

3 INTEGRATING HUMAN AND MACHINE INTELLIGENCE IN GALAXY MORPHOLOGY
4 CLASSIFICATION TASKS

5 MELANIE BECK¹, CLAUDIA SCARLATA¹, LUCY F. FORTSON¹, CHRIS J. LINTOTT^{2, 3}, MELANIE A. GALLOWAY¹, KYLE W.
6 WILLETT¹, B. D. SIMMONS^{2,4,7}, HUGH DICKINSON¹, KAREN L. MASTERS⁵, PHILIP J. MARSHALL⁶, AND DARRYL WRIGHT²

¹Minnesota Institute for Astrophysics, University of Minnesota, Minneapolis, MN 55455, USA; beck@astro.umn.edu

²Oxford Astrophysics, Denys Wilkinson Building, Keble Road, Oxford OX1 3RH, UK

³New College, Oxford OX1 3BN, UK

⁴Center for Astrophysics and Space Sciences, Department of Physics, University of California, San Diego, CA 92093, USA

⁵Institute of Cosmology and Gravitation, University of Portsmouth, Portsmouth, UK

⁶Kavli Institute for Particle Astrophysics and Cosmology, P.O. Box 20450, MS29, Stanford, CA 94309, U.S.A.

⁷Einstein Fellow

ABSTRACT

Quantifying galaxy morphology is a challenging yet scientifically rewarding task. As the scale of data continues to increase with upcoming surveys, traditional classification methods will struggle to handle the load. We present a solution through an integration of visual and automated classifications, preserving the best features of both human and machine. We demonstrate the effectiveness of such a system through a re-analysis of visual galaxy morphology classifications collected during the Galaxy Zoo 2 (GZ2) project (Replaced: ~~, reprocessed~~ replaced with: . We reprocess the top-level question of the GZ2 decision tree) with a classification aggregation algorithm dubbed SWAP, originally developed for the Space Warps gravitational lens project. Incorporating the SWAP algorithm increases the classification rate by a factor of 4.7, classifying 226,124 galaxies in 92 days of GZ2 project time. GZ2 classified \sim 48K in the same time period. This increased rate does not diminish the quality of classification as we maintain 95.7% accuracy when compared to GZ2 published data. We next combine this with a Random Forest machine learning algorithm that learns on a suite of non-parametric morphology indicators widely used for automated morphologies. We develop a decision engine that delegates tasks between human and machine and demonstrate that the combined system provides a factor of 11.4 increase in the classification rate, classifying 210,543 galaxies in just 32 days of GZ2 project time. Again, this impressive performance is achieved with 93.5% accuracy. (Added: Though we simplify our analysis to examine a binary classification scheme, and as) the Random Forest algorithm requires a minimal amount of computational cost, this result has important implications for galaxy morphology identification tasks in the era of *Euclid* and other large-scale surveys.

Keywords: galaxies: general — galaxies: morphology — methods: data analysis — methods: machine learning

1. INTRODUCTION

Astronomers have made use of visual galaxy morphologies to understand the dynamical structure of these systems for nearly ninety years (e.g., Hubble 1936; de Vaucouleurs 1959; Sandage 1961; van den Bergh 1976; Nair & Abraham 2010; Baillard et al. 2011). The division between early-type and late-type systems corresponds, for example, to a wide range of parameters from mass and luminosity, to environment, colour, and star formation history (e.g., Kormendy 1977; Dressler 1980; Strateva et al. 2001; Blanton et al. 2003; Kauffmann et al. 2003; Nakamura et al. 2003; Shen et al. 2003; Peng et al. 2010); while detailed observations of morphologi-

cal features such as bars and bulges provide information about the history of their host systems (e.g., review by Kormendy & Kennicutt 2004; Elmegreen et al. 2008; Sheth et al. 2008; Masters et al. 2011; Simmons et al. 2014). Modern studies of morphology divide systems into broad classes (e.g., Conselice 2006; Lintott et al. 2008; Kartaltepe et al. 2015; Peth et al. 2016), but a wealth of information can be gained from identifying new and often rare classes, such as low redshift clumpy galaxies (e.g., Elmegreen et al. 2013), polar-ring galaxies (e.g., Whitmore et al. 1990), and the green peas (Caldamone et al. 2009).

While the Galaxy Zoo project has provided a solution

that scales visual classification for current surveys (Lintott et al. 2008, 2011; Willett et al. 2013, 2017; Simmons et al. 2017), producing a prolific amount of scientific output (e.g., Land et al. 2008; Bamford et al. 2009; Darg et al. 2010; Schawinski et al. 2014; Galloway et al. 2015; Smethurst et al. 2016), upcoming surveys such as *LSST* and *Euclid* will require a different approach, imaging more than a billion new galaxies (LSST Science Collaboration et al. 2009; Laureijs et al. 2011). If detailed morphologies can be extracted for just 0.1% of this imaging, we will have millions of images to contend with. A project of this magnitude would take more than sixty years to classify at Galaxy Zoo’s current rate and configuration. Standard visual morphology methods will thus be unable to cope with the scale of data.

Another approach has been the use of automated morphologies with the development of parametric (Sersic 1968; Odewahn et al. 2002; Peng et al. 2002), and non-parametric (Abraham et al. 1994; Conselice 2003; Abraham et al. 2003; Lotz et al. 2004; Freeman et al. 2013) structural indicators. While these scale well to large samples (e.g., Simard et al. 2011; Griffith et al. 2012; Casteels et al. 2014; Holwerda et al. 2014; Meert et al. 2016), they often fail to capture detailed structure and can provide only statistical morphologies with large uncertainties (e.g., Abraham et al. 1996; Bershady et al. 2000).

Machine learning techniques are becoming increasingly popular for classification and image processing tasks. Another automated approach, these generally work by defining a set of features that describe the morphology in an N -dimensional space. The location in this morphology space defines a morphological type for each galaxy. Learning the morphology space can be achieved through algorithms such as Support Vector Machines (Huertas-Company et al. 2008) or Principal Component Analysis (Watanabe et al. 1985; Scarlata et al. 2007). Another approach is through deep learning, a machine learning technique that attempts to model high level abstractions. Algorithms like convolutional and artificial neural networks (CNNs, ANNs) have been used for galaxy morphology classification with impressive accuracy (Ball et al. 2004; Banerji et al. 2010; Dieleman et al. 2015; Huertas-Company et al. 2015). A drawback to all machine learning classification techniques is the need for standardized training data, with more complex algorithms requiring more data. Furthermore, that data must be consistent for each survey: differences in resolution and depth can be inherently learned by the algorithm making their application to disparate surveys challenging.

In this work we present a system that preserves the best features of both visual and automatic classifications, developing for the first time a framework that

brings both human and machine intelligence to the task of galaxy morphology to handle the scale and scope of next generation data. We demonstrate the effectiveness of such a system through a re-analysis of visual galaxy morphology classifications collected during the Galaxy Zoo 2 project, and combine these with a Random Forest machine learning algorithm that trains on a suite of non-parametric morphology indicators widely used for automated morphologies. (Added: As a proof of concept, we focus on the first question of the Galaxy Zoo decision tree in this paper.) (Replaced: Our replaced with: We demonstrate that our) method provides a factor of 11.4 increase in the rate of galaxy morphology classification, and a factor of 10 reduction in human effort, while maintaining at least 93.5% classification accuracy as compared to Galaxy Zoo 2 published data. We first present an overview of our framework, which also serves as a blueprint for this paper.

2. GALAXY ZOO EXPRESS OVERVIEW

The Galaxy Zoo Express (GZX) framework combines human and machine to increase morphological classification efficiency, both in terms of the classification rate and required human effort. Figure 1 presents a schematic of GZX including section numbers as a short-cut for the reader. We note that transparent portions of the schematic represent areas of future work which we explore in Section 7. Any system combining human and machine classifications will have a set of generic features: a group of human classifiers, at least one machine classifier, and a decision engine which determines how these classifications should be combined.

In this work we demonstrate our system through a re-analysis of Galaxy Zoo 2 (GZ2) classifications. This allows us to create simulations of human classifiers (described in Section 3). These classifications are used most effectively when processed with SWAP, a Bayesian code described in Section 4, first developed for the Space Warps gravitational lens discovery project (Marshall et al. 2016). These subjects provide the machine’s training sample.

In Section 5, we incorporate a machine classifier. We have developed a Random Forest algorithm that trains on measured morphology indicators such as Concentration, Asymmetry, Gini coefficient and M_{20} (Added: , well-suited for the top-level question of the GZ2 decision tree, as discussed below). After a sufficient number of subjects have been classified by humans, the machine is trained and its performance assessed through cross-validation. This procedure is repeated nightly and the machine’s performance increases with size of the training sample, albeit with a performance limit. Once the machine reaches an acceptable level of performance it is applied to the remaining galaxy sample.

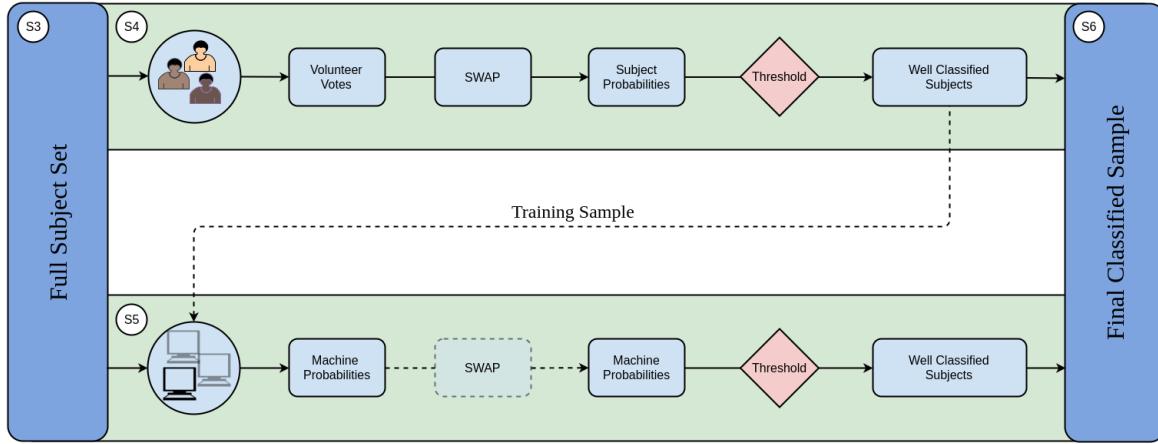


Figure 1. Schematic of our hybrid system. Humans provide classifications of galaxy images via a web interface. We simulate this with the Galaxy Zoo 2 classification data described in Section 3. Human classifications are processed with an algorithm described in Section 4. Subjects that pass a set of thresholds are considered human-retired (fully classified) and provide the training sample for the machine classifier as described in Section 5. The trained machine is applied to all subjects not yet retired. Those that pass an analogous set of machine-specific thresholds are considered machine-retired. The rest remain in the system to be classified by either human or machine. This procedure is repeated nightly. Our results are reported in Section 6.

Even with this simple description, one can see that the classification process will progress in three phases. First, the machine will not yet have reached an acceptable level of performance; only humans contribute to subject classification. Second, the machine's performance will improve; both humans and machine will be responsible for classification. Finally, machine performance will slow; remaining images will likely need to be classified by humans. These results are explored in Section 6. This blueprint allows even modest machine learning routines to make significant contributions alongside human classifiers and removes the need for ever-increasing performance in machine classification.

3. GALAXY ZOO 2 CLASSIFICATION DATA

Our simulations utilize original classifications made by volunteers during the GZ2 project. These data¹ are described in detail in Willett et al. (2013), though we provide a brief overview here. The GZ2 subject sample consists of 285,962 galaxies identified as the brightest 25% (r -band magnitude < 17) residing in the SDSS North Galactic Cap region from Data Release 7 and included subjects with both spectroscopic and photometric redshifts out to $z < 0.25$.

