations and studying about 100 different CNN models in total, we shortlisted four models with fairly different architectures that achieved good performance. The structures of these models are listed in Table 4.

Some of the factors guiding the initial architectures were:

1. Presence of Pooling layers (models 1–3 in Tables 4) as opposed to the use of kernel striding (model 4) to reduce image size during convolution.
2. Addition of multiple convolutional and batch normalization layer blocks before implementing pooling layers (models 2, 3).
3. Presence of dropout layers after batch normalization layers.

We explored kernel sizes from 2 to 8 (since the image size is 51×51) and convolutional layer sizes from 10 to 400. As we approached models with good classification performance, we lowered the learning rate to obtain better classification. The entire code was written in python using the keras package (Chollet 2015). We used the numpy library (Harris et al. 2020) for computations and jupyter notebooks (Kluyver et al. 2016) for analysis and visualizations. While we do not claim to have explored every architecture, our search is reasonably thorough and we do find multiple CNN models that are efficient for the given classification problem.

4.2. Results with Original Labels

For the first part of the work, we split the data into *training* (50%), *validation* (5%) and *test* samples (5%). The validation data was used to assess the classification performance of trained models on unseen data. The test data was used to compare the ROC curves of the different models. At this stage, we kept about 40% of the data in reserve for further analysis. We also trained a random forest using the hyperparameters described in the Goldstein et al. (2015).

Table 1 compares the four best CNN models. All models have an AUC score close to 1.0. Figure 4 shows the FPR-MDR ROC curve for the best CNN models. The black squares represent the ROC curve of the random forest from Figure 7 of Goldstein et al. (2015). Here we see that our CNN models and

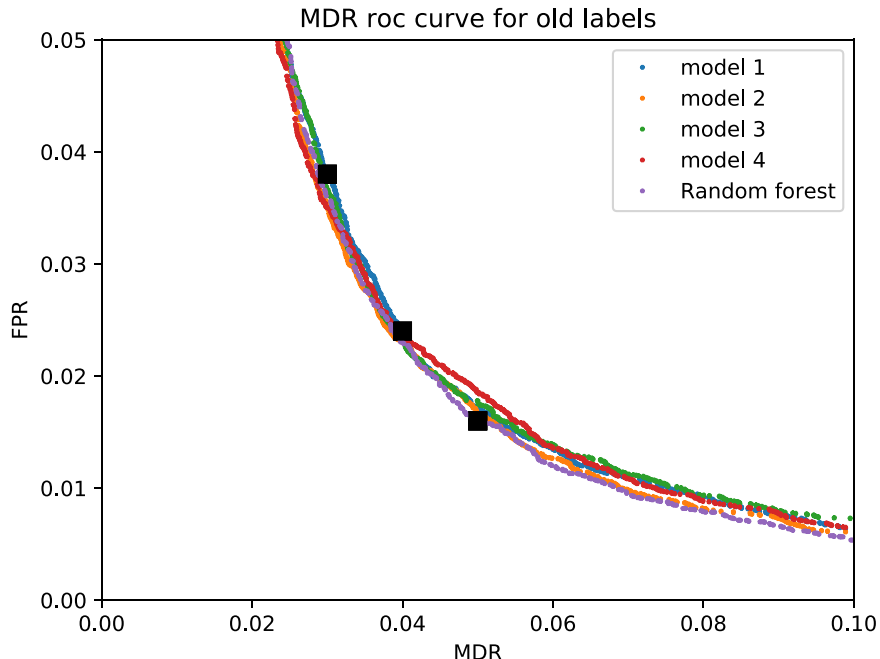


Figure 4. The ROC curve of FPR vs. MDR for the CNN models. The black squares show the points obtained in Goldstein et al. (2015) with a random forest. The ROC curves of all the models are adjacent to each other implying similar classification performance.

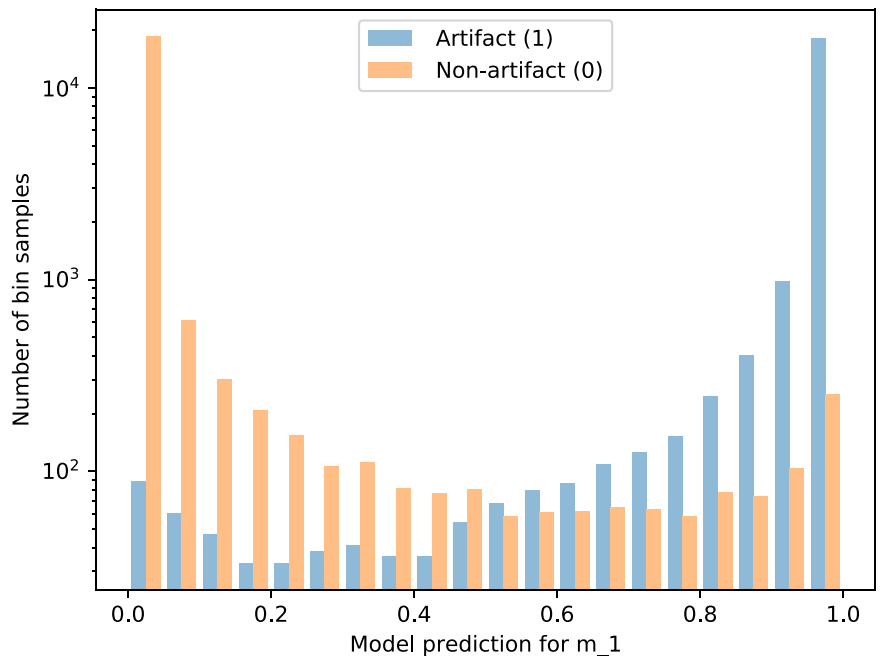


Figure 5. The prediction histograms for Model 1. Ideally, the classifier should have a prediction value 0 for all artifacts and 1 for all non-artifacts. Note the log scale on the y-axis. Most of the labels are predicted correctly. Also, there are relatively very few predictions in the intermediate region.

the random forest are comparable in performance to the previous work.

