

A transient search using combined human and machine classifications

Darryl E. Wright,¹★ Chris J. Lintott,¹ Stephen J. Smartt,² Ken W. Smith,²
Lucy Fortson,^{3,4} Laura Trouille,^{5,6} Campbell R. Allen,¹ Melanie Beck,³
Mark C. Bouslog,⁶ Amy Boyer,⁶ K. C. Chambers,⁷ Heather Flewelling,⁷
Will Granger,⁶ Eugene A. Magnier,⁷ Adam McMaster,¹ Grant R. M. Miller,¹
James E. O’Donnell,¹ Brooke Simmons,⁸† Helen Spiers,¹ John L. Tonry,⁷
Marten Veldthuis,¹ Richard J. Wainscoat,⁷ Chris Waters,⁷ Mark Willman,⁷
Zach Wolfenbarger⁶ and Dave R. Young²

¹*Department of Physics, University of Oxford, Denys Wilkinson Building, Keble Road, Oxford OX1 3RH, UK*

²*Astrophysics Research Centre, School of Mathematics and Physics, Queen’s University Belfast, Belfast BT7 1NN, UK*

³*Minnesota Institute for Astrophysics, University of Minnesota, Minneapolis, MN 55454, USA*

⁴*School of Physics and Astronomy, University of Minnesota, Minneapolis, MN 55455, USA*

⁵*Center for Interdisciplinary Exploration and Research in Astrophysics (CIERA) and Department of Physics and Astronomy, Northwestern University, Evanston, IL 60208, USA*

⁶*Citizen Science Department, The Adler Planetarium, Chicago, IL 60605, USA*

⁷*Institute for Astronomy, University of Hawaii, 2680 Woodlawn Drive, Honolulu, HI 96822, USA*

⁸*Center for Astrophysics and Space Sciences, University of California, San Diego, La Jolla, CA 92093, USA*

Accepted 2017 July 17. Received 2017 July 3; in original form 2016 December 16

ABSTRACT

Large modern surveys require efficient review of data in order to find transient sources such as supernovae, and to distinguish such sources from artefacts and noise. Much effort has been put into the development of automatic algorithms, but surveys still rely on human review of targets. This paper presents an integrated system for the identification of supernovae in data from Pan-STARRS1, combining classifications from volunteers participating in a citizen science project with those from a convolutional neural network. The unique aspect of this work is the deployment, in combination, of both human and machine classifications for near real-time discovery in an astronomical project. We show that the combination of the two methods outperforms either one used individually. This result has important implications for the future development of transient searches, especially in the era of Large Synoptic Survey Telescope and other large-throughput surveys.

Key words: methods: data analysis – methods: statistical – techniques: image processing – surveys – supernovae: general.

1 INTRODUCTION

The detection and identification of transient sources has long been an important part of astronomical observation. New surveys such as LSST (Large Synoptic Survey Telescope; Ivezić et al. 2008) will increase the number of transient candidates detected by many orders of magnitude, leading to renewed attention being paid to the methods used by transient searches. To extract the most scientific value from surveys, we want to follow the entire evolution of transients from the time of outburst to the point at which they fade below the detection limit. This requires a rapid processing of data to

enable a fast decision on whether or not to expend valuable follow-up resources for each potential candidate extracted by a transient survey’s image processing pipeline. The first problem is deciding if a source, flagged by the pipeline, is a detection with real astrophysical significance or an artefact of the detector or image processing. We want to promote the former for a decision on whether to follow up and to reject the latter without further consideration. In preparation for LSST and to deal with the data volumes of present surveys, much effort has been invested in developing systems that automatically reject false positives with supervised learning. Using large volumes of past observations that have been identified as real or ‘bogus’, the aim is to train a machine to make predictions about future observations (Romano, Aragon & Ding 2006; Bailey et al. 2007; Donalek et al. 2008; Bloom et al. 2012; Brink et al. 2013; du Buisson et al. 2015; Goldstein et al. 2015; Masci et al. 2017).

* E-mail: darryl@zooniverse.org

† Einstein Fellow.

Providing these data typically requires manual inspection of individual detections by human experts to mitigate label contamination that would confound the learning algorithm. This quickly becomes unwieldy given that performance of a machine learning solution has been shown to depend on the quantity of labelled training data (Banko & Brill 2001).

The requirement for sufficiently large and representative training sets is prohibitive for the largest surveys or for small research teams, and is particularly problematic for rarer classes of transients. An alternative to relying on expert labelled training sets for machine learning is therefore to expand the population providing labels. For example, Melchior et al. (2016) describe a crowd sourcing platform for vetting image quality from the Dark Energy Survey (DES; Dark Energy Survey Collaboration 2005) relying on the consensus of ~ 100 volunteers from a team of professional astronomers. For surveys of the future, a few hundred volunteers will not be enough. Instead, we must cast our net wider increasing the population of the crowd beyond those directly involved in the survey. The obvious path is to engage citizen scientists. *Galaxy Zoo Supernovae* (Smith et al. 2011) and *Snapshot Supernova* (Campbell et al. 2015) are two projects to have taken this approach for transient surveys, using data from the Palomar Transient Factory (Law et al. 2009; Rau et al. 2009) and SkyMapper (Keller et al. 2007), respectively. Both projects were facilitated through the *Zooniverse* citizen science platform (see description in Marshall, Lintott & Fletcher 2015), and asked volunteers to assess the target, reference and difference images for each detection and answer a series of questions that led to a classification of real or bogus.

For classification tasks, humans and machines have complementary strengths. Human classifiers are good at rapidly making abstract judgements about data, allowing them to see only a small number of examples before making decisions about novel images. Machines can consume large quantities of data and make more systematic judgements based on complex relationships between the features provided. In the *Space Warps* citizen science project (Marshall et al. 2016), all example gravitational lenses provided to volunteers appeared blue, yet despite this, volunteers were able to identify a gravitationally lensed hyperluminous infrared radio galaxy that appeared red (Geach et al. 2015). In contrast, a machine would need to be provided examples of these in the training data. On the other hand, computer vision techniques allow images to be examined systematically, with relationships between different features used for classification. Combining machine classifications with those of experts will trivially be expected to improve performance, but the situation in which classifications from volunteer citizen scientists are used is less clear. If the noisier data sets provided by citizen science are combined with machines, is performance of the system improved? If machines are included in classification, does it relieve some of the burden on citizen scientists? Answering these questions is critically important for surveys where even a substantial number of experts will not be able to review all the data promoted by a machine classifier.