Subjects were shown as colour composite images via a web-based interface² wherein volunteers answered a series of questions pertaining to the morphology of the subject. With the exception of the first question, subsequent queries were dependent on volunteer responses from the previous task creating a complex decision tree

(Added: ³). Using GZ2 nomenclature, a *classification* is the total amount of information about a subject obtained by completing all tasks in the decision tree. A subject is *retired* after it has achieved a sufficient number of classifications.

For our current analysis, we choose the first task in the tree: “Is the galaxy simply smooth and rounded, with no sign of a disk?” to which possible responses include “smooth”, “features or disk”, or “star or artifact”. This (Added: choice) serves two purposes: 1) this is one of only two questions in the GZ2 decision tree that is asked of every subject thus maximizing the amount of data we have to work with, and 2) our analysis assumes a binary task and this question is simple enough to cast as such. (Added: By combining “star or artifact” responses with “features or disk” responses, we obtain a binary task.)

(Replaced: By combining the “star or artifact” vote fraction, f_{artifact} , with the “features or disk” vote fraction, f_{features} we obtain a binary response. Here, a vote fraction is simply the fraction of volunteers who voted for a particular response. We define a label for each GZ2 subject as the majority vote fraction, that is, if $f_{\text{features}} + f_{\text{artifact}} > f_{\text{smooth}}$, the galaxy is labelled ‘Featured’, otherwise it is labelled ‘Not’. replaced with: In order to compare our classification output with GZ2 we assign each subject a descriptive label. GZ2 classifications are comprised of volunteer vote fractions (f_{response}) for each response to every task in the decision tree, where vote fractions are derived from

¹ data.galaxyzoo.org

² www.galaxyzoo.org

³ A visualization of this decision tree can be found at https://data.galaxyzoo.org/gz_trees/gz_trees.html

the fraction of volunteers who voted for a particular response (more on this below). GZ2 classifications are thus continuous. A common technique is to place a threshold on these vote fractions to select samples with an emphasis on purity or completeness, depending on the science case. For our current analysis we choose a threshold of 0.5, that is, if $f_{\text{featured}} + f_{\text{artifact}} > f_{\text{smooth}}$, the galaxy is labelled ‘Featured’, otherwise it is labelled ‘Not’. We note that this threshold is not much different from the suggested 0.430 threshold in Willett et al. (2013) that produces a well-sampled subset of ‘Featured’ galaxies. Though naive, we will demonstrate throughout this paper that this threshold produces adequate results, though a more sophisticated mechanism will be explored in a future publication.) We note that only 512 subjects in the GZ2 catalog have a majority f_{artifact} , contributing less than half a percent contamination (Added: by combining the “star or artifact” with “features or disk” responses).

The GZ2 catalog (Replaced: assigns every subject three types of volunteer vote fractions: replaced with: publishes three types of vote fractions for each subject:) raw, weighted, and debiased. Debiased vote fractions are calculated to correct for redshift bias, a task that GZX does not perform. The weighted vote fractions account for inconsistent volunteers. (Deleted: , a task we perform as well.) (Replaced: However, because our mechanism is entirely different from GZ2, replaced with: The SWAP algorithm (described below) also has a mechanism to weight volunteer votes, however, the two methods are in stark contrast. In order to ensure equal footing,) we derive labels from the raw vote fractions (GZ2_{raw}) (Added: that have received no post-processing whatsoever.) In total, the data consist of over 16 million classifications from 83,943 individual volunteers.

4. EFFICIENCY THROUGH INTELLIGENT HUMAN-VOTE AGGREGATION

Galaxy Zoo 2 had a brute-force subject retirement rule whereby each galaxy was to receive approximately forty independent classifications. Once the project reached completion, inconsistent volunteers were down-weighted (Willett et al. 2013), a process that does not make efficient use of those who are exceptionally skilled. To intelligently manage subject retirement and increase classification efficiency, we adapt an algorithm from the Zooniverse project Space Warps (Marshall et al. 2016), which searched for and discovered several gravitational lens candidates in the CFHT Legacy Survey (More et al. 2016). Dubbed SWAP (Space Warps Analysis Pipeline), this algorithm predicted the probability that an image contained a gravitational lens given volunteers’ classifications and experience after being shown a training

sample consisting of simulated lensing events. We provide a brief overview here.

The algorithm assigns each volunteer an *agent* which interprets that volunteer’s classifications. Each agent assigns a 2×2 confusion matrix to their volunteer which encodes that volunteer’s probability of correctly identifying feature A given that the subject (Deleted: actually) exhibits feature A ; and the probability of correctly identifying the absence of feature A (Added: (denoted N)) given that the subject does not exhibit that feature. The agent updates these probabilities by estimating them as

$$P("X|X, \mathbf{d}) \approx \frac{\mathcal{N}_X}{\mathcal{N}_X} \quad (1)$$

(Added: where X is the true classification of the subject and “ X ” is the classification made by the volunteer upon viewing the subject. Thus) \mathcal{N}_X is the number of classifications the volunteer labelled as type X , \mathcal{N}_X is the number of subjects the volunteer has seen that were actually of type X , and \mathbf{d} represents the history of the volunteer, i.e., all subjects they have seen. (Added: Therefore the confusion matrix for a single volunteer goes as)

$$\mathcal{M} = \begin{bmatrix} P("A|N, \mathbf{d}) & P("A|A, \mathbf{d}) \\ P("N|N, \mathbf{d}) & P("N|A, \mathbf{d}) \end{bmatrix} \quad (2)$$

(Added: where probabilities are normalised such that $P("A|A) = 1 - P("N|A)$.)

Each subject is assigned a prior probability that it exhibits feature A : $P(A) = p_0$. When a volunteer makes a classification, Bayes’ theorem is used to (Replaced: derive replaced with: compute) how that subject’s prior probability should be updated into a posterior using elements of the agent’s confusion matrix. As the project progresses, (Replaced: each subject’s probability is continually updated replaced with: each subject’s posterior probability is updated after every volunteer classification), nudged higher or lower depending on volunteer input. (Replaced: Probability thresholds can be set such that subjects crossing a threshold are highly likely to exhibit the feature of interest or the absence thereof. These subjects are then considered retired. replaced with: Upper and lower probability thresholds can be set such that when a subject’s posterior crosses the upper threshold it is highly likely to exhibit feature A , while if it crosses the lower threshold it is highly likely that feature A is absent. Subjects whose posteriors cross either of these thresholds are considered retired.)

4.1. Volunteer Training Sample

A key feature of the original Space Warps project was the training of individual volunteers through the use of

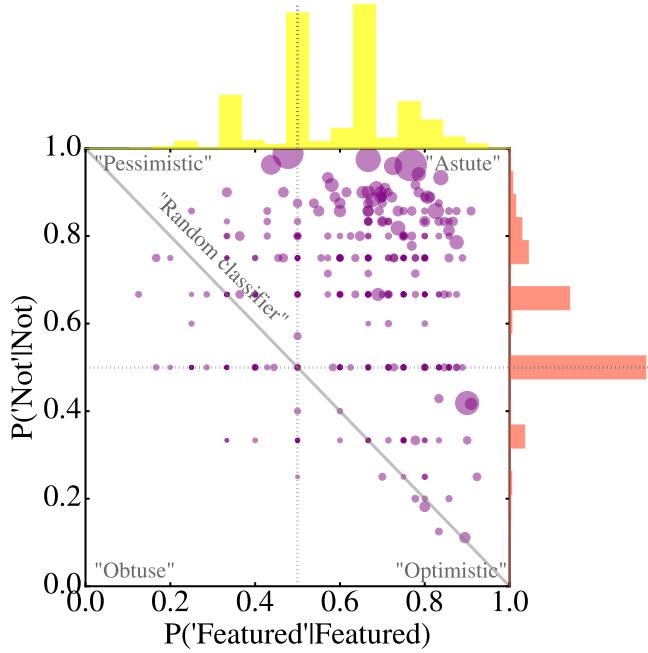


Figure 2. Confusion matrices for 1000 randomly selected GZ2 volunteers after fiducial SWAP assessment. Circle size is proportional to the number of gold standard subjects that volunteer classified. The histograms on top and right represent the distribution of each component of the confusion matrix for all volunteers. A quarter of GZ2 volunteers are “Astute”; (Replaced: they are adept at correctly identifying both ‘Featured’ and ‘Not’ subjects; replaced with: they correctly identify both ‘Featured’ and ‘Not’ subjects more than 50% of the time.) The peaks at 0.5 in both distributions are due to volunteers who see only one training image; only half of their confusion matrix is updated.

simulated images. These were interspersed with real imaging and were predominantly shown at the beginning of a volunteer’s association with the project, allowing that volunteer’s agent time to update before classifying real data. Volunteers were provided feedback in the form of a pop-up comment after classifying a training image. GZ2 did not train volunteers in such a way, which presents a challenge when applying SWAP to GZ2 classifications. (Replaced: We describe how we engineer the GZ2 data to mimic the Space Warps system. replaced with: Though we cannot retroactively train GZ2 volunteers in such a manner, we develop a gold standard sample in lieu of simulated data, and arrange the classification order of gold standard data in order to mimic the Space Warps system.)

We create a gold standard sample by selecting 3496 SDSS galaxies representative of the relative abundance of T-Types, a numerical index of a galaxy’s stage along the Hubble sequence, at $z \sim 0$ by considering galaxies that overlap with the Nair & Abraham (2010) catalogue, a collection of $\sim 14K$ galaxies classified by eye

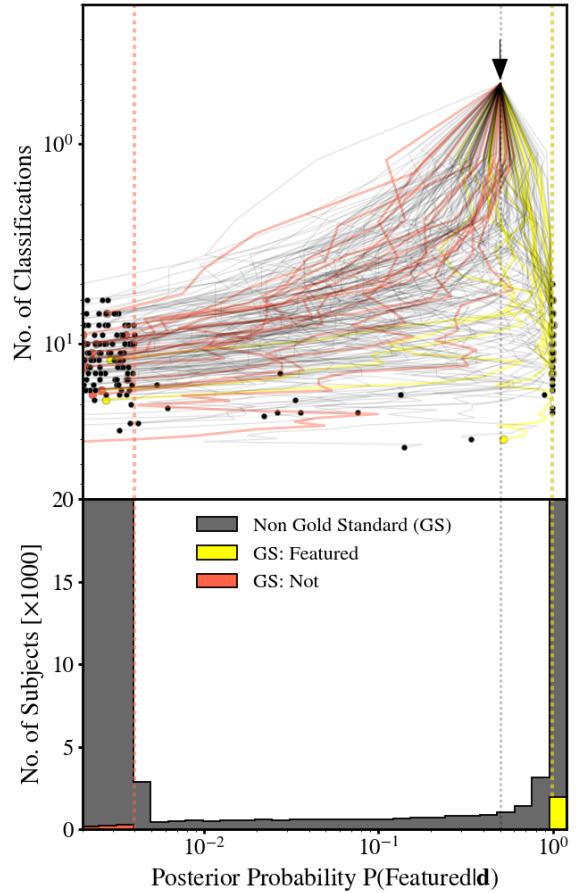


Figure 3. Posterior probabilities for GZ2 subjects. The top panel depicts the probability trajectories of 200 randomly selected GZ2 subjects. All subjects begin with a prior of 0.5 denoted by the arrow. Eaprobabilitych subject’s probability is nudged back and forth with each classification. From left to right the dotted vertical lines show the retirement threshold, prior probability, and detection threshold. Different colours denote different types of subjects. The bottom panel shows the distribution in probability for all GZ2 subjects by the end of our simulation.

into T-Types. (Replaced: Expert classifications were obtained replaced with: We must generate expert labels for these galaxies that are consistent with the labels we defined for GZ2 classifications. These are obtained) through the Zooniverse platform⁴ from 15 professional astronomers, including members of the Galaxy Zoo science team. The question posed was identical to the original GZ2 question and at least five experts classified each

⁴ The Project Builder template facility can be found at <http://www.zooniverse.org/lab>.