Figure 5 shows the prediction histogram for Model 1. Noting the log-scale on the y-axis, it is clear that most of the samples are

classified correctly as either artifacts or non-artifacts. Also, very few samples have prediction scores in the intermediate 0.2–0.8 range, which is the hallmark of a good classifier. The prediction histograms for the other three models also look very similar.

Table 2
Dividing the Samples Into Categories, Depending on Model Predictions

Category	Original Label	Prediction range	Description	Model 1	Model 3	Random Forest
1	1	0-0.1	Strongly mis-classified Artifact	0.35%	0.76%	0.15%
2	1	0.1-0.5	Weakly mis-classified Artifact	0.9%	1.2%	1.0%
3	1	0.5-1.0	Correctly classified Artifact	48.2%	47.7%	48.4%
4	0	0-0.5	Correctly classified Non-artifact	48.4%	49%	48.4%
5	0	0.5-0.9	Weakly mis-classified Non-artifact	1.2%	0.8%	1.6%
6	0	0.9-1.0	Strongly mis-classified Non-artifact	0.8%	0.7%	0.4%

Note. Categories 1 and 6 correspond to the strongly mis-classified samples. The CNN models have more strongly mis-classified samples.

Table 3

The Table on the Left Shows How the Different Models Categorize the 149 Samples of a Test Data Set That are in Category 1 (Artifact Classified as Non-artifact) for Model 1

Points in category 1 for model 1			
Category	1	2	3
CNN 2	65.8%	20.1%	14.1%
CNN 3	72.5%	18.8%	8.7%
CNN 4	72.5%	17.4%	10.1%
Random forest	26.8%	45.0%	28.2%

Points in category 6 for model 1			
Category	4	5	6
2	7.9%	18.4%	73.7%
3	12.4%	11.9%	75.7%
4	9.9%	29.3%	60.7%
Random forest	15.8%	44.6%	39.5%

Note. The CNNs place over 66% of these samples in category 1. The random forest places only about 27% of these, instead placing many of the rest in category 2. Similarly, the table on the right shows how the different models categorize the 354 samples of a test data set that are in category 6 for Model 1. From these, it can be inferred that, the different CNN models all strongly misclassify the same set of points.

4.3. Mislabeled Images

Since the CNNs use the entire information from the images, one would expect well trained CNNs to achieve optimal classification performance. While the CNNs in Figure 4 do perform very well, the fact that they do not improve upon the performance of the random forest, prompted us to explore their classification in more detail.

In the bottom left part of Figure 5, it can be seen that a significant number of artifacts with input label 1 are accorded a prediction value very close to 0. Similarly, a large number of non-artifacts with input label 0 have prediction values close to 1, as seen in the bottom right. We see a similar pattern for the other CNN models. This seems to imply that the CNN models are strongly mis-classifying a few samples, thus affecting the quality of their ROC curves.

To better understand this issue of strong mis-classification, we divided the samples into six categories depending on the original label and the predicted value for the model. For

example, category 1 corresponds to samples that are labeled as artifacts (label = 1), but have prediction values between 0 and 0.1. We then compare the number of samples in the different categories for the different models. The results and description of the categories are summarized in Table 2. It can be seen that the CNN models have more points in categories 1 and 6, corresponding to strongly mis-classified samples.

To better understand this, we look at how the samples classified by Model 1 into category 1 are categorized by other models. In Table 3, we show the prediction values for the other three models and random forest, for that were classified by Model 1 to be in categories 1 and 6. It can be seen that about 66% of the samples are also placed in the same category by the other 3 CNN models. However, the random forest model places a much smaller proportion of these samples in category 1. This is further confirmation that the four CNN models seem to strongly mis-classify the same set of images.