In this paper, we report some initial findings from the *Supernova Hunters* project,¹ a new citizen science project similar in spirit to those mentioned above but applied to the Pan-STARRS Survey for Transients (PSST). In Section 2, we describe the Pan-STARRS1 telescope, PSST survey and the Supernova Hunters project and citizen science platform. Section 3 shows the relative performance

of humans and machines on data uploaded to Supernova Hunters during the first two months of the project. We also describe and measure the performance of a simple method for combining the classifications of citizen scientists and the current PSST machine classifier. We further discuss a mechanism to take advantage of metadata associated with each detection to boost classification performance. In Section 4, we conclude and discuss potential avenues for future improvements. This paper therefore represents the first study of combined citizen science and machine classifications within a live astronomical survey.

2 METHOD

2.1 Pan-STARRS1

Pan-STARRS1 comprises a 1.8 m primary mirror (Kaiser et al. 2010) and 60 detectors with 4800 pixels, constructed from $10\ \mu\text{m}$ pixels subtending 0.258 arcsec (Magnier et al. 2013) and a field of view of 3.3 deg. The filter set consists of g_{P1} , r_{P1} , i_{P1} , z_{P1} (similar to SDSS *griz*; York et al. 2000), y_{P1} extending redwards of z_{P1} and the ‘wide’ w_{P1} -band filter extending over g_{P1} to i_{P1} (Tonry et al. 2012). Between 2010 and 2014, Pan-STARRS1 was operated by the PS1 Science Consortium performing two major surveys. The Medium Deep Survey (MDS) was allocated 25 per cent of observing time for high-cadence observations of the 10 Medium Deep fields and the 3π survey allocated 56 per cent observing time to observe the entire sky north of -30° declination with four exposures per year in each of g_{P1} , r_{P1} , i_{P1} , z_{P1} and y_{P1} for each pointing.

The 3π survey was completed in mid-2014 and since then the telescope has been carrying out a NASA-funded wide-field survey for near-Earth objects through the NEO Observation Program operated by the Pan-STARRS Near Earth Object Science Consortium (PSNSC). The NASA PSNSC survey is similar to the 3π survey but optimized for NEO discoveries. Observations are in w_{P1} in dark time and combinations of i_{P1} , z_{P1} and y_{P1} during bright time. The PSST survey (Inserra et al. 2013; Huber et al. 2015) searches the data for static transients, releasing these publicly within 12–24 h.

Typically, a single field is imaged four times in a night with exposures separated by 10–20 min called transient time interval (TTI) exposures to allow for the discovery of moving objects. The quads of exposures are not dithered or stacked, meaning that cross-talk ghosts, readout artefacts and problems of fill factor are inherent in the data (see Denneau et al. 2013 for some examples). Individual exposures are differenced (Alard & Lupton 1998; Bramich 2008) with the 3π all-sky reference stack, and sources in the resulting difference images are catalogued.

A series of pre-ingest cuts are performed before the catalogues are ingested into a MySQL data base at Queen’s University Belfast (QUB). These cuts are based on the detection of saturated, masked or suspected defective pixels within the point spread function area in addition to flag checks for ghost detections and rejecting detections within $\pm 5^\circ$ galactic latitude. Detections passing these cuts are grouped into transient candidates if they are spatially coincident within 0.5 arcsec and the rms scatter is < 0.25 arcsec. Post-ingest cuts are applied on detection quality, convolution checks and a check for proximity to bright objects. Additional cross-talk rules have been identified and implemented at QUB to reject ghosts not flagged at the pre-ingest stage. Remaining detections are cross-matched with the Minor Planet Center ephemeris data base to identify any asteroids not removed by the rms cut. Remaining transient candidates are passed to our machine classifier described in the next section.

¹ <https://www.zooniverse.org/projects/dwright04/supernova-hunters>

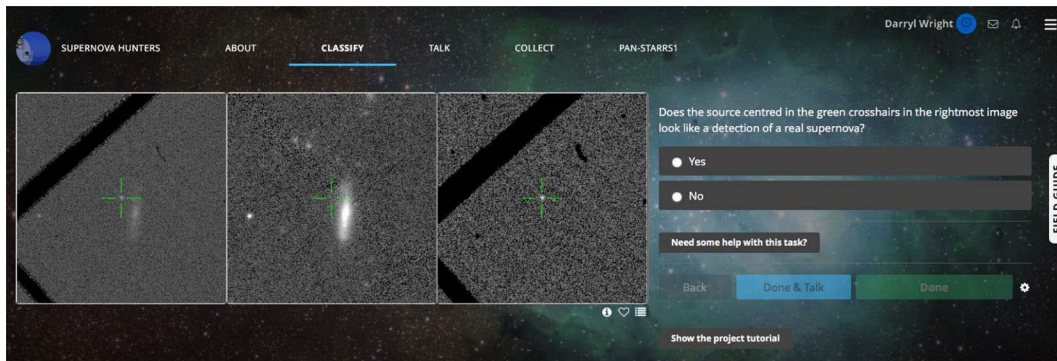


Figure 1. The classification interface presented to citizen scientists. The leftmost image is the *target* image taken during the previous week. In the centre is the equivalent 3π *reference* image and on the right is the *difference* image. Volunteers are asked to decide whether or not they think the detection in the green cross hairs in the difference image is a detection of a real transient.

2.2 Convolutional neural network

In Wright et al. (2015), we developed a machine classifier for real–bogus classification in the PS1 MDS. However, we found that this approach performs poorly for PSST because of the greater variety of artefacts in PSST data (a consequence of differencing individual exposures) and the difficulty obtaining a representative labelled training set at the beginning of a new survey. Instead, we turned to convolutional neural networks (CNNs) that maintain the advantages of operating solely on the pixel data but at a higher computational cost in deployment.

The training set for this classifier was drawn from 3π survey data between 2013 June 1 and 2014 June 20. The sample of real detections are taken from spectroscopically confirmed transients or detections of objects that have been labelled by experts as high-probability real transients. Bogus detections are taken from a random subsample of detections discarded by post-ingest cuts or human inspection. The training set consists of 6916 examples with an additional 2303 detections held out for testing with both data sets containing twice as many bogus detections to real. Each example was manually inspected in order to limit label contamination; not all detections associated with a spectroscopically confirmed transient are necessarily real for example.