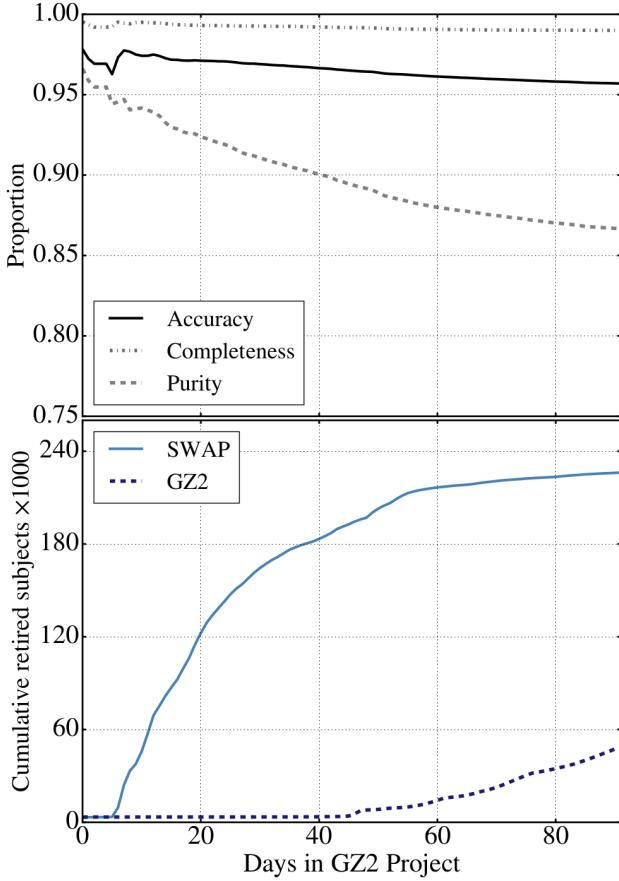


Figure 4. Fiducial SWAP simulation demonstrates a factor of 4-5 increase in the rate of subject retirement as a function of GZ2 project time (bottom panel, light blue) compared with the original GZ2 project (dashed dark blue). After 92 days, SWAP retires over 225K subjects, while GZ2 retires \sim 48K. The top panel displays the quality metrics (greys). These are calculated by comparing SWAP-assigned labels to GZ2_{raw} labels (Section 3) for the subject sample retired by that day of the simulation. Thus, on the final day, SWAP retires 226,124 subjects with 95.7% accuracy, and with completeness and purity of ‘Featured’ subjects at 99% and 86.7% respectively. The decrease in purity as a function of time is due, in part, to the fact that more difficult to classify subjects are retired later in the simulation.

galaxy. Votes are aggregated and a simple majority provides an expert label for each subject. (Added: This ensures that our expert labels are defined in exactly the same manner as the labels for all SDSS galaxies in the GZ2 sample.) Our final dataset consists of the GZ2 classifications made by those volunteers who classify at least one of these gold standard subjects. We thus retain for our simulation 12,686,170 classifications from 30,894 unique volunteers. When running SWAP, classifications of gold standard subjects are always processed first.

4.2. Fiducial SWAP simulation

Before we run a simulation, a number of SWAP parameters must be chosen: the initial confusion matrix

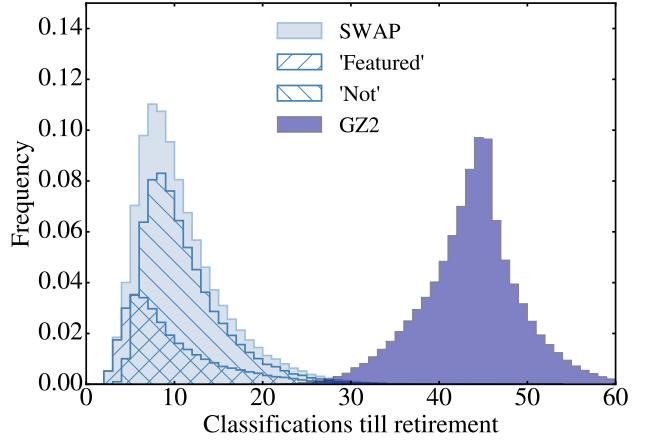


Figure 5. SWAP requires 4-5 times less human effort than GZ2 as evidenced by the distribution of the number of classifications a subject requires for retirement for the 226,124 subjects retired during our fiducial run. The GZ2 distribution peaks around 45 classifications per subject with 98.6% having at least 30 volunteer votes. In contrast, most subjects need only 9 classifications when processing with SWAP. Furthermore, ‘easy’ subjects can reach retirement in as few as 3-4 classifications.

for each volunteer’s agent, the subject prior probability, and the retirement thresholds. For our fiducial simulation, we initialize all confusion matrices at $(0.5, 0.5)$, and set the subject prior probability, $p_0 = 0.5$. We set the ‘Featured’ threshold, t_F , i.e., the minimum probability for a subject to be retired as ‘Featured’, to 0.99. Similarly, we set the ‘Not’ threshold, $t_N = 0.004$. In Appendix A we show that varying these parameters has only a small affect on the SWAP output. To simulate a live project, we run SWAP on a time step of $\Delta t = 1$ day, during which SWAP processes all volunteer classifications with timestamps within that range. This is performed for three months worth of GZ2 classification data.

Figure 2 (adapted from Figure 4 of Marshall et al. 2016) demonstrates the volunteer assessment we achieve, and shows confusion matrices for 1000 randomly selected volunteers. The circle size is proportional to the number of gold standard subjects each volunteer classified. The histograms represent the distribution of each component of the confusion matrix for all volunteers. Nearly 25% of volunteers are considered “Astute” indicating they correctly identify both ‘Featured’ and ‘Not’ subjects more than 50% of the time. (Replaced: are generally good at correctly identifying both ‘Featured’ and ‘Not’ subjects. replaced with: correctly identify both ‘Featured’ and ‘Not’ subjects more than 50% of the time.) The spikes at 0.5 in the histograms are due to volunteers who see only one gold standard subject (i.e., ‘Featured’), leaving their probability in the other (‘Not’) unchanged. Additionally, 4% of volunteers have a confusion matrix of $(0.5, 0.5)$ indicating these volunteers classified two gold standard sub-

371 jects of the same type, one correctly and one incorrectly.

372 (Added: Figure 3 (adapted from Figure 5 of Marshall
 373 et al. 2016) demonstrates how subject posterior proba-
 374 bilities are updated with each classification. The arrow
 375 in the top panel denotes the prior probability, $p_0 = 0.5$.
 376 With each classification, that prior is updated into a
 377 posterior probability thus creating a trajectory through
 378 probability space for each subject. The yellow and or-
 379 ange lines show the trajectories of a random sample of
 380 ‘Featured’ and ‘Not’ subjects from our gold standard
 381 sample, while the black lines show the trajectories of a
 382 random sample of GZ2 subjects that were not part of the
 383 gold standard sample. The vertical yellow and orange
 384 dashed lines correspond to the retirement thresholds,
 385 t_F and t_N . The lower panel shows the full distribution
 386 of GZ2 subject posteriors at the end of our simulation.
 387 An overwhelming majority of subjects cross one of these
 388 retirement thresholds.)

389 Our goal is to increase the efficiency of galaxy clas-
 390 sification. We therefore use as a metric the cumulative
 391 number of retired subjects as a function of the original
 392 GZ2 project time. We define a subject as GZ2-retired
 393 once it achieves at least 30 volunteer votes, encompass-
 394 ing 98.6% of GZ2 subjects. In contrast, a subject is
 395 considered SWAP-retired once its posterior probability
 396 crosses either of the retirement thresholds defined above.

397 However, it is important not to prioritize efficiency
 398 at the expense of quality. We thus also consider the
 399 metrics of accuracy, purity and completeness as a func-
 400 tion of GZ2 project time. These are defined as fol-
 401 lows: accuracy is the number of correctly identified
 402 subjects divided by the total number retired; comple-
 403 teness is the number of correctly identified ‘Featured’ sub-
 404 jects divided by the number of actual ‘Featured’ retired;
 405 and purity is the number of correctly identified ‘Fea-
 406 tured’ subjects divided by the number of subjects re-
 407 tired as ‘Featured’. Thus, a complete sample has no false
 408 negatives whereas a pure sample has no false positives.
 409 We compute these metrics by comparing the SWAP-
 410 assigned labels of the cumulatively retired subject set to
 411 the GZ2_{raw} labels for each day of the simulation. For
 412 example, by Day 20, SWAP retires 120K subjects with
 413 96% accuracy, 99.7% completeness, and 92% purity.

414 Figure 4 and Table 1 detail the results of our fiducial
 415 SWAP simulation compared to the original GZ2 project.
 416 The bottom panel shows the cumulative number of re-
 417 tired subjects as a function of GZ2 project time. By
 418 the end of our simulation, GZ2 (dashed dark blue) re-
 419 tires ~50K subjects while SWAP (solid light blue) re-
 420 tires 226,124 subjects. We thus classify 80% of the entire
 421 GZ2 sample in three months. (Replaced: The original
 422 GZ2 project took approximately one year to classify as
 423 many subjects, representing a factor of four increase in

424 the classification rate.) replaced with: Here we are con-
 425 sidering simply the number of classifications logged each
 426 day and not the length of time spent on a single classi-
 427 fication. Under the assumption that collapsing the GZ2
 428 decision tree to a single question would not decrease
 429 the number of classifications collected each day dur-
 430 ing the GZ2 project, processing volunteer classifications
 431 through SWAP presents a promising increase in classi-
 432 fication efficiency.) The top panel of Figure 4 demon-
 433 strates the quality of those classifications as a function
 434 of time and establishes that our full SWAP-retired sam-
 435 ple is 95.7% accurate, 99% complete, and 86.7% pure.
 436 (Added: We discuss these discrepancies in the next sec-
 437 tion.)

438 There is also a reduction in the human effort required
 439 to perform this classification task. Figure 5 shows the
 440 distribution of the number of volunteer classifications
 441 per subject achieved through SWAP (light blue) and
 442 GZ2 (dark blue) for the 226K subjects retired in our
 443 fiducial run. GZ2’s distribution peaks at ~45 indi-
 444 cating that, on average, 45 unique volunteers classify
 445 each subject. On the other hand, SWAP’s distribution
 446 peaks around 9 classifications per subject. Furthermore,
 447 subjects that are ‘easy’ to classify (i.e., ‘Featured’) re-
 448 quire even fewer classifications to reach strong consen-
 449 sus. More precisely, SWAP processes 2.3×10^6 volunteer
 450 classifications while GZ2 records $\sim 10^7$ for the same sub-
 451 ject set. SWAP reduces human effort by more than a
 452 factor of four.

453 As we demonstrate in Appendix A, varying the initial
 454 SWAP parameters from the fiducial values does not sub-
 455 stantially change the results presented here. The largest
 456 influence comes from choosing unrealistic subject prior
 457 probabilities, which can mildly degrade the quality of
 458 the resulting classifications. More importantly, none of
 459 these effects significantly alters our human and machine
 460 integration in Section 6.

461 4.3. Disagreements between SWAP and GZ2

462 Galaxy Zoo’s strength comes from the consensus of
 463 dozens of volunteers voting on each subject. Process-
 464 ing votes with SWAP reduces the number of classifica-
 465 tions to reach consensus. Though we typically recover
 466 the GZ2_{raw} label, SWAP disagrees about 5% of the time.
 467 We thus examine the false positives (subjects SWAP la-
 468 bels as ‘Featured’ but GZ2_{raw} labels as ‘Not’) and false
 469 negatives (subjects SWAP labels as ‘Not’ but GZ2_{raw} la-
 470 bels as ‘Featured’). (Added: We explore these subjects
 471 in redshift, magnitude, physical size, and concentration.
 472 We find no correlation with any of these variables, sug-
 473 gesting there are no physical reasons why SWAP’s label
 474 disagrees with GZ2_{raw}.)