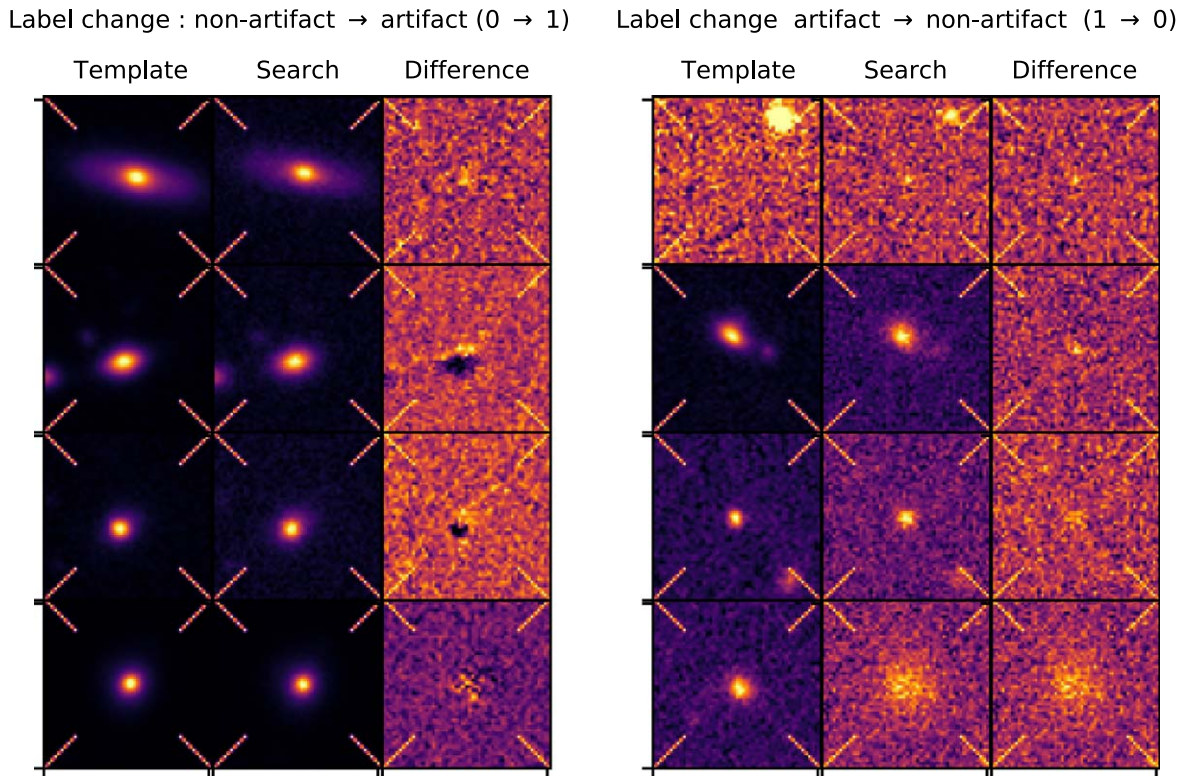


Figure 6. The left figure show four samples that were incorrectly labeled as non-artifacts, while the right figure shows four samples that were incorrectly labeled as artifacts. In almost all cases of a True-positive being mislabeled, it is due to the image stamp being placed on a saturated star or galaxy or on a bad part of the CCD. The mislabeling of the false negatives is well understood, and expected, as this sample in Goldstein et al. (2015) was taken from all candidates and some true astrophysical transients crept into the data.

4.4. Relabeling

The fact that different optimally trained CNN models strongly mis-categorize the same set of samples, points to the possibility of a case of mislabeling. To confirm this, we went back to the set of images in categories 1 and 6 for Model 1. Upon performing a visual inspection for a small subset of these samples, we found that about 90% of these samples were indeed given the wrong label during the process of data preparation. As the artifact sample was randomly drawn from all detections, while the non-artifact came from injected transient candidates this is not too surprising. In fact, what we saw in our visual inspection is that some of the injected non-artifacts fell on bad parts of the CCD detector or were located near saturated stars while several of the artifacts were in fact heretofore unknown astrophysical transients. This can be seen in Figure 6, where we show falsely classified non-artifacts in the left figure and falsely classified artifacts in the figure to the right.

For the purposes of relabeling, we developed a GUI tool in python using the *Tkinter* (Van Rossum 2020) library that enabled us to view blocks of images with their labels and allow an expert to quickly mark the images that require relabeling.

As a first step, we performed a relabeling by inspecting roughly 750 samples that were classified into categories 1 and 6 by Model 1. Using the same trained models and their predictions, and just using the new labels, the resulting ROC curves are shown in Figure 7. It is clear that the models are doing significantly better with the newer labels. In addition, the CNN models are now performing better than the random forest. Thus despite all models being trained with a data set containing some mislabeled points, the CNNs are doing a better job at classification.

Having convinced ourselves of the efficacy of the relabeling process, we then obtained the predictions of Model 1 on the entire data set. Collecting the samples in categories 1 and 6, we then inspected this subset of 8093 samples. We found that 7402 (91%) of these had to be relabelled. In all, 0.8% of the samples were relabelled (7402 out of 898,963 samples).

4.5. Results with New Labels

We trained all four CNN models and the random forest on the relabeled data set, after splitting the data into training (70%), validation (10%), and test samples (20%). The resulting ROC curves are shown in Figure 8. It is clear that the random

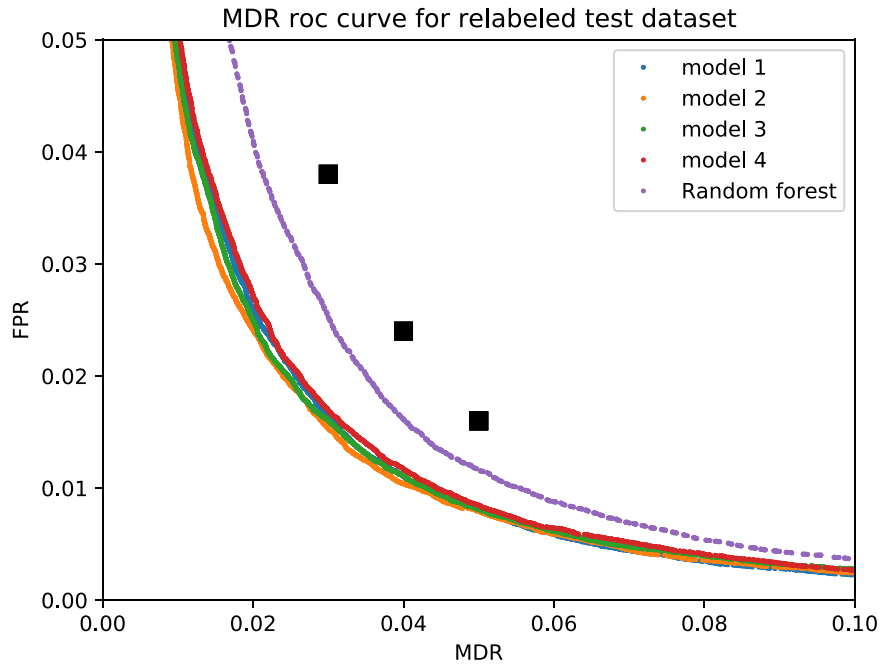


Figure 7. The ROC curve of FPR vs. MDR for the previously trained models using new labels for the test data set. The black squares show the points obtained in Goldstein et al. (2015) with random forest. The CNN models and random forest show significant improvement.