Given the small data set, to avoid overfitting we limit the CNN to a single convolution layer with 400 kernels and a pooling layer followed by a binary softmax classifier. We also perform unsupervised pre-training with sparse filtering (Ngiam et al. 2011) using unlabelled images from the STL-10 (Coates, Lee & Ng 2011) data set. The classifier is applied to nightly PSST data producing a score for each TTI exposure for every candidate passing the cuts in the previous section. The score, or *hypothesis*, is a function $h(x)$ of the input feature representation, x (the output of the convolution and pooling layers). For each candidate, we simply combine the TTI exposure hypotheses by taking the median, resulting in a single ‘real–bogus factor’ for each transient candidate that we take as the machine equivalent of $P(\text{real})$ below. To automatically reject candidates, we must choose a decision boundary on $h(x)$ such that any candidate with a hypothesis lying below the decision boundary is considered a bogus detection and discarded. This inevitably leads to a trade-off between false positives and false negatives (or missed detections). If the decision boundary is set too high, we will discard many real detections of supernovae; if it is too low and we will be inundated with artefacts (see for example Fig. 2). We chose the decision boundary based on the expected performance measured on the test set. For example, using our CNN to generate hypotheses for

each detection in the test set, we can choose the decision boundary that corresponds to a false positive rate (FPR) of 1 per cent at $h(x) = 0.842$, where the FPR is the number of false positives divided by the total number of bogus candidates in the sample. However, although the number of artefacts promoted would be low, based on the test set this decision boundary would be expected to result in a missed detection rate (MDR) of ~ 21 per cent for future data and is therefore not a sensible choice. We instead opted for a decision boundary at $h(x) = 0.436$ with expected FPR and MDR of 5 per cent and ~ 5.2 per cent, respectively. As detailed in Section 3.1, clearly the decision boundary can be scaled to take advantage of available human effort; lowering the decision boundary beyond 0.436 would increase the FPR requiring more human screening but at the same time reduce the MDR such that humans could recover real supernova detections with low $h(x)$ that would otherwise be automatically rejected.

2.3 Citizen science platform

Supernova Hunters was launched on 2016 July 12 (MJD 57581). As of 6th December 2016, the project has accumulated 1082 170 classifications from 5845 citizen scientists with a few tens of volunteers submitting thousands of classifications. Citizen scientists are presented with the interface shown in Fig. 1 and asked to classify individual TTI observations (see Section 2.1). So far volunteers have classified 117 693 individual images of 46 277 individual PS1 objects. As guidance, we provide a ‘Field Guide’ that provides a description and examples of the different artefact types we expect. Once a week ~ 5800 new subjects are uploaded to the project consisting of the previous week’s detections that pass our machine cuts. The arrival of the data is announced to existing volunteers via email.² We require at least seven citizen scientist classifications before a subject is considered classified and subsequently retired from the project. The choice of seven classifications is simply motivated by a trade-off between speed of data processing and accuracy, and we did not see significant gains by requiring 10 classifications during the beta test. Since launch the project averages ~ 21 000 classifications in the first 24 h after the data are released and ~ 8200 classifications in the following 24 h by which time all subjects are normally retired.

² Though enthusiasm for the project is such that traffic to the site now increases significantly *ahead* of this email alert!

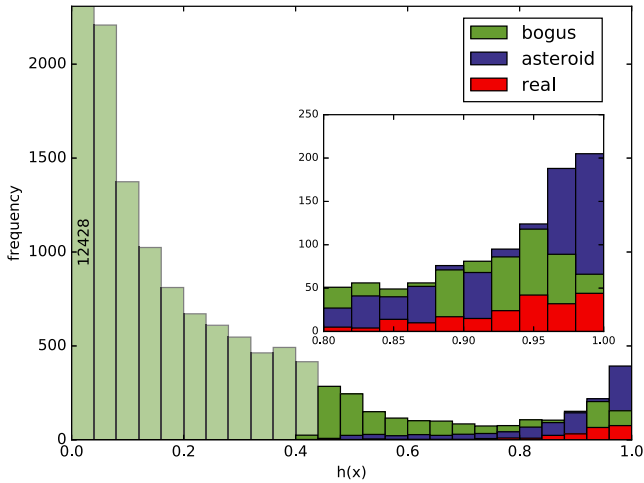


Figure 2. The distribution of hypotheses, $h(x)$, from the current 3π machine classifier for detected objects between MJD 57570 and MJD 57586. The light green shows the distribution of objects with $h(x) \leq 0.436$ that are automatically rejected. The remaining objects promoted for human screening even at high values of $h(x)$ contain many false positives. The first interval has a frequency of 12 428, but the plot is truncated for clarity. Inset: zoom-in of the region with $h(x) > 0.8$.

We calculate a ‘human score’ simply by taking the fraction of all the volunteers who saw a detection and classified it as real. High-confidence (typically $P(\text{real}) > 0.8$, see Section 3) supernova candidates are screened by experts to remove a small number of false positives (~ 10 per cent) before the targets are submitted to the Transient Name Server (TNS). To date, citizen scientists have discovered over 450 supernova candidates that have been submitted to the TNS and two confirmed supernovae including SN 2016els (Mattila et al. 2016) and a superluminous supernova Type I. The classification spectra were obtained by PESSTO (Smartt et al. 2015).

3 PERFORMANCE

We present the results of our machine classifier (Section 2.2) in Fig. 2 on PS1 data uploaded to Supernova Hunters between MJD 57570 and MJD 57586. This data set includes classifications from an initial beta test of the project prior to launch on MJD 57581. A major contaminant is the presence of asteroids. These appear in the difference image as identical to supernovae, and are in that sense ‘correctly classified’ but are identified here via cross-matching with the Minor Planet Center.

The results of the machine learning were additionally reviewed by at least one expert member of the team (normally DEW or KWS) to identify genuine supernovae. Candidates were divided into ‘real’ and ‘bogus’ categories based on these expert classifications. We note that a future improvement to the project would be to inject fake real and bogus detections into the data as a means to track performance. The Andromeda Project (Johnson et al. 2015) and Space Warps are examples of two citizen science projects to have taken this approach, while Goldstein et al. (2015) used faked detections to augment DES training data for their real–bogus classifier. With this approach we would no longer be reliant on the assumption that every expert label is correct, but we must be careful to ensure that injected fake sources are truly representative of our observations.

Candidates with high scores as assigned by the machine are more likely to be real. However, although the machine successfully rejects the majority of bogus candidates, the sample produced by the

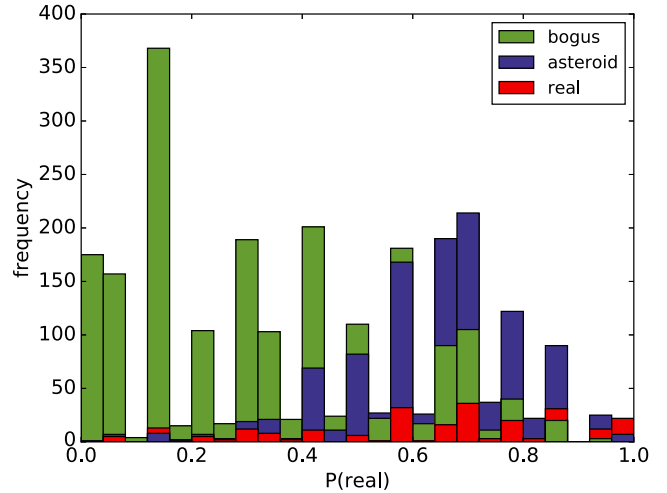


Figure 3. The distribution of $P(\text{real})$ from Supernova Hunters for objects detected between MJD 57570 and MJD 57586. Compared with the machine $h(x)$ in Fig. 2, the objects at the extremes are pure. There are no real detections with $P(\text{real}) < 0.04$ and few bogus detections above 0.92.

simple cut on hypothesis is far from pure: 1403 real candidates from 3384 in the sample. Higher cutoffs run the risk of rejecting an increasing number of real candidates; requiring a 1 per cent FPR will result in an MDR of 60.3 per cent. In order to improve this performance, candidates that exceed the $h(x) = 0.436$ threshold (see Section 2.2) were also classified by volunteers via the *Supernova Hunters* project. The results of this analysis are shown in Fig. 3.