475 (Replaced: We find the majority of these

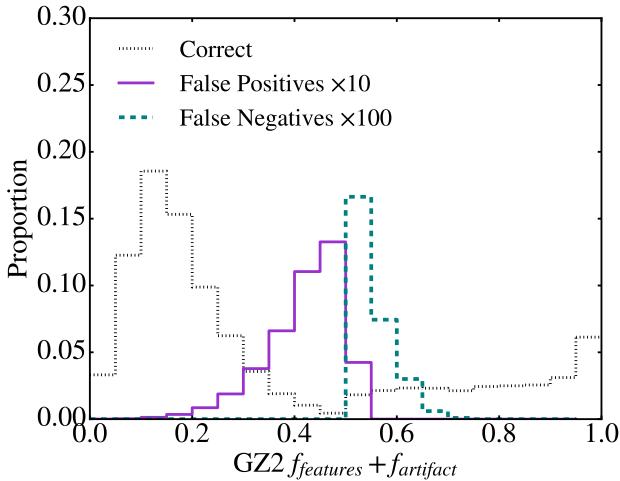


Figure 6. Distribution of GZ2 $f_{\text{features}} + f_{\text{artifact}}$ vote fractions for subjects correctly identified by SWAP (dotted grey), along with those identified as false positives (solid purple), and false negatives (dashed teal). The false positives and false negatives are scaled by factors of 10 and 100 respectively for easier comparison. From Section 3, subjects with values > 0.5 are defined as ‘Featured’, however, the teal distribution indicates that SWAP labels them as ‘Not’. This is not a flaw of SWAP: 68.9% of incorrectly identified subjects have $0.4 \leq f_{\text{features}} + f_{\text{artifact}} \leq 0.6$ suggesting that GZ2_{raw} labels are simply too uncertain. The overlap between the false positives and negatives is due to subjects that are exactly 50-50; by default these are labelled ‘Not’.

disagreements are due to uncertainties in the GZ2_{raw} label. replaced with: Instead we consider the errors associated with the GZ2 vote fraction, which can be estimated as binomial. Let success be a response of “smooth” and failure be any other response. The 68% confidence interval on a subject with $f_{\text{smooth}} = 0.5$ is then (0.42, 0.57) assuming 40 classifications, each with a probability of 0.5.) Figure 6 shows the distribution of $f_{\text{features}} + f_{\text{artifact}}$ for the false positives (solid purple), and the false negatives (dashed teal) compared to the subjects where SWAP and GZ2 agree (dotted grey). Recall that if this value is greater than 0.5, the subject is labelled ‘Featured’. The majority of (Replaced: incorrectly labelled replaced with: disagreements between SWAP and GZ2 are for) subjects that have $0.4 \leq f_{\text{features}} + f_{\text{artifact}} \leq 0.6$. (Replaced: indicating that the GZ2 raw vote fractions are simply too uncertain to provide high-quality labels. replaced with: It is thus unsurprising that SWAP and GZ2 disagree most within the errors of GZ2 vote fractions.) We note that the distribution overlap is due to subjects that do not have a majority; these are labelled ‘Not’ by default.

Two other effects contribute to the disagreement between SWAP and GZ2. First, as the number of classifications used to retire a galaxy decreases, the likelihood of misclassification by random chance increases.

Second, disagreement arises due to expert-level volunteers whose confusion matrices are close to 1.0. These volunteers are essentially more strongly weighted, allowing that subject’s posterior to cross a retirement threshold in as few as two classifications. In rare cases, despite training, some expert-level volunteers get it wrong (Added: compared to expert classifications). These issues can be mitigated by requiring each subject reach a minimum number of classifications before allowing its probability to cross a threshold, thus combining the best qualities of GZ2 and SWAP.

4.4. Summary

We demonstrate a factor of four or more increase in classification efficiency while maintaining 95% accuracy, nearly perfect completeness of ‘Featured’ subjects, and with a purity that can be controlled by careful selection of input parameters to be better than 90% (see Appendix A). Exploring those subjects wherein SWAP and GZ2 disagree, we conclude that the majority of this disagreement stems from uncertainty in GZ2_{raw} labels. We now turn our focus towards incorporating a machine classifier utilizing these SWAP-retired subjects as a training sample.

5. EFFICIENCY THROUGH INCORPORATION OF MACHINE CLASSIFIERS

We construct the full Galaxy Zoo Express by incorporating supervised learning, the machine learning task of inference from labelled training data. The training data consist of a set of training examples, and must include an input feature vector and a desired output label. Generally speaking, a supervised learning algorithm analyzes the training data and produces a function that can be mapped to new examples. An optimized algorithm will correctly determine class labels for unseen data. By processing human classifications through SWAP, we obtain a set of binary labels by which we can train a machine classifier. We briefly outline the technical details of our machine below, turning towards the decision engine we develop in Section 5.4.

5.1. Random Forests

We use a Random Forest (RF) algorithm (Breiman 2001), an ensemble classifier that operates by bootstrapping the training data and constructing a multitude of individual decision tree algorithms, one for each subsample. An individual decision tree works by deciding which of the input features best separates the classes. It does this by performing splits on the values of the input feature that minimize the classification error. These feature splits proceed recursively.

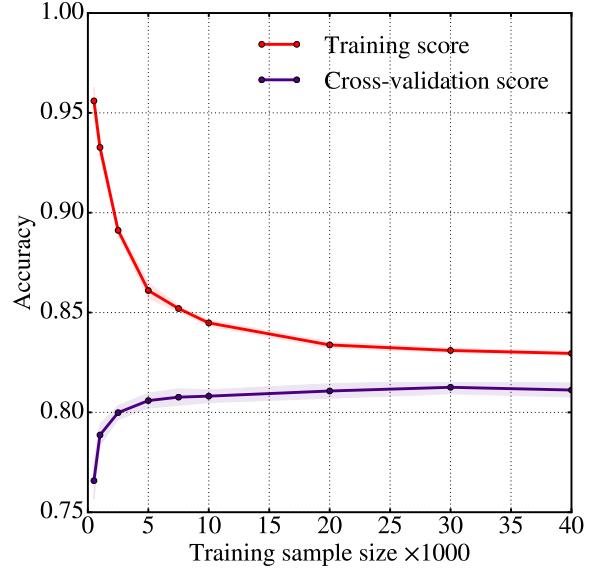
552 Decision trees alone are prone to overfitting, precluding them from generalizing well to new data. Random
 553 Forests mitigate this effect by combining the output labels from a multitude of decision trees. Specifically,
 554 we use the `RandomForestClassifier` from the Python
 555 module `scikit-learn` (Pedregosa et al. 2011).

558 5.2. Grid Search and Cross-validation

559 Of fundamental importance is the task of choosing an algorithm’s hyperparameters, values which determine how the machine learns. For a RF, key
 560 quantities include the maximum depth of individual trees (`max_depth`), the number of trees in the forest
 561 (`n_estimators`), and the number of features to consider
 562 when looking for the best split (`max_features`). The
 563 goal is to determine which values will optimize the machine’s performance and thus these values cannot be chosen *a priori*. We perform a grid search with k -fold cross-validation whereby the training sample is split into k
 564 subsamples. One subsample is withheld to estimate the machine’s performance while the remaining data is used
 565 to train the machine. This is performed k times and the average performance value is recorded. The entire process is repeated for every combination of the hyperparameters in the grid space and values that optimize the output are chosen. In this work we let $k = 10$, however, we leave this as an adjustable input parameter. In the interest of computational speed, we set `n_estimators` = 30 and perform the grid search for `max_depth` over the range [5, 16], and `max_features` over the range [\sqrt{D} , D], where D is the number of features in the feature vector, described below.

583 5.3. Feature Representation and Pre-Processing

584 The feature vector on which the machine learns is composed of D individual numeric quantities associated with the subject that the machine uses to discern that subject from others in the training sample. To segregate ‘Featured’ from ‘Not’, we draw on ZEST (Scar-
 585 lata et al. 2007) and compute concentration, asymmetry, Gini coefficient, and M_{20} , the second-order moment of light for the brightest 20% of galaxy pixels as measured from SDSS DR12 *i*-band imaging (see Appendix B). Coupled with SExtractor’s measurement of ellipticity (Bertin & Arnouts 1996), we provide the machine with a $D = 5$ dimensional morphology parameter space. These non-parametric diagnostics have long been used to quantify galaxy morphology in an automated fashion (e.g., Abraham et al. 1996; Bershady et al. 2000; Conselice et al. 2000; Abraham et al. 2003; Conselice 2003; Lotz et al. 2004; Snyder et al. 2015). Because the RF algorithm handles a variety of input formats, the only pre-processing step we perform is the removal of poorly-measured morphological indicators, i.e. catas-



566 **Figure 7.** Learning curve for a Random Forest with fixed hyperparameters. (Replaced: The training score is the accuracy of the trained machine applied to its own training sample. The cross-validation score is the accuracy of the machine computed during the cross-validation process. replaced with: These curves show the mean accuracy computed during cross-validation and on the training sample, where the shaded regions denote the standard deviation.) When the training sample size is small, the machine accurately identifies its own training sample but is unable to generalize to unseen data as evidenced by a low cross-validation score. As the training sample size increases, the cross-validation score increases. This behavior plateaus indicating that larger training samples provide little in additional performance.

604 trophic failures.

605 5.4. Decision Engine

606 A number of decisions must be addressed before attempting to train the machine. In particular, which subjects should be designated as the training sample? When should the machine attempt its first training session? When has the machine’s performance been optimized such that it will successfully generalize to unseen subjects? The field of machine learning provides few hard rules for answering these questions, only guidelines and best practices. Here we briefly discuss our approach for the development of our decision engine.

607 As discussed in detail in Section 4, SWAP yields a probability that a subject exhibits the feature of interest. While some machine algorithms can accept continuous input labels, the RF requires distinct classes. We thus use only those subjects which have crossed either of the retirement thresholds. Though we find that SWAP consistently retires 35-40% ‘Featured’ subjects on any given day of the simulation, a balanced ratio of ‘Featured’ to ‘Not’ isn’t guaranteed. Highly unbalanced training samples should be resampled to correct

the imbalance; however, as we exhibit only a mild lopsidedness, we allow the machine to train on all SWAP-retired subjects.

SWAP retires a few hundred subjects during the first days of the simulation. In principle, a machine can be trained with such a small sample, but will be unable to generalize to unseen data. We estimate a minimum number of training samples and the machine’s ability to generalize by considering a learning curve, an illustration of a machine’s performance with increasing sample size for fixed hyperparameters. Figure 7 demonstrates such a curve wherein we plot the accuracy from both the 10-fold cross-validation, and the trained machine applied to its own training sample for a random sample of GZ2 subjects required to be balanced between ‘Featured’ and ‘Not’. We fix the RF’s hyperparameters as follows: `max_depth = 8`, `n_estimators = 30`, and `max_features = 2`. When the sample size is small, the cross-validation score is low and the training score is high, a clear sign of over-fitting. However, as the training sample size increases, the cross-validation score increases and eventually plateaus, indicating that larger training sets will yield little additional gain.

We estimate this plateau begins when the training sample reaches 10,000 subjects and require SWAP retire at least this many before the machine attempts its first training. We estimate the machine has trained sufficiently if the cross-validation score fluctuates by less than 1% for three consecutive nights of training to ensure we have reached the plateau. This requires that we record the machine’s training performance each night, including how well it scores on the training sample, the cross-validation score, and the best hyperparameters.