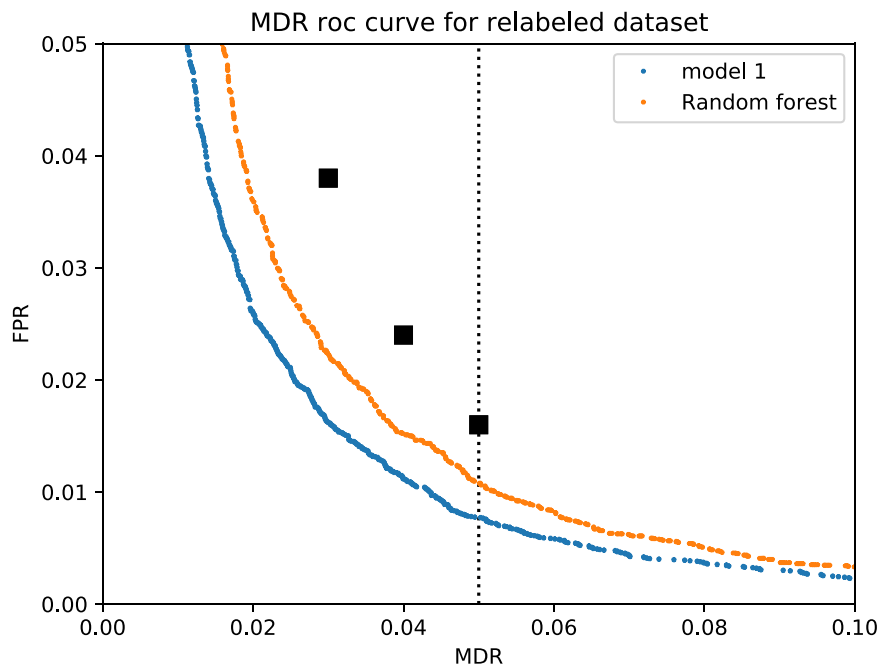


Figure 8. The ROC curves of the best chosen CNN model: model 1 and random forest that were trained on the relabeled data set. The black squares show the points obtained in Goldstein et al. (2015) with the random forest. It is clear that Model 1 outperforms the random forest even on the relabeled data set. The dotted line represents an MDR value of 0.05. For this MDR value, the correspond FPR values of model1 and random forest are 0.008 and 0.011 respectively. Thus the FPR value of the CNN model is only 73% of the FPR value of the random forest.

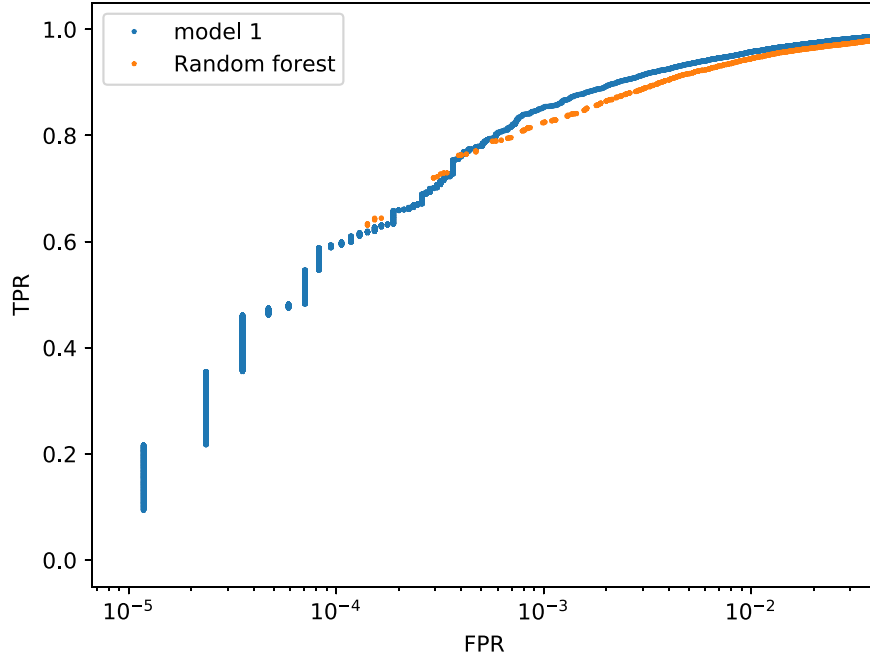


Figure 9. The figure shows the ROC curves for Model 1 and random forest, plotting the true positive rate (TPR) with the false positive rate (FPR).

forest is performing better with the new, relabeled data set, and the four CNN are comparable in performance, but better than the random forest. Figure 9 shows the ROC comparing the true positive rate with the false positive rate as defined in Equation (1). Based on the two ROC curves, we choose Model 1 as our best model, although the classification of the four CNN models are fairly similar. To quantify the improvement in performance, we compare the FPR values for a fixed MDR value of 0.05 as shown in the Figure 8. The corresponding FPR values for the CNN model 1 and random forest at 0.008 and 0.011 respectively. Thus the CNN model 1 lowers the FPR value by 27%.

We present the learning curve of Model 1 in Figure 10 and its detailed structure in Table 4. The confusion matrix is presented in Figure 11.

We would like to reiterate that the relabeling procedure has been conducted by visual inspection, with the machine learning method being used to only shortlist the suspected mislabeled images. It might be argued that our relabeling process might be biasing the performance of Model 1, since we chose the points to relabel, based on the predictions of Model 1. However, in Section 4.2, we trained with only 50% of the samples, keeping the rest of the data in reserve. For this final analysis, we built the data set with a different random seed and used 70% of the data for training, thus mitigating the bias. Table 5 shows the number of points in various categories for Model 1 for the three cases: train-test with old labels, train with old label, test with new labels, train and test with new labels. It is clear that values in column 2 are lowest, due to the bias. But the values in

column 3 are intermediate, indicating that the bias has been mitigated. Hence, we are confident that our final CNN models are indeed better at classification than the random forest.

4.6. Results in an Ongoing Survey

We are currently using the CNN to provide a real/bogus score for an ongoing survey, the DECam Deep Drilling Fields (DDF) program (M. L. Graham et al., 2022, in preparation). This is being run at the Blanco 4 m telescope at Cerro Tololo-Inter-American Observatory as part of the DECam Alliance for Transients (DECAT), a consortium of time-domain DECam programs. We have been using real/bogus scores produced by the CNN to decide which detections are sent out in alerts. Throughout 2021, we used the original CNN trained for this paper. Starting in 2022, and ultimately for the analysis of the entire data from the survey, we will be using a retrained CNN as described below for extragalactic ($|b| > 20^\circ$) events. Training of a CNN for galactic events is in progress.