Volunteer classifications were combined using the simplest possible metric; the fraction of volunteers who identified a detection as real is assumed to be an estimate of the probability of that candidate being real, denoted by $P(\text{real})$. Despite this simple procedure, the results show that volunteers could effectively distinguish between real and bogus classifications. However, the structure of the resulting distribution is strikingly different from that of the machine classifier. Whereas for machine classification, a threshold could be chosen to give a complete but not pure sample, with volunteer classifications it is easier to construct a pure sample of candidates that are highly likely to be supernovae, but this sample is far from complete. There are candidates judged ‘real’ by experts even at low probabilities, although there were no real candidates assigned $P(\text{real}) < 0.04$.

There are two routes that might be expected to improve this performance. First, we could improve on the naive combination of volunteer votes described above. To this end, citizen science projects typically explore methods to weight volunteer contributions [see for example, Schwamb et al. (2012), Willett et al. (2013) and Marshall et al. (2016)]. Secondly, given that we have a human and machine score for every detection, we could seek a combination of the two in the hope of benefiting from the different capabilities of both.

3.1 Combining human and machines

In a companion paper analysing performance of Galaxy Zoo, Beck et al. (2017) used a machine classifier running in parallel to a simulation of the Galaxy Zoo 2 project (Willett et al. 2013). They showed that the addition of such a classifier, which retired subjects classified above a certain level of confidence at the end of each day, retraining each time, can greatly accelerate the speed of classification in a data set. In their work, images are retired by either machine or human, whereas we set out in this section to use a combination. While they

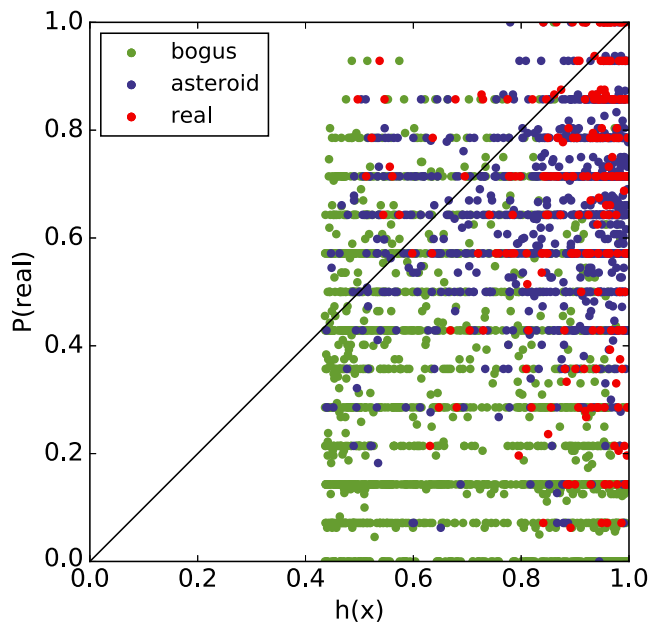


Figure 4. The $P(\text{real})$ from Supernova Hunters against the machine $h(x)$ for 3384 detected objects between MJD 57570 and MJD 57586. $P(\text{real})$ and $h(x)$ are combined by projecting the data on to the solid black line in the Euclidean sense. A Spearman rank correlation test shows the correlation between $P(\text{real})$ and $h(x)$ to be 0.237.

retrain their machine with volunteer input as it accumulates, we use a static training set derived from expert classifications.

Fig. 4 shows the combination of human and machine classifications. It is immediately apparent from the figure that no single threshold on either machine or human classification can outperform the combination of the two. This is an important result; it is the first time that the benefits of combining classification from both machines and volunteers have been clearly demonstrated using data from a live astronomical survey.

How should the two independent classifications be combined? We simply apply a decision boundary of the form $\tau = (x + y)/2$ on the 2D surface, where $0 \leq \tau \leq 1$. For a constant value of τ , a candidate is classified as bogus if $[h(x) + P(\text{real})]/2 \leq \tau$ and classified real otherwise. This is equivalent to projecting the data on to $P(\text{real}) = h(x)$, producing a new scalar score for each detection and classifying candidates as bogus if the combined score is less than or equal to τ .

As an independent test, we apply the same method to data between MJD 57587 and MJD 57627 in Fig. 5. For Figs 4 and 5, a Spearman rank correlation test gives 0.237 and 0.122, respectively, showing $P(\text{real})$ and $h(x)$ to be only weakly correlated. Between MJD 57609 and MJD 57615, we relaxed our cut on $h(x)$ from 0.436 to 0.3 uploading any objects passing this cut to Supernova Hunters. This allowed us to explore the subjects wrongly classified by the machine and resulted in the recovery of SN 2016fev, a Type Ia supernova that would have been automatically rejected with $h(x) = 0.39$, but which received a $P(\text{real})$ of 1.0 from Supernova Hunters. The performance of the combination method on this data set is shown in the receiver operator characteristic (ROC) curve and purity–completeness (precision–recall) curves plotted in Fig. 6.

For any choice of FPR, the combination of classifications produced a lower MDR. Equally, for any required purity or completeness, the combination provides a better trade-off.

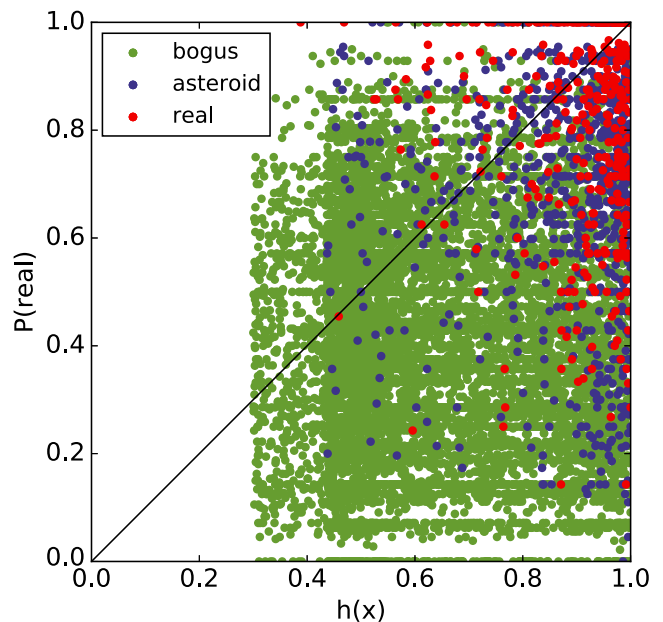


Figure 5. The same as Fig. 4 but on a new sample of 10 908 objects detected between MJD 57587 and MJD 57627. For one week during this period, we relaxed our cut on $h(x)$ to 0.3, which allowed us to recover a supernova with $h(x) = 0.39$, but which achieved a $P(\text{real}) = 1.0$ from Supernova Hunters. In this case, the Spearman rank correlation is found to be 0.122.