5.5. The Machine Shop

We can now describe a full GZX simulation, which begins with human classifications processed through SWAP for several days. Once at least 10K subjects have been retired, their feature vectors are passed to the machine for its inaugural training. A suite of performance metrics are recorded by a machine agent, similar in construction to SWAP’s agents. This agent determines when the machine has trained sufficiently by assessing the variation in performance metrics for all previous nights of training. Once the machine has been optimized, the agent introduces it to the test sample consisting of any subject that has not yet reached retirement through SWAP and is not part of the gold standard sample.

Analogous to SWAP, we generate a retirement rule for machine-classified subjects. In addition to the class prediction, the RF algorithm computes the probability for each subject to belong to each class. This probability is simply the average of the probabilities of each individ-

ual decision tree, where the probability of a single tree is determined as the fraction of subjects of class X on a leaf node. Only subjects that receive a class prediction with $p_{\text{machine}} \geq 0.9$ are considered retired. The remaining subjects have the possibility of being classified by humans or the machine on a future night of the simulation. This constitutes the core of our passive feedback mechanism. Subjects that are not retired by the machine can instead be retired by humans, thus providing the machine a more fully sampled morphology parameter space on future training sessions.

6. RESULTS

We perform a full GZX simulation incorporating our RF with the fiducial SWAP run discussed in Section 4.2. The machine attempts its first training on Day 8 with an initial training sample of $\sim 20K$ subjects. It undergoes several additional nights of training, each time with a larger training sample. By Day 12, SWAP has provided over 40K subjects for training and the machine’s agent has deemed the machine optimized. The machine predicts class labels for the remaining 230K GZ2 subjects. Of those, the machine retires over 70K, dramatically increasing the subset of retired subjects. We end the simulation after 32 days, having retired $\sim 210K$ subjects as detailed in Table 1.

We present these results in Figure 8 where subject retirement with GZX (red) is compared to our fiducial SWAP-only run (light blue) and GZ2 (dashed dark blue). Using the GZ2_{raw} labels as before, we compute our usual quality metrics on the full sample of GZX-retired subjects; reported in Table 1. Accuracy and purity remain within a few percent of the SWAP-only run at 93.5% and 84.2% respectively. Instead we see a 5% decline in the completeness. While the SWAP-only run identified 99% of ‘Featured’ subjects, incorporation of the machine seems to miss a significant portion thus dropping GZX completeness to 94.3%. We discuss this behavior below.

By dynamically generating a training sample through a more sophisticated analysis of human classifications coupled with a machine classifier, we retire more than 200K GZ2 subjects in just 27 days. Visual classification through SWAP alone retires as many in 50 days, while GZ2 requires a full year. (Replaced: ~~GZX thus provides an order of magnitude increase in the rate of classification over the traditional crowd-sourced approach.~~ replaced with: Though our analysis considers the only the top-level task of GZ2’s decision tree, GZX suggests a tantalizing potential to increase the classification rate by an order of magnitude over the traditional crowd-sourced approach.) We next explore the composition of those classifications.

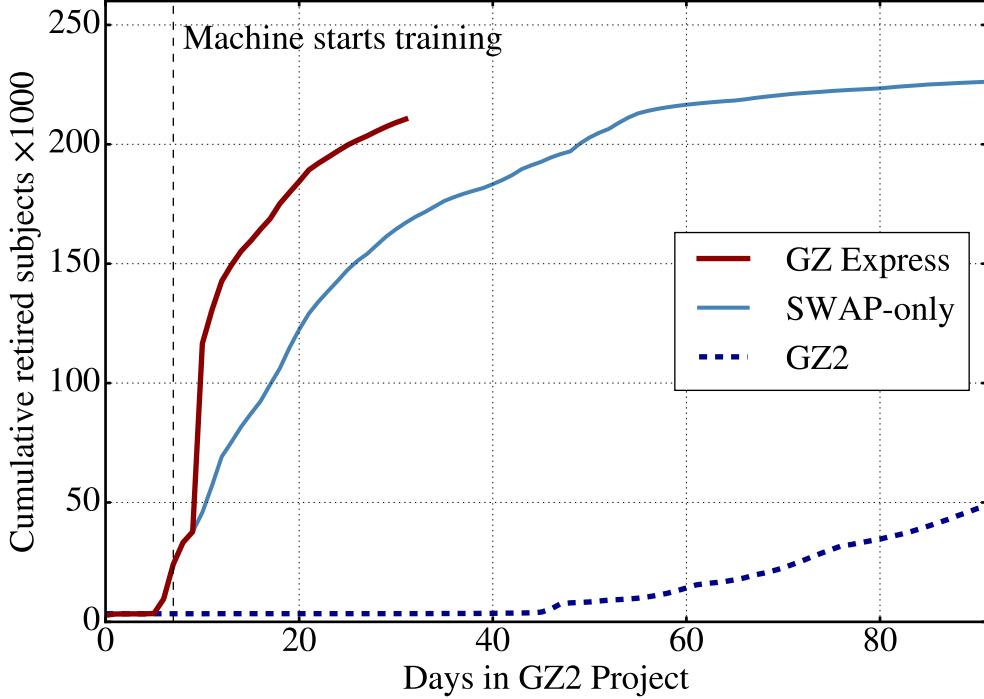


Figure 8. By incorporating a machine classifier, GZX (red) increases the classification rate by an order of magnitude compared to GZ2 (dashed dark blue) and out-performs the SWAP-only run (light blue), retiring more than 200K subjects in just 27 days of GZ2 project time. The dashed black line marks the first night the machine trains. After several additional nights of training, it is deemed optimized and allowed to retire subjects. Both humans and machine then contribute to retirement. We end the simulation after 32 days having retired over 210K galaxies. See Table 1 for details.

Table 1. Summary of key quantities for GZ2 and our various simulations. All quality metrics are calculated using GZ2_{raw} labels.

Simulation Summary						
	Days	Subjects Retired	Human Effort (classifications)	Accuracy (%)	Purity (%)	Completeness (%)
Galaxy Zoo 2	430	285962	16,340,298	–	–	–
SWAP only	92	226124	2,298,772	95.7	86.7	99.0
SWAP+RF	32	210543	932,017	93.5	84.2	94.3

6.1. Who retires what, when?

In the top panel of Figure 9 we explore the individual contributions to GZX subject retirement from the RF (dash-dotted teal) and SWAP (dashed orange). The solid black line shows the total GZX retirement (SWAP+RF), while the dotted grey line depicts the fiducial SWAP-only run from Section 4.2 for reference. Two things are immediately obvious. First, each component shoulders approximately half of the retirement burden with the machine and SWAP responsible for $\sim 100\text{K}$ and $\sim 110\text{K}$ subjects respectively. Secondly, the rate of retirement exhibited by the two components is in stark contrast. SWAP retires at a relatively constant rate while the machine retires dramatically at the beginning of its application, quickly surpassing the human contri-

bution, and plateaus thereafter.

We thus clearly see three epochs of subject retirement, as we presumed. In the first phase, humans are the only contributors to subject retirement. Once the machine is optimized, it immediately contributes more to retirement than humans. However, the machine's performance plateaus quickly; the third phase is again dominated by human classifications.

In the bottom panels of Figure 9, we consider the class composition of subjects retired by SWAP and the RF. The left (right) panel shows the retired fraction of GZ2 subjects identified as ‘Featured’ (‘Not’) according to their GZ2_{raw} labels as a function of GZ2 project time. Overall, GZX retires 73.6% of the GZ2 subject sample and this is evenly distributed between ‘Featured’ and ‘Not’ subjects as indicated by the solid black

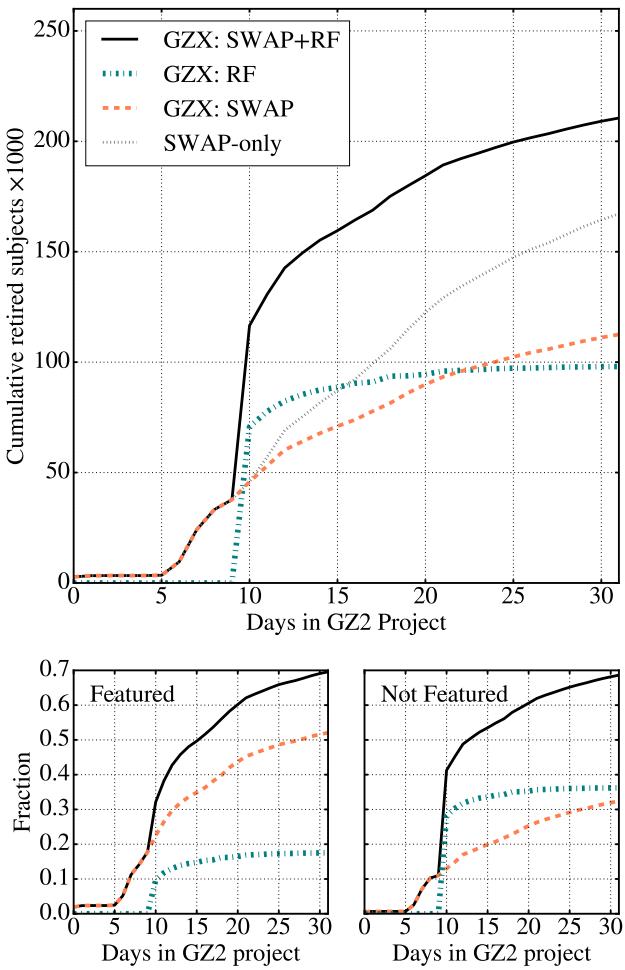


Figure 9. Contributions to subject retirement by both classifying agents of GZX: human (SWAP) and machine. The top panel shows cumulative subject retirement for GZX as a whole (solid black), along with that attributed to the RF (dash-dotted teal), and SWAP (dashed orange). The dotted grey line shows the fiducial SWAP-only run for comparison. Retirement totals for humans and machine are nearly equal over the course of the simulation but display different behaviors: SWAP's retirement rate is almost constant while the RF contributes substantially after its initial application and then plateaus. The bottom panels show what fraction of GZ2 subjects are retired, separated by class label. Overall, GZX retires 73.6% of the entire GZ2 sample in 32 days, retiring the same proportion of ‘Featured’ and ‘Not’ subjects as indicated by the black lines. However, humans retire 30% more ‘Featured’ subjects than the machine, while both components retire a similar proportion of ‘Not’ subjects.

lines in both panels. However, SWAP retires more than 50% of all ‘Featured’ subjects while the machine retires only 18%. This divergence does not exist for ‘Not’ subjects where each component contributes 33-37%.

What is the source of this discrepancy? Each night the machine trains on a sample composed consistently of 30-40% ‘Featured’ subjects but does not retire a similar proportion, indicating that the 30% of non-retired ‘Featured’ subjects do not receive high p_{machine} . In the fol-

lowing section we explore whether this is an artifact of our choice in machine or in the human-machine combination implemented here.

6.2. Machine performance

Throughout our analysis we have defined ‘Featured’ and ‘Not’ subjects by their GZ2_{raw} labels as this was the most compatible choice for comparison with SWAP output. However, the machine does not learn in the same way, nor is it presented with the same information. We argue that the machine classifications are valid and complimentary to human classifications.

Of the 6127 subjects that were deemed false positives, i.e., galaxies retired by the machine as ‘Featured’ that have ‘Not’ GZ2_{raw} labels, we visually examine several hundred and assess that, to the expert eye, a majority are, in fact, ‘Featured’. A random sample is shown in Figure 10, where the value in the lower left corner is the raw GZ2 smooth vote fraction, f_{smooth} ; the fraction of volunteers who classified that subject as ‘Not’. This small sample consists predominantly of edge-on disks and disk galaxies with low surface brightness features.