Because this is real, incoming data, rather than a simulation, we do not have the absolute truth as to what’s a genuine astronomical detection on the difference image, and what’s a subtraction artifact or other “bogus” event. In order to produce an evaluation data set, several observers manually tagged events from this survey as “real” or “bogus”. Participating observers were all trained on examples of good and bad events. The observers included two with decades of experience in vetting supernova candidates in searches like this, and three undergraduate student assistants. Each observer was given

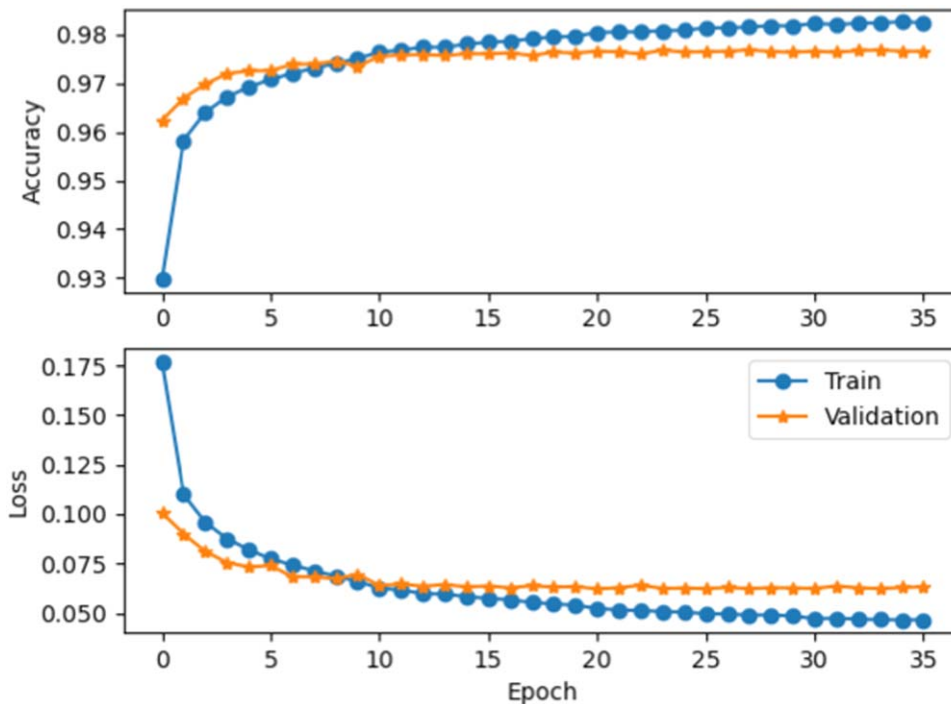


Figure 10. The figure shows the learning curves for CNN Model 1. The top figure provides the training and validation accuracy values at the end of each epoch, while the bottom figure provides the corresponding loss values. The validation loss and accuracy stabilize after about 10 epochs.

events randomly chosen from the few million events that had been found by the data pipeline. Once several observers had rated enough events, they were given events randomly chosen from those that had already been rated by others. In this way, we were able to build up a set of $\approx 25,000$ events that had been tagged by three or more observers.

It is important to emphasize that what we are trying to do here is different from what is described in the rest of this paper. The development of the CNN described in most of this paper aimed to correctly identify simulated transient candidates. Here, we are using the same CNN architecture in an attempt to reproduce in bulk, the messier process of human scanning of transient candidates from real data. The training and validation sets cannot be as clean in this case. Of the 25,000 events tagged by three or more observers, we selected those where the number of observers in the majority was at least two greater than the number in the minority. (Effectively, this means that for those only rated by three observers, the tags would have been unanimous.) Of this subset, about 1700 were tagged by the majority as good and 19,000 as bad. There was a unanimous agreement on 75% of the events that the majority deemed to be good; there was a consensus on 95% of the events deemed to be bad by the majority. The choice to use the majority tags rather than the consensus tags represents a greater emphasis on reducing missed detections as opposed to reducing false positives.

The CNN trained on the simulated transient candidates (whose results are in Section 4.5) did not perform particularly well in reproducing human scanning of the live data set. In particular, there was a high missed detection rate (for any reasonable r/b cutoff) of ~ 0.5 for candidates that were unanimously agreed to be good by three or more observers; this high MDR was present even when limiting to events with a high S/N ratio. To allow the CNN to better model the visual scanning of this survey, we re-trained the model using the majority-tagged events described above as a training and validation set. This retrained model performed much better than the original model on this new data set, yielding a MDR of ~ 0.055 and a FPR of ~ 0.04 (similar to the performance of the originally model). The results for the retrained CNN are shown in Figures 12 and 13. We cannot expect the ROC curve (left plot of Figure 12 here to be as good as the ones in Figure 8 because, as mentioned above, the training and validation set is not nearly as clean.

There is some evidence that the non-consensus events were at least sometimes more marginal cases, based on the r/b scores produced by the CNN retrained on the majority rankings. The training of the CNN was just given a single real or bogus flag based on the majority of the human rankings; it had no knowledge of what was a consensus-good or consensus-bad event. Despite this, the retrained CNN showed different statistics for consensus versus non-consensus events.

Table 4
Structures of the 4 Best CNN Models

Model 1		
Layer	Output Shape	No. of Parameters
Input	$51 \times 51 \times 3$	0
Conv2D	$51 \times 51 \times 80$	2240
BatchNorm	$51 \times 51 \times 80$	320
MaxPooling	$25 \times 25 \times 80$	0
Conv2D	$25 \times 25 \times 80$	57680
BatchNorm	$51 \times 25 \times 25$	320
MaxPooling	$12 \times 12 \times 80$	0
Conv2D	$12 \times 12 \times 80$	57680
BatchNorm	$12 \times 12 \times 80$	320
MaxPooling	$6 \times 6 \times 80$	0
Flatten	2880	0
Dropout	2880	0
Dense	51	146931
BatchNorm	51	204
Dense	1	52

Total trainable parameters **265,165**

Model 2

Layer	Output Shape	No. of Parameters
Input	$51 \times 51 \times 3$	0
Conv2D	$51 \times 51 \times 80$	3920
BatchNorm	$51 \times 51 \times 80$	320
Conv2D	$51 \times 51 \times 80$	102480
BatchNorm	$51 \times 51 \times 80$	320
MaxPooling	$17 \times 17 \times 80$	0
Conv2D	$17 \times 17 \times 80$	102480
BatchNorm	$17 \times 17 \times 80$	320
Conv2D	$17 \times 17 \times 80$	102480
BatchNorm	$17 \times 17 \times 80$	320
MaxPooling	$5 \times 5 \times 80$	0
Flatten	2000	0
Dropout	2000	0
Dense	51	102051
BatchNorm	51	204
Dense	1	52