We have chosen to implement one of the simplest methods for combining human and machine classifications to demonstrate how they complement one another, but it is easy to think of more complex combination methods. For example, we trained a linear support vector machine (SVM) on the data presented in Fig. 4 and found, unsurprisingly, that the performance measured on the data in Fig. 5 was typically within 1 per cent of the values reported in Tables 1 and 2. Although the gains in this example are negligible, if we wish to incorporate additional information from an ensemble of machine classifiers for example, and the input space becomes higher dimensional, such methods become important. It is unlikely that the combination method presented here will work for higher dimensional data, but we can expect that an SVM may take advantage of the additional information.

3.2 Improving $P(\text{real})$ with SWAP

We also expect to gain from improving $P(\text{real})$. To demonstrate this, we implemented the Space Warps Analysis Pipeline (SWAP; Marshall et al. 2016), an algorithm designed to improve the sample of good gravitational lens candidates from the Space Warps³ citizen science project in images from the Canada–France–Hawaii Telescope Legacy Survey (Gwyn 2012). SWAP does not treat all classifiers equally, but quantifies the value of their contribution in terms of the information gained. SWAP assigns a software agent to individual citizen scientists. Classifications are weighted according to the agent’s estimate of likely performance based on each volunteer’s past performance measured on gold standard data. The agent maintains a confusion matrix that monitors the fraction of gold standard data correctly or incorrectly classified by a volunteer for each class (‘LENS’ or ‘NOT’ in the case of Space Warps). We

³ <https://spacewarps.org>

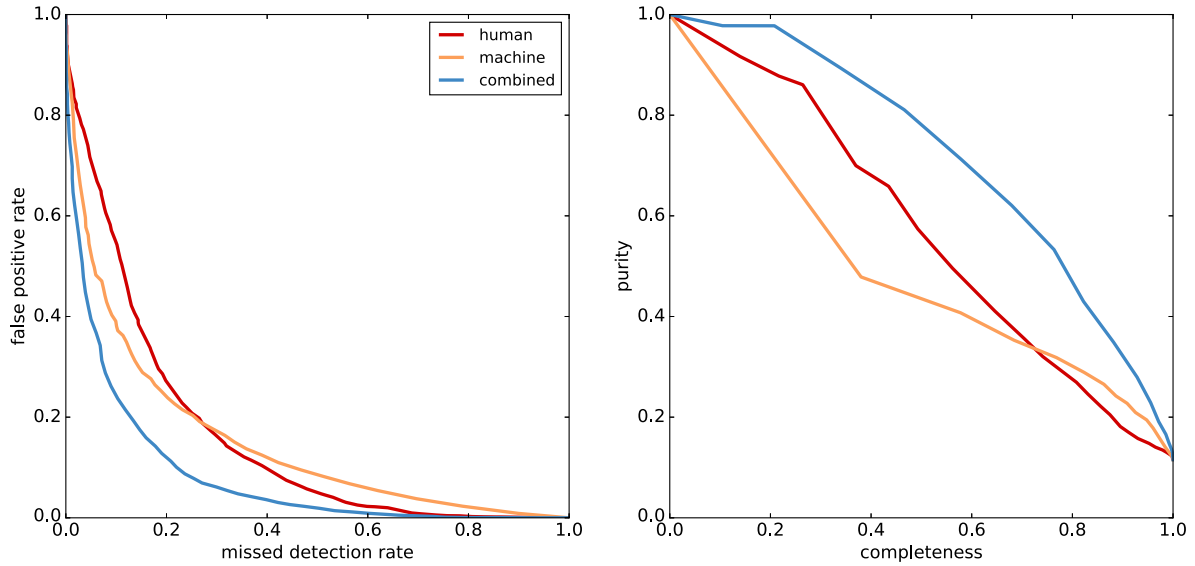


Figure 6. Left: ROC curve showing performance measured on data in Fig. 5 for human (red), machine (yellow) and the combination of human and machine classifications (blue). Right: the equivalent purity–completeness curve. Both plots show that the combination always outperforms humans and the machine individually.

Table 1. MDR recorded for a choice of FPR, based on expert classifications.

FPR	Human	Machine	Combined	SWAP	SWAP (mag.)	SWAP (seeing)	Combined SWAP	Combined SWAP (mag.)	Combined SWAP (seeing)
1 per cent	73.9 per cent	90.1 per cent	58.7 per cent	66.6 per cent	64.5 per cent	64.2 per cent	54.5 per cent	54.1 per cent	53.5 per cent
5 per cent	56.3 per cent	69.7 per cent	35.8 per cent	41.0 per cent	39.7 per cent	40.3 per cent	30.2 per cent	29.7 per cent	28.8 per cent
10 per cent	45.6 per cent	46.7 per cent	23.8 per cent	31.2 per cent	30.0 per cent	30.5 per cent	20.5 per cent	18.9 per cent	20.1 per cent

Table 2. FPR recorded for a selection of MDR, based on expert classifications.

MDR	Human	Machine	Combined	SWAP	SWAP (mag.)	SWAP (seeing)	Combined SWAP	Combined SWAP (mag.)	Combined SWAP (seeing)
1 per cent	92.5 per cent	85.9 per cent	69.3 per cent	100.0 per cent	100.0 per cent	100.0 per cent	66.2 per cent	68.7 per cent	71.2 per cent
5 per cent	75.1 per cent	52.8 per cent	41.8 per cent	100.0 per cent	100.0 per cent	100.0 per cent	29.6 per cent	31.0 per cent	31.7 per cent
10 per cent	53.8 per cent	39.1 per cent	26.5 per cent	49.9 per cent	47.9 per cent	49.1 per cent	20.6 per cent	19.6 per cent	20.2 per cent

use the online implementation of SWAP considering each classification in turn. We make several small adjustments to the original SWAP implementation. First, our gold standard labels are ‘real’ and ‘bogus’ rather than ‘LENS’ or ‘NOT’. Secondly, we set the prior probability that each detection is a real supernova ρ_0 to be 0.01, roughly the expected ratio between real and bogus detections each night determined by expert classifications. Finally, in our analysis, we do not set rejection and detection thresholds, instead we continue to require that at least seven volunteers classify each subject. As in Marshall et al. (2016), we set the initial confusion matrix for the i th volunteer to

$$M^i = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}.$$

This initialization of the confusion matrix corresponds to that of a random classifier, but will be quickly modified as we observe classifications of gold standard data by this volunteer. The results of applying SWAP to the citizen scientist classifications on data between MJD 57587 and MJD 57627 (Fig. 5) are shown as the purple line in Fig. 7. Compared with the human performance in