That the machine can identify ‘Featured’ galaxies that humans classify as ‘Not’ has two contributing factors: 1) the first task of the GZ2 decision tree asks a very specific question that does not necessarily correlate with a split between early- and late-type galaxies, and 2) the machine learns on morphology diagnostics that are very different from visual inspection. Regarding the first point, the full sample of false positives has $\langle f_{\text{smooth}} \rangle = 0.645 \pm 0.106$ with 56.7% having $f_{\text{smooth}} \leq 0.65$. This indicates that volunteers have not reached a strong consensus for a majority of these subjects; behavior that could be modified by providing actual training images and live feedback as performed in Marshall et al. (2016). We note that 71.7% of these galaxies are labelled ‘Featured’ after GZ2’s debiasing process.

The second point suggests that, in some cases, the morphology indicators we measure are sufficient for the machine to recognize ‘Featured’ galaxies regardless of the labels humans provide. Figure 11 shows the distribution of each morphology indicator for all subjects the machine retires as ‘Featured’ (yellow) and ‘Not’ (not yellow) compared to the full GZ2 subject set. The difference between ‘Featured’ and ‘Not’ is stark in all but the M_{20} distribution. This can be seen explicitly in Figure 12 where we show the RF’s ranked feature importances with large values indicating higher importance. Feature importance is computed as how much each feature decreases the impurity of a split in a tree. The impurity decrease from each feature is then averaged over all trees and ranked. We show the feature importance averaged over all nights of training with black bars indicating the standard deviation. The machine finds the

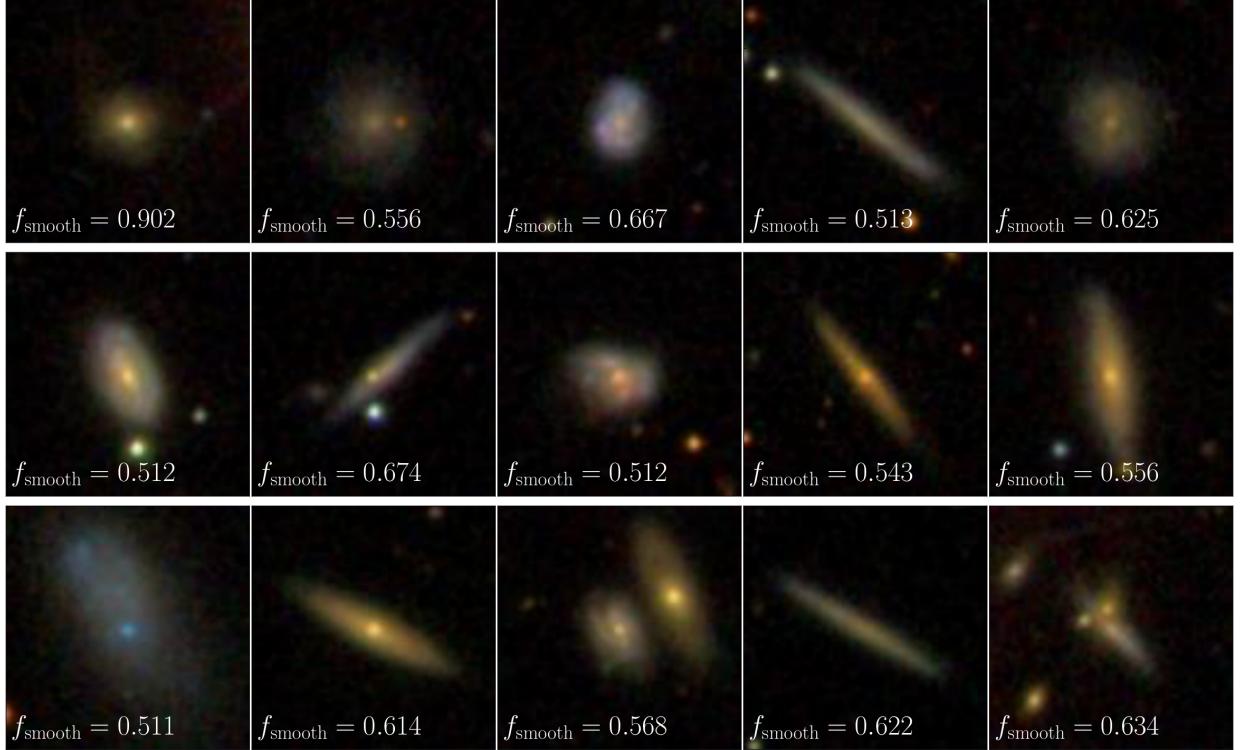


Figure 10. A random subsample of subjects identified as false positives: labelled by machine as ‘Featured’, but as ‘Not’ according to GZ2_{raw}. We display f_{smooth} in the lower left corner, that is, the fraction of volunteers who classified the subject as ‘smooth’ (‘Not’). Values are typically between 0.5 and 0.65 indicating that GZ2 volunteers did not reach a strong consensus. Fortunately, the machine is able to identify these subjects as ‘Featured’ due to their measured morphology diagnostics.

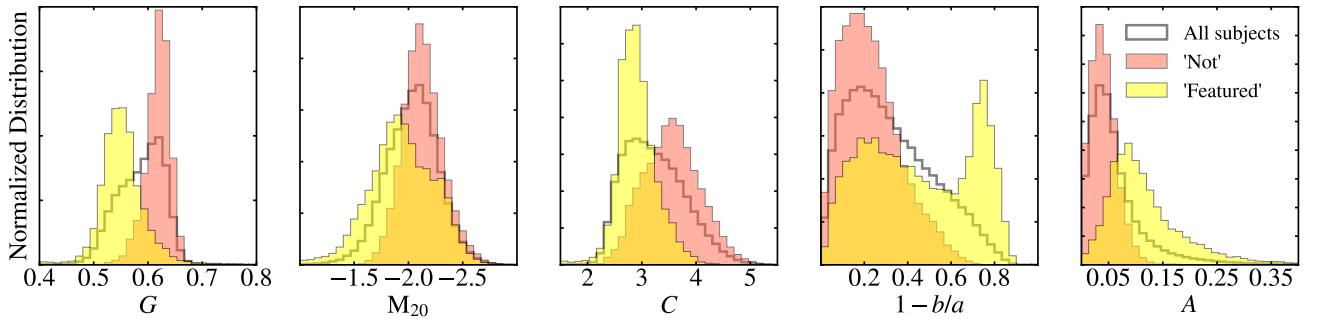


Figure 11. The RF is trained on a 5-dimensional morphology parameter space. We show the distribution of each morphology indicator for machine-retired ‘Featured’ (yellow) and ‘Not’ (not yellow) subjects compared to the full GZ2 subject sample (black). The difference between ‘Featured’ and ‘Not’ subjects is in stark contrast for all distributions except, perhaps, M_{20} .

824 Gini coefficient most important for class prediction, and
 825 places little emphasis on M_{20} . It is well known that the
 826 Gini coefficient is more sensitive to noise than other di-
 827 agnostics, however, we point out that when a machine is
 828 faced with two or more correlated features, any of them
 829 can be used as the predictor. Once chosen, the impor-
 830 tance of the others is reduced. This explains why Con-
 831 centration is ranked much lower than Gini even though
 832 they are strongly correlated as seen in Figure A2. That
 833 the machine relies heavily on these two morphology diag-

834 nistics is unsurprising as concentration has long been an
 835 automated predictor between early- and late-type galax-
 836 ies (Abraham et al. 1994, 1996; Shen et al. 2003).

837 The complementary nature of human and machine
 838 classification can best be utilized by a feedback mech-
 839 anism in which a portion of machine-retired subjects
 840 are reviewed by humans. Subjects that display exces-
 841 sive disagreement should be verified by an expert (or
 842 expert-user). In the same way that humans increase the
 843 machine’s training sample over time, subjects that the

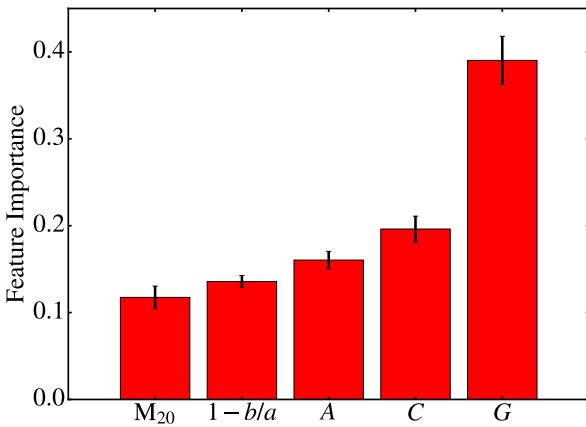


Figure 12. The RF’s ranked feature importance averaged over all nights of training with black bars indicating the standard deviation. A larger value corresponds to higher importance. The machine computes feature importance according to how much each feature increases the purity of the resulting split averaged over all trees in the forest. The RF places great importance in the Gini coefficient though we note that it can under-represent the importance of highly correlated features such as concentration.

844 machine properly identifies can become part of the humans’ training sample.

7. LOOKING FORWARD

845 We have demonstrated the first practical framework
846 for combining human and machine intelligence in galaxy
847 morphology classification tasks. While we focus below
848 on a brief discussion of our next steps and potential
849 applications to large upcoming surveys, we note that
850 our results have implications for the future of citizen
851 science and Galaxy Zoo in particular.

852 GZX is perhaps one of the simplest ways to combine
853 human and machine intelligence and its impressive per-
854 formance motivates a higher level of sophistication. A
855 first step will be an implementation of SWAP that can
856 handle a complex decision tree. In addition, we envision
857 multiple forms of active feedback in addition to our pas-
858 sive feedback mechanism. SWAP allows us to leverage
859 the most skilled volunteers to review galaxies difficult
860 for either human or machine to classify. Additionally,
861 machine-retired subjects should contribute to the train-
862 ing sample for humans in an analogous fashion to what
863 we have already implemented.

864 Secondly, our RF can be improved by providing it in-
865 formation equal to what humans receive: multi-band
866 morphology diagnostics will be included in our future
867 feature vector. However, the Random Forest algorithm
868 is not easily adapted to handle measurement errors or
869 class labels with continuous distributions. To fully uti-
870 lize the information provided by SWAP, sophisticated
871 algorithms such as deep convolutional neural networks

872 (CNN) or Latent Dirichlet allocation (LDA), an algo-
873 rithm that is frequently used in document processing,
874 should be considered. Furthermore, there is no reason
875 to limit to a single machine. As hinted at in Figure 1,
876 several machines could train simultaneously, their pre-
877 dictions aggregated through SWAP, creating an on-the-
878 fly machine ensemble.

879 With the above upgrades implemented, we expect per-
880 formance of both the classification rate and quality to
881 further increase. However, even our current implemen-
882 tation can cope with upcoming data volumes from large
883 surveys. By some estimates, *Euclid* is expected to ob-
884 tain measurable morphology with its visual instrument
885 (VIS) for approximately $10^6 - 10^7$ galaxies (Laureijs
886 et al. 2011). Visual classification at the rate achieved
887 with Galaxy Zoo today would require 12–120 years to
888 classify.⁵ If the *Euclid* sample is on the high end, GZX as
889 currently implemented could classify the brightest 20%
890 during the six years of its observing mission. As cur-
891 rently implemented, we obtain accuracy around 95% po-
892 tentially leaving hundreds of thousands of galaxies with
893 unreliable classifications. In a companion paper that
894 seeks to identify supernovae, Wright et al. (submitted)
895 demonstrate a dramatic increase in accuracy through an
896 entirely different human-machine combination whereby
897 the scores from human and machine are averaged to-
898 gether with the combined score yielding the most reli-
899 able classification. Again, a combination of both ap-
900 proaches will allow us to take full advantage of legacy
901 output from large scale surveys.