Total trainable parameters **414,205**

Model 3

Layer	Output Shape	No. of Parameters
Input	$51 \times 51 \times 3$	0
Conv2D	$51 \times 51 \times 120$	5880
BatchNorm	$51 \times 51 \times 120$	480
Conv2D	$51 \times 51 \times 120$	230520
BatchNorm	$51 \times 51 \times 120$	480
MaxPooling	$17 \times 17 \times 120$	0
Conv2D	$17 \times 17 \times 120$	230520
BatchNorm	$17 \times 17 \times 120$	480
Conv2D	$17 \times 17 \times 120$	230520
BatchNorm	$17 \times 17 \times 120$	480
MaxPooling	$5 \times 5 \times 120$	0
Flatten	3000	0
Dropout	3000	0

Table 4
(Continued)

Model 1		
Layer	Output Shape	No. of Parameters
Dense	51	153051
BatchNorm	51	204
Dense	1	52

Total trainable parameters **851,605**

Model 4

Layer	Output Shape	No. of Parameters
Input	$51 \times 51 \times 3$	0
Conv2D	$26 \times 26 \times 40$	4360
BatchNorm	$26 \times 26 \times 40$	160
Dropout	$26 \times 26 \times 40$	0
Conv2D	$13 \times 13 \times 60$	86460
BatchNorm	$13 \times 13 \times 60$	240
Dropout	$13 \times 13 \times 60$	0
Conv2D	$13 \times 13 \times 80$	172880
BatchNorm	$13 \times 13 \times 80$	320
Dropout	$13 \times 13 \times 80$	0
Flatten	13520	0
Dropout	13520	0
Dense	51	689571
BatchNorm	51	204
Dense	1	52

Total trainable parameters **953,785**

Note. As mentioned in 3.1, each CNN reads a batch of input images of dimensions $51 \times 51 \times 3$, with the 3 corresponding to the three types of images *Template*, *Search* and *Difference*. As the different layers of the CNN are applied to each image, the dimensions of the image array change. The above table lists these details, with the first column describing the layer and the second column denoting the dimensions of the intermediate image array. The third column gives the number of parameters in each layer. The terms Layers Conv2D, Batch Norm, MaxPooling, Flatten, Dropout and Dense represent the standard CNN operations convolution, batch normalization, maxpooling, flattening, dropout and dense respectively. More information about these can be found in the keras layers API [documentation](#).

For the majority-good events, the r/b scores of the consensus-good events were on average higher than the events on which there was disagreement (0.87 versus 0.75). For the majority-bad events, the r/b scores of the consensus-bad events were on average lower (0.05 versus 0.25). Note that the fraction of events that had a consensus did not appreciably change when limiting to only high S/N events; those events that were marginal cases were not simply low-S/N cases, but represented cases where visual appearance of the residual might have had some suggestion of being an artifact, but was not clearly an artifact.

In conclusion, the model obtained using a CNN architecture optimized for the original data set works fairly well with

		True labels	
		Non-artifact (=0)	Artifact (=1)
Predicted Labels	Non-artifact (=0)	97.8% (81916)	2.4% (1855)
	Artifact (=1)	2.2% (2002)	97.6% (82864)

Figure 11. The confusion matrix, with both normalized percentages and total number of triplets for used in training (in parens), for Model 1 based on a threshold cut of 0.5.

Table 5

Comparing the Points in various Categories for the three Cases: 1: Train and Test with old Labels, 2: Train with old Labels, but Test with new Labels, 3: Train and Test with new Labels

Category	1: Train and Test with Old Labels	2: Train with Old Labels, Test with New Labels	3: Train-test with New Labels
1	0.35%	0.2%	0.36%
2	0.75%	0.82%	0.83%
3	48.4%	49%	49.1%
4	48.4%	48.8%	48.5%
5	1.2%	1.02%	0.77%
6	0.84%	0.23%	0.33%

Note. Comparing columns 1, 2 and 3, it can be seen that column 2 has few points in categories 1 and 6, since the test data set was biased by Model 1. Column 3 has slightly higher values in these categories compared to column 2, since the models were re-trained on a bigger data set. Overall, column 3 has fewer points than column 1 in category 6, implying that the relabeling procedure had the intended effect.

another, similar data set, after a re-training of the weights. This points to our model architecture being fairly generic and hence more broadly applicable for data sets of this type. To improve upon the performance on this new data set, we would have to perform another architecture search with a subset of this new data. We aim to address this in a future publication.

5. Conclusions and Discussion

5.1. Inference

We have discussed automating the identification of transient detections obtained in astronomical imaging data using machine learning. Here we developed CNN models trained directly on the raw image data. The best CNN models match the performance of the previously used random forest method. In addition, using the CNNs predictions, we were able to identify that some of the images were mislabeled in the original data. After performing a relabeling of 0.9% of the data set, we re-trained the best CNN models. The resultant models outperform the original random forest method. We also find that the CNNs are more robust to mislabeled samples in the training data.

Since we have only relabeled a small subset of the data, there is still the possibility that many other points are mislabelled. However, the significant increase in classification performance suggests that we have identified and relabelled most of the mislabeled images.

5.2. Discussion

There are two main benefits of using CNNs over random forest for image classification: classification efficiency and ease of use.