Fig. 6, SWAP considerably improves $P(\text{real})$. By combining this improved score from human classifications with the machine, using the same method as above, we once again observe an improvement in the measured performance (pink line from Fig. 6). In fact, apart from a small exception at about 97 per cent purity and 25 per cent completeness (likely an artefact of the specific data set), the combination outperforms all previous methods for classifying this data set, demonstrating that an improvement in either $P(\text{real})$ or $h(x)$ can lead to valuable performance gains when combined. SWAP combined with the machine produces a few per cent improvement over all the other methods tested so far for any figure of merit in Tables 1 and 2.

Considering the scenario where we are reliant on citizen scientists to label training data for machine learning algorithms, it seems prudent to consider how we might further improve $P(\text{real})$ in the future. One observation we make is that citizen scientists perform differently on detections with varying degrees of signal-to-noise. In general, citizen scientists are very good at recovering bright supernovae but find it more and more difficult as detections get closer to the detection limit. Given the history of classifications of gold

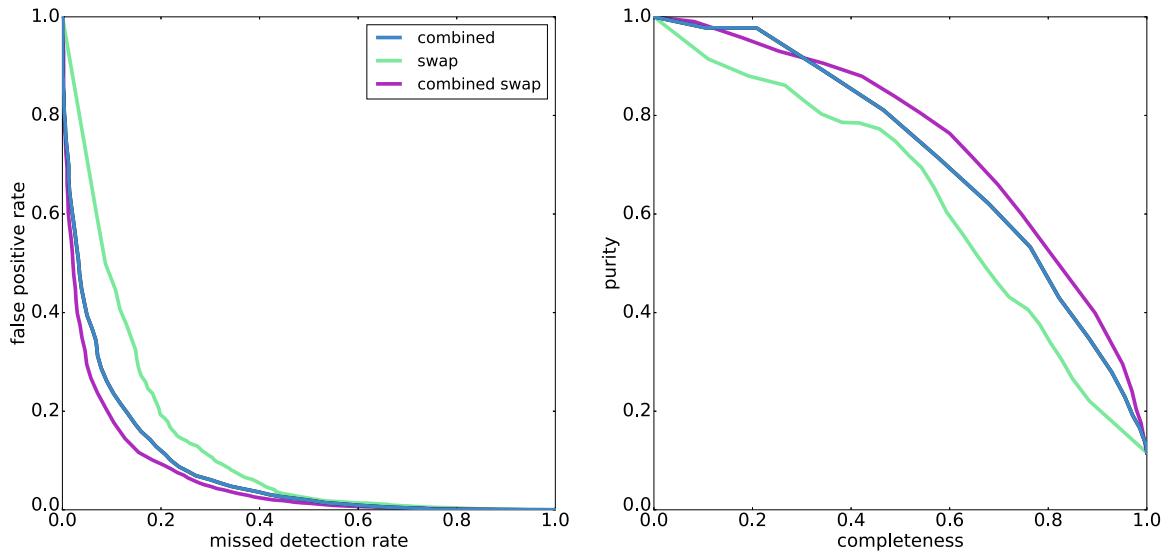


Figure 7. Left: ROC curve for different approaches to classification. Right: the corresponding purity–completeness curves. The blue lines are the same as in Fig. 6 showing the simple combination of the machine with the basic human score described in Section 3. The green line shows the results of implementing SWAP, which improves on the basic human score (red line in Fig. 6). Combining this improved score with the machine produces the pink line.

standard data a volunteer has submitted, we can calculate the probability the volunteer will classify a detection as real given that the true label is real, $P(\text{'real'}|\text{real})$, by simply taking the fraction of the gold standard detections labelled ‘real’ the volunteer correctly classified. We can similarly calculate $P(\text{'bogus'}|\text{bogus})$. These are the probabilities tracked by the SWAP agent confusion matrix and that we have plotted on the unit plane for the 3158 volunteers participating in the project at this time for different magnitude bins in Figs 8 and 9. A volunteer who correctly classifies every gold standard example they have seen is a perfect classifier and will lie at (1,1) in the ‘astute’ quadrant of the plot. The size of the point corresponds to the quantity of gold standard data classified by that volunteer in a given magnitude bin. Most volunteers have submitted only a few classifications and cannot be seen on the plot. Clearly, volunteers are more ‘astute’ at classifying brighter sources, tending to lie closer to the top right of the plot. This information is lost in the basic SWAP implementation. A volunteer’s classification could now be weighted according to the magnitude bin that a given detection falls in. This results in a simple modification of the SWAP calculation of $P(\text{real})$. In this case, we are taking advantage of additional ‘metadata’ that is available for each detection and assuming that it has some effect on classifiers’ behaviour. We also expect that humans and machines are good at classifying different types of images. In our case, the CNN is good at classifying detections around 20th magnitude, but often misclassifies detections brighter than 17th magnitude. This is because the training set is dominated by fainter detections. The machine has not been able to learn a function that accurately maps between the pixels and the classification for bright examples because it has not ‘seen’ many during training. In Tables 1 and 2, we include the results of running SWAP on the magnitude splits and the combination with the machine scores. We find that, for the most part, this gives a ~ 1 per cent improvement in performance over not using metadata splits depending on the figure of merit. There are other metadata parameters we could explore, both properties of the survey (seeing) and properties of the detection (proximity to a galaxy). Tables 1 and 2 also show the effect of using seeing for the metadata splits, and we observe similar performance to magnitude. Although the gains are small, this may