7.1. Conclusions

902 In this paper we design and test Galaxy Zoo Express,
903 an innovative system⁶ for the efficient classification of
904 galaxy morphology tasks that integrates the native abil-
905 ity of the human mind to identify the abstract and novel
906 with machine learning algorithms that provide speed
907 and brute force. We demonstrate for the first time that
908 the SWAP algorithm, originally developed to identify
909 rare gravitational lenses in the Space Warps project, is
910 robust for use in galaxy morphology classification. We
911 show that by implementing SWAP on GZ2 classifica-
912 tion data we can increase the rate of classification by a fac-
913 tor of 4–5, requiring only 90 days of GZ2 project time to
914 classify nearly 80% of the entire galaxy sample.

915 Furthermore, we have implemented and tested a Ran-
916 dom Forest algorithm and developed a decision engine
917 that delegates tasks between human and machine. We

⁵ We note that the classification rate of GZ2 was 4 times higher than GZ’s current steady rate.

⁶ Our code can be found at <https://github.com/melaniebeck/GZExpress>

921 show that even this simple machine is capable of providing significant gains in the classification rate when
 922 combined with human classifiers: GZX retires over 70%
 923 of GZ2 galaxies in just 32 days of GZ2 project time. This
 924 represents a factor of 11.4 increase in the classification
 925 rate as well as an order of magnitude reduction in hu-
 926 man effort compared to the original GZ2 project. This
 927 is achieved without sacrificing the quality of classifica-
 928 tions as we maintain accuracy well above 90% through-
 929 out our simulations. Additionally, we have shown that
 930 training on a 5-dimensional parameter space of tradi-
 931 tional non-parametric morphology indicators allows the
 932 machine to identify subjects that humans miss, provid-
 933 ing a complementary approach to visual classification.
 934 The gain in classification speed allows us to tackle the
 935 massive amounts of data soon to be forthcoming from
 936 large surveys like *LSST* and *Euclid*.

938 MB thanks Steven Bamford and Boris Häußler for in-
 939 sightful discussions on citizen science and Galaxy Zoo;
 940 and John Wallin and Marc Huertas-Company for sev-
 941 eral enlightening conversations on machine learning and
 942 classification. We are grateful to Elisabeth Baeten, Mi-
 943 caela Bagley, Karlen Shahinyan, Vihang Mehta, Steven
 944 Bamford, Kevin Schawinski, and Rebecca Smethurst

945 for providing expert classifications in addition to those
 946 provided by the authors. PJM acknowledges Aprajita
 947 Verma and Anupreeta More for their ongoing collabora-
 948 tion on the Space Warps project.

949 MB, CS, LF, KW, and MG gratefully acknowledge
 950 support from the US National Science Foundation Grant
 951 AST-1413610. MB acknowledges additional support
 952 through New College and Oxford University's Balzan
 953 Fellowship as well as the University of Minnesota Doc-
 954 toral Dissertation Fellowship. Travel funding was sup-
 955 plied to MB, in part, by the University of Minnesota
 956 Thesis Research Travel Grant. CJL recognizes support
 957 from a grant from the Science & Technology Facilities
 958 Council (ST/N003179/1). BDS acknowledges support
 959 from Balliol College, Oxford, and the National Aeronau-
 960 tics and Space Administration (NASA) through Einstein
 961 Postdoctoral Fellowship Award Number PF5-160143 is-
 962 sued by the Chandra X-ray Observatory Center, which is
 963 operated by the Smithsonian Astrophysical Observatory
 964 for and on behalf of NASA under contract NAS8-03060.
 965 The work of PJM is supported by the U.S. Department
 966 of Energy under contract number DE-AC02-76SF00515.

967 *Software:* scikit-learn ([Pedregosa et al. 2011](#)), As-
 968 trophy ([Astropy Collaboration et al. 2013](#)), TOPCAT
 969 ([Taylor 2005](#))

APPENDIX

A. EXPLORING SWAP'S PARAMETER SPACE

971 In this Appendix we explore the SWAP parameter space and assess the effects on subject retirement.
 972 **Initial agent confusion matrix.** In our fiducial simulation each volunteer was assigned an agent whose confusion
 973 matrix was initialized at (0.5, 0.5), which presumes that volunteers are no better than random classifiers. We perform
 974 two simulations wherein we initialize agent confusion matrices as (0.4, 0.4), slightly obtuse volunteers; and (0.6, 0.6),
 975 slightly astute volunteers, with everything else remaining constant. Results of these simulations compared to the
 976 fiducial run are shown in the left panel of Figure A1. We find that SWAP is largely insensitive to the initial confusion
 977 matrix both in terms of the subject retirement rate and classification quality.

978 We retire $\sim 225K \pm 3.5\%$ subjects as shown by the light blue shaded region in the bottom left panel of Figure A1,
 979 where the dashed blue line denotes the fiducial run. Predictably, when the confusion matrix probabilities are low, we
 980 retire fewer subjects than when these probabilities are high for a given period of time. This is easy to understand since
 981 it takes longer for volunteers to become astute classifiers when they are initially given values denoting them as obtuse.
 982 Regardless, most volunteers become astute classifiers by the end of the simulation. The top left panel demonstrates
 983 our usual quality metrics as computed in Section 4.2. The dashed lines again denote the fiducial run. We maintain
 984 $\sim 95\%$ accuracy, 99% completeness, and $\sim 84\%$ purity; and no metric changes by $> 2\%$ regardless of initial confusion
 985 matrix values.

986 This spread is due to three effects: 1) subjects can receive an alternate SWAP label in different simulations, 2)
 987 subjects can be retired in a different order, and 3) the set of retired subjects is not guaranteed to be common to all
 988 runs. We find SWAP to be highly consistent: more than 99% of retired subjects are the same among all simulations,
 989 and, of these, 99% receive the same label. Instead we find that the order in which subjects are retired changes between
 990 runs. When the confusion matrix is low, subjects take longer to classify compared to the fiducial run (i.e., they retire
 991 on a later date in GZ2 project time). Likewise, subjects retire sooner when the confusion matrix is high. This can
 992 cause quality metrics to vary since they are calculated on a day to day basis. These effects each contribute less than
 993 one per cent variation and thus we see a high level of consistency between simulations.

994 Of interest, perhaps, is that the quality metrics for these simulations are not symmetric about the fiducial run. How-

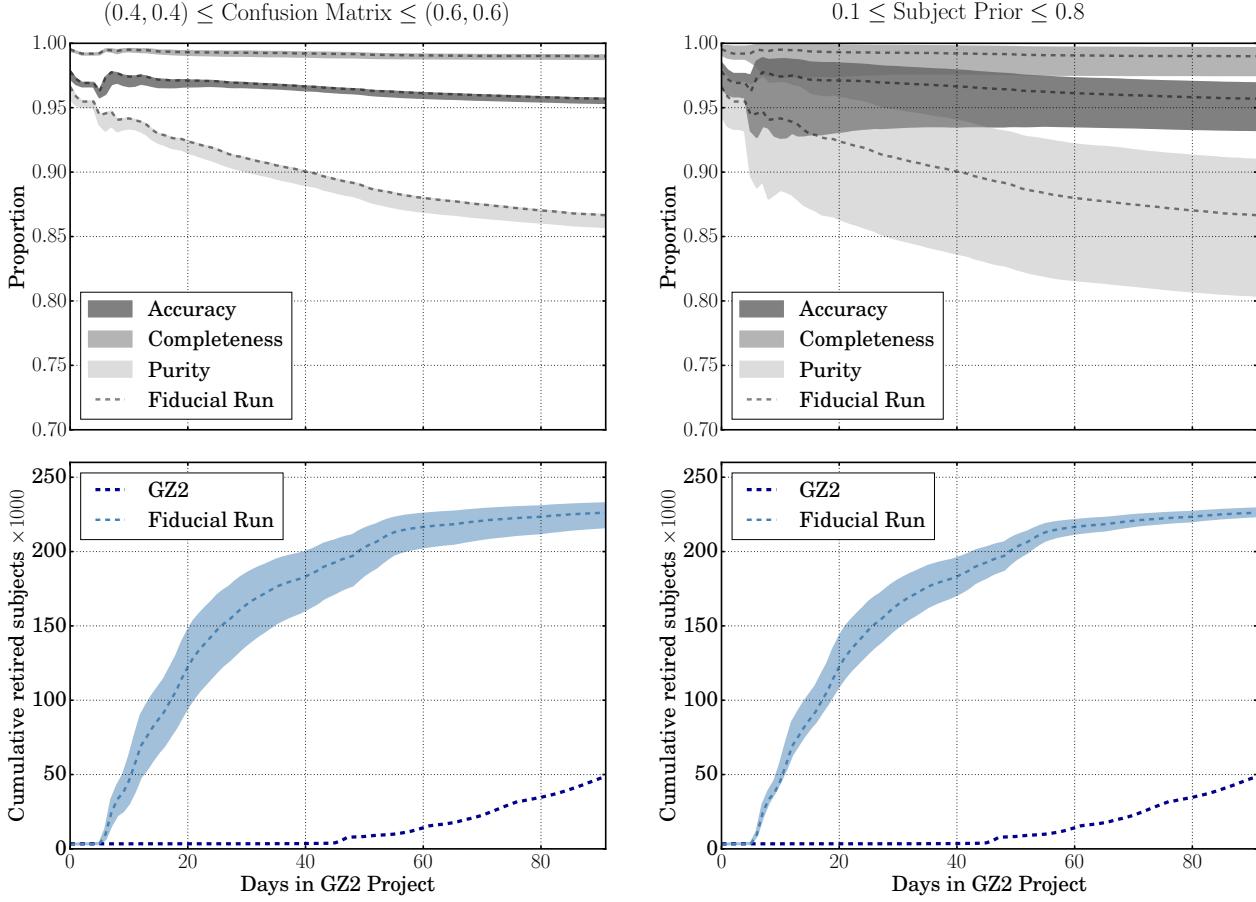


Figure A1. SWAP performance does not dramatically change even with a range of input parameters as compared to the fiducial run of Section 4.2 (dashed lines). *Left.* The quality (top) and retirement rate (bottom) when the confusion matrix is initialized as $(0.4, 0.4)$ and $(0.6, 0.6)$, with all other input parameters remaining constant. *Right.* Same as the left panel but allowing the subject prior probability, $p_0 = 0.2, 0.35$ and 0.8 . Changes in the confusion matrix have little impact on the quality of the labels but varies the total number of subjects retired. In contrast, changing the subject prior is more likely to affect the classification quality rather than the total number of subjects retired.

ever, in the Bayesian framework of SWAP, an agent with confusion matrix $(0.4, 0.4)$ contributes as much information as an agent with confusion matrix $(0.6, 0.6)$. The quality metrics computed are thus within a per cent of each other. In either case, we find that initializing agents at $(0.5, 0.5)$ provides optimal performance for the ‘training’ we simulate with our current approach. Further assessment would require a live project with real-time training and feedback.

Subject prior probability, p_0 . The prior probability assigned to each subject is an educated guess of the frequency of that characteristic in the scope of the data at hand. For galaxy morphologies, this number should be an estimate of the probability of observing a desired feature (bar, disk, ring, etc.). In our case, we desire simply to find galaxies that are ‘Featured’; however, this is dependent on mass, redshift, physical size, etc. The original GZ2 sample was selected primarily on magnitude and redshift. As there was no cut on galaxy size (with the exception that each galaxy be larger than the SDSS PSF), the sample includes a large range of masses and sizes. Designating a single prior is not clear-cut; we thus explore how various p_0 values effect the SWAP outcome.