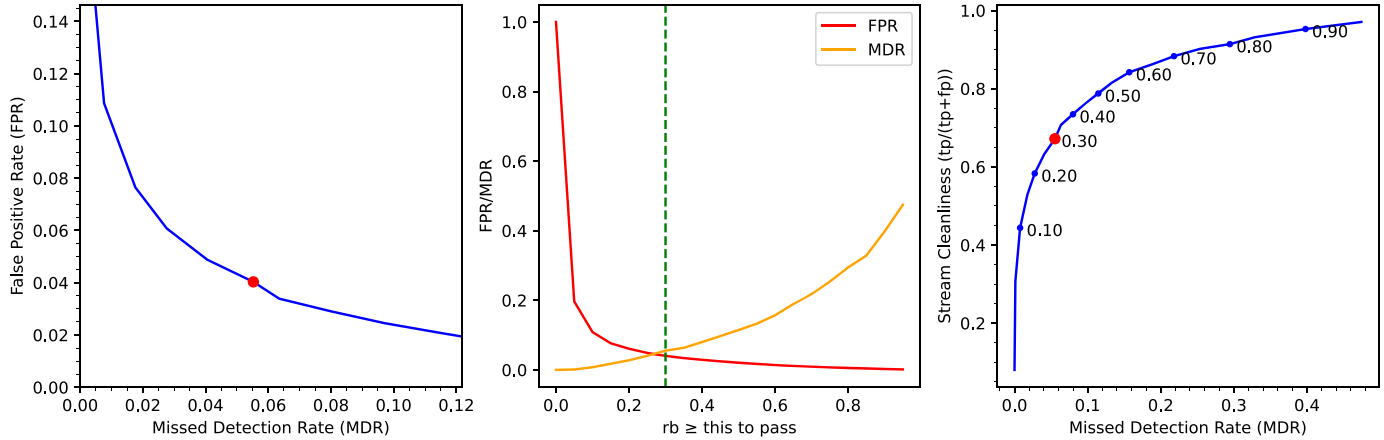


Figure 12. Results for the CNN model trained against manual vetting of the ongoing DECAT/DDF survey. Left: ROC curve. Middle: false positive rate (FPR) and missed detection rate (MDR) as a function of real/bogus score cutoff. Right: stream cleanliness vs. MDR. Stream cleanliness is the fraction of passed objects that are real objects. This is different from $1 - \text{FPR}$ because in the real data set, there are a factor of 7 more bogus events than real events. In the left and right plots, the red dot indicates the chosen real/bogus threshold of 0.3 that will be used in determining if an alert should be sent out for the detection. The dashed green vertical line in the middle plot is the same cutoff.

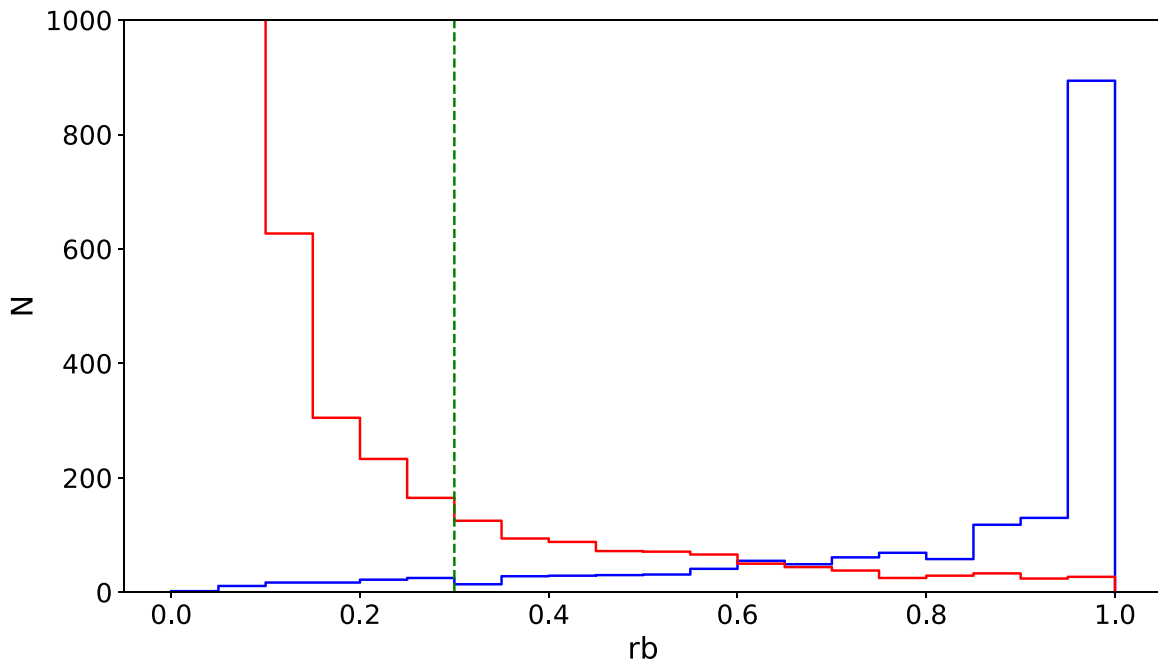


Figure 13. Histogram of real/bogus produced by the retrained CNN for the ongoing DECAT/DDF survey. The blue histogram are events labeled as “real” by at least two out of three visual inspections, and the red histogram are events labeled as “bogus”. The vertical dashed green line is the real/bogus threshold of 0.3 that will be used for generating alerts. Compared to Figure 5; as discussed in the text, we cannot expect the histogram to be as cleanly separated with this data set as we can for the data set used for the bulk of the paper.

Currently in astronomy, random forest methods are the most common method to auto-identify transients in image subtractions. In Wright et al. (2015), the authors compared the performance of random forests with neural networks and demonstrated that their random forest was most efficient. However, in our work, we have demonstrated that our CNN,

obtained by performing a detailed architecture search, outperforms our random forest. For example, comparing the FPR value for a fixed MDR, we find that the CNN model lowers the FPR by 27% compared to the random forest. Since the data sets are different, one cannot perform a direct comparison of ROC curves. However, the fact that our CNN outperforms the our

random forest on the same data set clearly demonstrates the benefit of using CNNs.

Recently Acero-Cuellar et al. (2022) used CNN's for transient discovery that is the most comparable to our work as both used the data sets from Goldstein et al. (2015) to train the CNN's. While the focus of their paper is on performing transient discovery without image subtraction, they do present a confusion matrix for a similar design as ours. Both their false-positives and true-negatives are a factor of ~ 2 larger than ours.