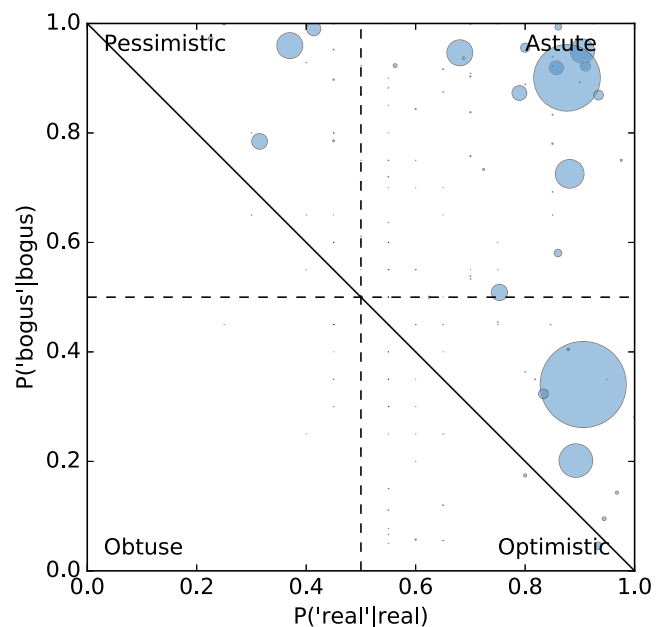


Figure 8. Confusion matrix elements for volunteers classifying detections from Fig. 5 with apparent magnitude $13 \leq m < 18$. Each point represents an individual volunteer, where the size of the point corresponds to the quantity of gold standard data in this magnitude bin classified by the volunteer. Larger points therefore correspond to more experienced classifiers. Many volunteers tend to lie high in the top right of the ‘astute’ region; a perfect classifier would lie at (1,1). This shows that citizen scientists can accurately classify high signal-to-noise real and bogus detections.

indicate that a more thorough analysis of metadata is worth pursuing. For instance, we could more carefully define where the splits on metadata should be made and consider the effect of combining multiple metadata parameters. In future, it would be interesting to train an ensemble of machines with each only consuming retired data from a single metadata split. When adequately trained, the resulting machines would specialize in specific regions of parameter

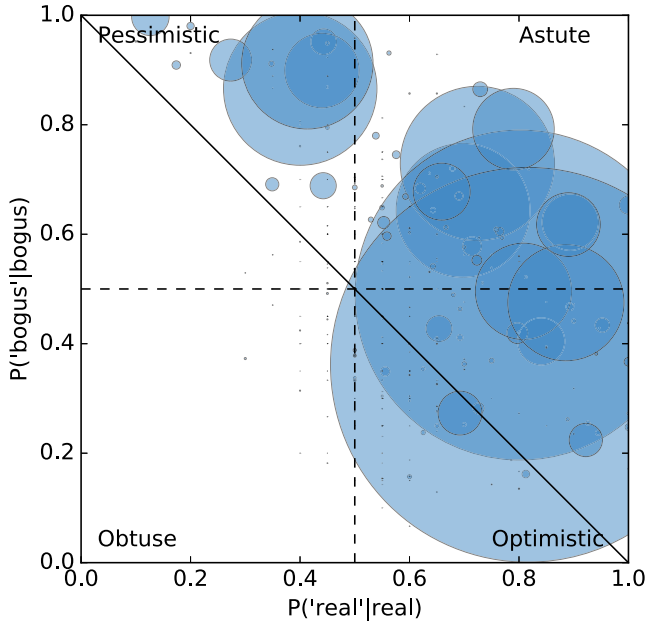


Figure 9. Similar to Fig. 8 but for detections with apparent magnitude $19 \leq m < 20$. In this case, volunteers are further from the top right and tend towards ‘pessimistic’ or ‘optimistic’ classifiers.

space and could eventually be folded back into the project. We could also incorporate these machines into the SWAP analysis with their classifications weighted just as for volunteers.

4 CONCLUSIONS

In this paper, we introduced a new citizen science project, Supernova Hunters, built entirely with the off-the-shelf Zooniverse Project Builder⁴ requiring no custom features or additional development. The project aims to classify detections of potential supernova candidates from the PSST survey either as real transients or bogus detections of image processing or instrumentation artefacts. To be uploaded to the Supernova Hunters project, a detection must first pass a series of cuts based on catalogue information and secondly be promoted by our machine learning algorithm, a CNN. With this approach, we expect that only about 5 per cent of the false positives passing the catalogue cuts make it into the project, greatly reducing the number of objects we ask volunteers to screen. Citizen scientists excel at mining these data for a very pure sample of high- $P(\text{real})$ supernova candidates, typically those of higher signal-to-noise and offset from a galaxy. But compared with expert labels, many less obvious candidates are missed. Rather than immediately consider methods of weighting classifications of individual citizen scientists, we instead applied a simple combination of the scores provided by humans and machines. We showed that the new combined score achieved better performance than either individually for any choice of purity or completeness.

We expect that there are many ways to improve. As an example, we generalized the SWAP to Supernova Hunters and applied it to the volunteer classifications and again combined the resulting scores with those from the machine. We found that this resulted in a few per cent performance gain compared to the next best method. To further improve $P(\text{real})$, we could invest more effort into educating citizen

scientists to identify more subtle artefact indicators. For example, the Gravity Spy⁵ project has implemented a training regimen where volunteers are provided with feedback on classifications of gold standard data and can progress to more advanced levels performing more complex tasks (Zevin et al. 2017). This could help address the many bogus detections with $P(\text{real}) > 0.5$, but relatively few real detections below 0.5 in Fig. 3. Another improvement could be with our machine classifier, which was trained at the beginning of the PSST survey and the algorithm was specifically chosen to learn from the limited amount of training data available. Given the large volume of data accumulated since, we could train more sophisticated algorithms that can learn more complex relationships between the features, though extracting robust labels for these additional data is a challenge that still needs to be addressed.

This effort offers hope for dealing with the large data volumes from all-sky surveys such as LSST, ZTF (Bellm 2014), ATLAS (Tonry et al. 2016) and Pan-STARRS2. We used machines to reject the vast majority of false positives and then combined the machine hypotheses with classifications from a few thousand citizen scientists for the remaining candidates. With Supernova Hunters we have not actively sought additional citizen scientists, beyond the $\sim 30\,000$ volunteers on the Zooniverse beta testing e-mail list, who were asked to review the project before launch. New volunteers must ‘discover’ the project on the Zooniverse projects page to participate. Given that we could actively seek the participation of $\sim 10^6$ registered Zooniverse volunteers and assuming that 10 per cent chose to participate, with the current classification rate (21 000 in the first 24 h each week from ~ 6000 volunteers), we could achieve $\sim 350\,000$ classifications per night. This provides 0.35 classifications for the $\sim 10^6$ transient alerts expected from LSST at the beginning of the survey (Ridgway et al. 2014). If the FPR is an order of magnitude more than the transient alert rate [perhaps overly pessimistic given the expected $\sim 500\text{--}2200$ false positives per field per visit (Becker et al. 2013) with a 5σ detection threshold] and assuming that we can discard 90 per cent of those with machine learning, we can expect to achieve 0.175 citizen science classifications per promoted detection. Assuming that we will require roughly 10 classifications per detection before considering it classified, we are roughly two orders of magnitude short. Making up this deficit may be achievable with continued improvements to difference imaging (Zackay, Ofek & Gal-Yam 2016), automated real–bogus classification, encouraging greater participation from a growing community of citizen scientists and more efficient use of their classifications.