Duev et al. (2019) has created CNN's for the classification of transients in the Zwicky Transient Facility (ZTF). It should be noted that ZTF is slightly different than our survey in that $\sim 30\%$ of the images they take are undersampled (Bellm et al. 2018), thus a direct comparison to our work is not exactly correct. That said, their confusion matrix is very similar to ours in quality (their true-positive's are slightly more pure while their false-negatives are slightly worse). As the number of validation triplets was low in their study, the uncertainties on these numbers are on the order of $\sim 2\%$. Both these studies confirm the utility and benefit of using CNNs for transient candidate detection and that they are superior to random forest methods.

Another benefit of CNNs is their ease of implementation. One of the drawbacks of the original random forest method is the need to identify a set of important features to use. This process is fairly painstaking, and also involves performing many, often computationally expensive, operations on the raw images. The CNN method, on the other hand works directly on the raw image data. Although it requires an architecture search, different models with reasonably high complexity perform well in classification. Their computational cost is also quite reasonable as they run efficiently on GPUs.

Hence, we believe the CNN method is more suitable for implementing automation of image subtraction classification in astronomy. Such methods could, and should be explored in upcoming transient surveys such as the Rubin Observatory (Street et al. 2020) and the La Silla Schmidt Southern Survey (P. E. Nugent et al. 2022, in preparation) among others.

This research used resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operated under Contract No. DE-AC02-05CH11231. V.A.'s work was supported by the Computational Center for Excellence, a Computational HEP program in the Department of Energy Science Office of High Energy Physics (grant #KA2401022).

P.E.N. and R.A.K. acknowledge support from the DOE under grant DE-AC02-05CH11231, Analytical Modeling for Extreme-Scale Computing Environments.

ORCID iDs

Venkitesh Ayyar,  <https://orcid.org/0000-0001-9081-1840>

Robert Knop, Jr.  <https://orcid.org/0000-0002-3803-1641>

Autumn Awbrey,  <https://orcid.org/0000-0001-5287-3004>

Alexis Andersen,  <https://orcid.org/0000-0001-8983-4893>

Peter Nugent,  <https://orcid.org/0000-0002-3389-0586>

References

- Acero-Cuellar, T., Bianco, F., Dobler, G., Sako, M., & Qu, H. 2022, arXiv:2203.07390
- Ayyar, V., Bhimji, W., Gerhardt, L., Robertson, S., & Ronaghi, Z. 2020, in 24th Int. Conf. Computing in High Energy and Nuclear Physics (Les Ulis: EDP Sciences), 06003
- Bailey, S., Aragon, C., Romano, R., et al. 2007, *ApJ*, **665**, 1246
- Bellm, E. C., Kulkarni, S. R., Graham, M. J., et al. 2018, *PASP*, **131**, 018002
- Bhimji, W., Farrell, S. A., Kurth, T., et al. 2018, *J. Phys. Conf. Ser.*, **1085**, 042034
- Bloom, J., Richards, J. W., Nugent, P. E., et al. 2012, *PASP*, **124**, 1175
- Cabrera-Vives, G., Reyes, I., Förster, F., Estévez, P. A., & Maureira, J.-C. 2016, in 2016 Int. Joint Conf. on Neural Networks (IJCNN) (Piscataway, NJ: IEEE), 251
- Cabrera-Vives, G., Reyes, I., Förster, F., Estévez, P. A., & Maureira, J.-C. 2017, *ApJ*, **836**, 97
- Chollet, F. 2015, Keras, <https://keras.io>
- Ciregan, D., Meier, U., & Schmidhuber, J. 2012, in 2012 IEEE Conf. on Computer Vision and Pattern Recognition (Piscataway, NJ: IEEE), 3642
- Duev, D. A., Mahabal, A., Masci, F. J., et al. 2019, *MNRAS*, **489**, 3582
- Flaugher, B. 2005, *Int. J. Modern Physics A*, **20**, 3121
- Gieseke, F., Bloemen, S., van den Bogaard, C., et al. 2017, *MNRAS*, **472**, 3101
- Goldstein, D. A., D'Andrea, C. B., Fischer, J. A., et al. 2015, *AJ*, **150**, 82
- Gómez, C., Neira, M., Hernández Hoyos, M., Arbeláez, P., & Forero-Romero, J. E. 2020, *MNRAS*, **499**, 3130
- Harris, C. R., Millman, K. J., van der Walt, S. J., et al. 2020, *Natur*, **585**, 357
- He, K., Zhang, X., Ren, S., & Sun, J. 2015, CoRR, arXiv:1512.03385
- Jayaraman, R., Fausnaugh, M., & Ricker, G. 2021, *BAAS*, **53**, 1
- Killestein, T. L., Lyman, J., Steeghs, D., et al. 2021, *MNRAS*, **503**, 4838
- Kluyver, T., Ragan-Kelley, B., Pérez, F., et al. 2016, in Positioning and Power in Academic Publishing: Players, Agents and Agendas, ed. F. Loizides & B. Schmidt (Amsterdam: IOS Press), 90
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. 2012, in Advances in Neural Information Processing Systems, ed. F. Pereira, Vol. 25 (Red Hook, NY: Curran Associates, Inc.)
- Lecun, Y., Bottou, L., Bengio, Y., & Haffner, P. 1998, *Proc. IEEE*, **86**, 2278
- Mahabal, A., Rebbapragada, U., Walters, R., et al. 2019, *PASP*, **131**, 038002
- Ronneberger, O., Fischer, P., & Brox, T. 2015, CoRR, arXiv:1505.04597
- Street, R. A., Bianco, F. B., Bonito, R., et al. 2020, *RNAAS*, **4**, 41
- Szegedy, C., Liu, W., Jia, Y., et al. 2015, in 2015 IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (Piscataway, NJ: IEEE), 1
- Turpin, D., Ganet, M., Antier, S., et al. 2020, *MNRAS*, **497**, 2641
- Van Rossum, G. 2020, The Python Library Reference, release 3.8.2 (Python Software Foundation)
- Wright, D. E., Smartt, S. J., Smith, K. W., et al. 2015, *MNRAS*, **449**, 451