ACKNOWLEDGEMENTS

The Pan-STARRS1 Surveys have been made possible through contributions of the Institute for Astronomy, the University of Hawaii, the Pan-STARRS Project Office, the Max Planck Society and its participating institutes, the Max Planck Institute for Astronomy, Heidelberg and the Max Planck Institute for Extraterrestrial Physics, Garching, the Johns Hopkins University, Durham University, the University of Edinburgh, Queen’s University Belfast, the Harvard-Smithsonian Center for Astrophysics, the Las Cumbres Observatory Global Telescope Network Incorporated, the National Central University of Taiwan, the Space Telescope Science Institute, the National Aeronautics and Space Administration under Grant No. NNX08AR22G issued through the Planetary Science Division of

⁴ <https://www.zooniverse.org/lab>

⁵ <https://www.zooniverse.org/projects/zooniverse/gravity-spy>

the NASA Science Mission Directorate, the National Science Foundation under Grant No. AST-1238877, the University of Maryland, and Eotvos Lorand University (ELTE) and the Los Alamos National Laboratory. The development of the Zooniverse platform was supported by a Global Impact Award from Google. We also acknowledge support from STFC under grant ST/N003179/1. LF and MB were partially supported by the National Science Foundation under Grant No. AST-1413610. DW and LF were partially supported by the National Science Foundation under Grant No. IIS-1619177. Support for this work was provided by the National Aeronautics and Space Administration through Einstein Postdoctoral Fellowship Award Number PF5-160143 issued by the Chandra X-ray Observatory Center, which is operated by the Smithsonian Astrophysical Observatory for and on behalf of the National Aeronautics Space Administration under contract NAS8-03060. We also wish to acknowledge the dedicated effort of our citizen scientists who have made this work possible.

REFERENCES

- Alard C., Lupton R. H., 1998, *ApJ*, 503, 325
- Bailey S., Aragon C., Romano R., Thomas R. C., Weaver B. A., Wong D., 2007, *ApJ*, 665, 1246
- Banko M., Brill E., 2001, in Bonnie L. W., ed., Proc. 39th Annual Meeting on Association for Computational Linguistics. ACL '01. Association for Computational Linguistics, Stroudsburg, PA, p. 26
- Beck M. R. et al., 2017, *MNRAS*, submitted
- Becker A. et al., 2013, in Report on Late Winter 2013 Production: Image Differencing. LSST, Tucson, AZ
- Bellm E., 2014, in Wozniak P. R., Graham M. J., Mahabal A. A., Seaman R., eds, The Third Hot-wiring the Transient Universe Workshop. p. 27. Available at: <http://www.slac.stanford.edu/econf/C131113.1/proceedings.html>
- Bloom J. S. et al., 2012, *PASP*, 124, 1175
- Bramich D. M., 2008, *MNRAS*, 386, L77
- Brink H., Richards J. W., Poznanski D., Bloom J. S., Rice J., Negahban S., Wainwright M., 2013, *MNRAS*, 435, 1047
- Campbell H. et al., 2015, *Astron. Telegram*, 7254
- Coates A., Lee H., Ng A. Y., 2011, in Proc. 14th Int. Conf. Artif. Intell. Stat. (AISTATS), Vol. 15. Fort Lauderdale, FL
- Dark Energy Survey Collaboration, 2005, preprint ([astro-ph/0510346](https://arxiv.org/abs/astro-ph/0510346))
- Denneau L. et al., 2013, *PASP*, 125, 357
- Donalek C., Mahabal A., Djorgovski S. G., Marney S., Drake A., Glikman E., Graham M. J., Williams R., 2008, in Bailer-Jones C. A. L., ed., AIP Conf. Ser. Vol. 1082, Classification and Discovery in Large Astronomical Surveys. Am. Inst. Phys., New York, p. 252
- du Buisson L., Sivanandam N., Bassett B. A., Smith M., 2015, *MNRAS*, 454, 2026
- Geach J. E. et al., 2015, *MNRAS*, 452, 502
- Goldstein D. A. et al., 2015, *AJ*, 150, 82
- Gwyn S. D. J., 2012, *AJ*, 143, 38
- Huber M. et al., 2015, *Astron. Telegram*, 7153
- Inserra C. et al., 2013, *ApJ*, 770, 128
- Ivezić Z. et al., 2008, preprint ([arXiv:0805.2366](https://arxiv.org/abs/0805.2366))
- Johnson L. C. et al., 2015, *ApJ*, 802, 127
- Kaiser N. et al., 2010, *Proc. SPIE*, 7733, 77330E
- Keller S. C. et al., 2007, *PASA*, 24, 1
- Law N. M. et al., 2009, *PASP*, 121, 1395
- Magnier E. A. et al., 2013, *ApJS*, 205, 20
- Marshall P. J., Lintott C. J., Fletcher L. N., 2015, *ARA&A*, 53, 247
- Marshall P. J. et al., 2016, *MNRAS*, 455, 1171
- Masci F. J. et al., 2017, *PASP*, 129, 014002
- Mattila S. et al., 2016, *Astron. Telegram*, 9308
- Melchior P. et al., 2016, *Astron. Comput.*, 16, 99
- Ngiam J., Chen Z., Bhaskar S. A., Koh P. W., Ng A. Y., 2011, in Shawe-Taylor J., Zemel R., Bartlett P., Pereira F., Weinberger K., eds, Advances in Neural Information Processing Systems 24. Curran Associates, Inc., New York, p. 1125
- Rau A. et al., 2009, *PASP*, 121, 1334
- Ridgway S. T., Matheson T., Mighell K. J., Olsen K. A., Howell S. B., 2014, *ApJ*, 796, 53
- Romano R. A., Aragon C. R., Ding C., 2006, in 5th Int. Conf. on Machine Learning and Applications, ICMLA'06. IEEE Computer Society Washington, DC, p. 77
- Schwamb M. E. et al., 2012, *ApJ*, 754, 129
- Smartt S. J. et al., 2015, *A&A*, 579, A40
- Smith A. M. et al., 2011, *MNRAS*, 412, 1309
- Tonry J. L. et al., 2012, *ApJ*, 750, 99
- Tonry J., Denneau L., Stalder B., Heinze A., Sherstyuk A., Rest A., Smith K. W., Smartt S. J., 2016, *Astron. Telegram*, 8680
- Willett K. W. et al., 2013, *MNRAS*, 435, 2835
- Wright D. E. et al., 2015, *MNRAS*, 449, 451
- York D. G. et al., 2000, *AJ*, 120, 1579
- Zackay B., Ofek E. O., Gal-Yam A., 2016, *ApJ*, 830, 27
- Zevin M. et al., 2017, *Class. Quantum Grav.*, 34, 064003

This paper has been typeset from a $\text{\TeX}/\text{\LaTeX}$ file prepared by the author.