We run simulations allowing p_0 to take values $0.2, 0.35$, and 0.8 and compare these to the fiducial run, with everything else remaining constant. The results are shown in the right panels of Figure A1. We again find that SWAP is consistent in terms of subject retirement which varies by only 1%. However, as can be seen in the top panel, the variation in our quality metrics is more pronounced. Firstly, though we retire nearly the same number of subjects over the course of each simulation, they are less consistent than our previous runs. That is, only 95% of retired subjects are common to all simulations. Secondly, of those that are common, only 94% receive the same label from SWAP indicating that changing the prior is more likely to produce a different label for a given subject than changing the initial agent confusion matrix. Finally, there is also a larger spread for the day on which a subject is retired as compared to the

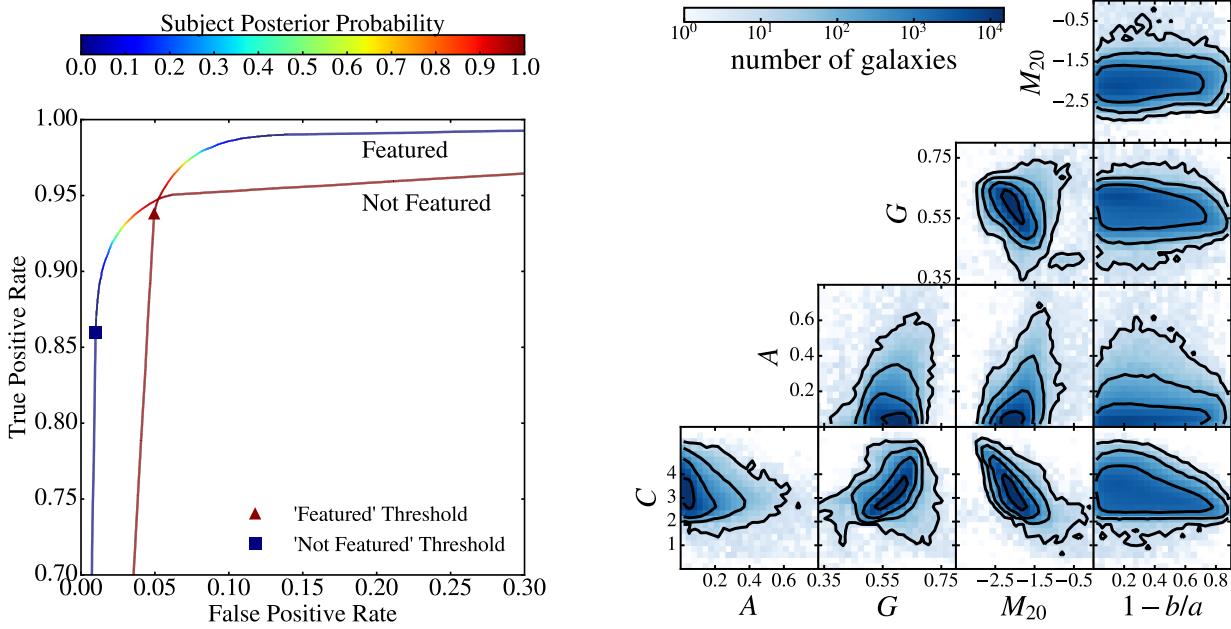


Figure A2. *Left.* Identifying ‘Featured’ subjects is independent of identifying ‘Not’ subjects. Both ROC curves use all subjects processed by SWAP where the score used to create the ROC curve is simply each subject’s achieved posterior probability. The Featured curve demonstrates how well we identify ‘Featured’ subjects with a threshold of 0.99, while the Not Featured curve demonstrates how well we identify ‘Not’ subjects with a threshold of 0.004. Typically, best performance is achieved by the score associated with the upper-left-most part of the curve. Our ‘Featured’ threshold is nearly optimal, while our ‘Not’ threshold could be improved since the blue square is not as close to the upper left hand corner as other possible values of the subject posterior. *Right.* Relation between measured morphology diagnostics for more than 280K SDSS galaxies. Most of these galaxies are processed through SWAP, receiving a posterior probability that estimates how likely each is to be ‘Featured’ or ‘Not’.

1015 fiducial run. These trends all contribute to a broader spread in accuracy, completeness, and purity as a function of
 1016 project time. We stress, however, that although more substantial than the previous comparison, these variations are
 1017 all within $\pm 5\%$.

1018 We can understand these variations more intuitively by considering the following. Recall that our retirement thresh-
 1019 olds, t_F and t_N , have not changed in these simulations. When p_0 is small, the subject’s probability is already closer
 1020 to t_N in probability space, and thus more subjects are classified as ‘Not’ compared to the fiducial run. Similarly,
 1021 when p_0 is large, some of these same subjects can instead be classified as ‘Featured’ because p_0 is already closer to t_F .
 1022 Obviously, both outcomes cannot be correct. We find that the simulation with $p_0 = 0.8$ performs the worst of any
 1023 run; this is a direct reflection of the fact that this prior is not suitable for this question or this dataset. Indeed,
 1024 the best performance is achieved when $p_0 = 0.35$. This reflects the distribution of ‘Featured’ subjects as determined
 1025 by GZ2_{raw} labels and is more characteristic of the expected proportion of ‘Featured’ galaxies in the local universe. As
 1026 a value far from the correct value can have a significant impact on the classification quality, it is important to choose
 1027 a prior wisely.

1028 **Retirement thresholds, t_F and t_N .** Retirement thresholds are directly related to the time that a subject will
 1029 spend in SWAP before retirement. If we lower t_F (and/or raise t_N), more subjects will be retired compared to the
 1030 fiducial run as each subject will have a smaller swath of probability space in which to fluctuate before crossing one of
 1031 these thresholds. On the other hand, if we raise t_F (and/or lower t_N), it will take longer for subjects to cross one of
 1032 these thresholds. This also increases the likelihood of some subjects never crossing either threshold, instead oscillating
 1033 indefinitely through probability space.

1034 What thresholds should one choose? To answer this question, we consider the left panel of Figure A2, which depicts
 1035 the receiver operating characteristic (ROC) curve for our fiducial simulation, an illustration of performance as a
 1036 function of a threshold for a binary classifier. ROC curves display the true positive rate against the false positive rate
 1037 for a discriminatory threshold or score with a perfect classifier achieving 100% true positives and no false positives.
 1038 The value of the threshold optimal for predicting class labels would be that which allows the ROC curve to reach the

upper-left-most point in the diagram. We have two thresholds to consider and thus we plot the curve twice: once under the assumption that “true positives” denote correctly identified ‘Featured’ subjects; and again under the assumption that “true positives” instead denote correctly identified ‘Not’ subjects. In both cases, the colour of the line corresponds to the subject posterior probability. We mark the location of $t_F = 0.99$ and $t_N = 0.004$ from our fiducial run with a red triangle and blue square respectively. We see that t_F is nearly optimal but t_N could be improved upon.

1044 B. MEASURING NONPARAMETRIC MORPHOLOGICAL DIAGNOSTICS ON SDSS STAMPS

1045 In order to train our Random Forest machine learning algorithm, we measure non-parametric morphology diagnostics
1046 for the GZ2 galaxy sample.

1047 We obtain i -band imaging from SDSS Data Release 12. Postage stamps are made from the SDSS fields for each
1048 galaxy with dimensions of 3 Petrosian radii. Galaxies located within 3 Petrosian radii of the edge of a field were
1049 excluded. Postage stamps undergo a cleaning process whereby nearby sources are identified with SExtractor (ver.
1050 2.8.6; [Bertin & Arnouts 1996](#)) and their pixels replaced with values that mimic the background in that region. We
1051 compute the following widely adopted nonparametric measurements of the galaxy light distribution on the cleaned
1052 postage stamps:

1053 Concentration is computed as $C = 5 \log(r_{80}/r_{20})$ where r_{80} and r_{20} are the radii containing 80% and 20% of the
1054 galaxy light respectively. Small values of this ratio tend to indicate disky galaxies, while larger values correlate with
1055 early-type ellipticals.

Asymmetry quantifies the degree of rotational symmetry in the galaxy light distribution (not necessarily the physical
shape of the galaxy as this parameter is not highly sensitive to low surface brightness features). A correction for
background noise is applied (as in e.g. [Conselice et al. \(2000\)](#)), i.e.,

$$A = \frac{\sum_{x,y} |I - I_{180}|}{2 \sum |I|} - B_{180} \quad (\text{B1})$$

1056 where I is the galaxy flux in each pixel (x, y) , I_{180} is the image rotated by 180 degrees about the galaxy’s central pixel,
1057 and B_{180} is the average asymmetry of the background.

The Gini coefficient, G , ([Glasser 1962](#); [Abraham et al. 2003](#)) describes how uniformly distributed a galaxy’s flux is.
If G is 0, the flux is distributed homogeneously among all galaxy pixels.; if G is 1, the light is contained within a single
pixel. This term correlates with C , however, G does not require that the flux be in the central region of the galaxy.
We follow [Lotz et al. \(2004\)](#) by first ordering the pixels by increasing flux value, and then computing

$$G = \frac{1}{|\bar{X}|n(n-1)} \sum_i^n (2i - n - 1)|X_i| \quad (\text{B2})$$

1058 where n is the number of pixels assigned to the galaxy, and \bar{X} is the mean pixel value.

M_{20} ([Lotz et al. 2004](#)) is the second order moment of the brightest 20% of the galaxy flux. We compute it as

$$M_{tot} = \sum_i^n f_i[(x_i - x_c)^2 + (y_i - y_c)^2] \quad (\text{B3})$$

$$M_{20} = \log_{10}\left(\frac{\sum_i M_i}{M_{tot}}\right), \quad \text{while } \sum_i f_i < 0.2f_{tot} \quad (\text{B4})$$

1059 where M_{tot} , the total moment, is computed first and f_{tot} is the total flux. For centrally concentrated objects, M_{20}
1060 correlates with C but is also sensitive to bright off-centre knots of light.

1061 Finally, we use the ellipticity, $\epsilon = 1 - b/a$, of the light distribution as measured by SExtractor which computes the
1062 semi-major axis a and semi-minor axis b from the second-order moments of the galaxy light.

1063 In total, we measure morphological indicators for 282,350 SDSS galaxies. The relations between these diagnostics for
1064 the full sample is shown in the right panel of Figure A2. The code developed to clean and compute these morphology
1065 indicators is open source and can be found at https://github.com/melaniebeck/measure_morphology.

REFERENCES

1066 Abraham, R. G., Tanvir, N. R., Santiago, B. X., et al. 1996, [Abraham, R. G., Valdes, F., Yee, H. K. C., & van den Bergh, S.](#)

1067 MNRAS, 279, L47

1069 1994, ApJ, 432, 75

- 1070 Abraham, R. G., van den Bergh, S., & Nair, P. 2003, ApJ, 588,
1071 218
- 1072 Astropy Collaboration, Robitaille, T. P., Tollerud, E. J., et al.
1073 2013, A&A, 558, A33
- 1074 Baillard, A., Bertin, E., de Lapparent, V., et al. 2011, A&A, 532,
1075 A74
- 1076 Ball, N. M., Loveday, J., Fukugita, M., et al. 2004, MNRAS, 348,
1077 1038
- 1078 Bamford, S. P., Nichol, R. C., Baldry, I. K., et al. 2009,
1079 MNRAS, 393, 1324
- 1080 Banerji, M., Lahav, O., Lintott, C. J., et al. 2010, MNRAS, 406,
1081 342
- 1082 Bershadsky, M. A., Jangren, A., & Conselice, C. J. 2000, AJ, 119,
1083 2645
- 1084 Bertin, E., & Arnouts, S. 1996, A&AS, 117, 393
- 1085 Blanton, M. R., Hogg, D. W., Bahcall, N. A., et al. 2003, ApJ,
1086 594, 186
- 1087 Breiman, L. 2001, Machine Learning, 45, 5
- 1088 Cardamone, C., Schawinski, K., Sarzi, M., et al. 2009, MNRAS,
1089 399, 1191
- 1090 Casteels, K. R. V., Conselice, C. J., Bamford, S. P., et al. 2014,
1091 MNRAS, 445, 1157
- 1092 Conselice, C. J. 2003, ApJS, 147, 1
1093 —. 2006, MNRAS, 373, 1389
- 1094 Conselice, C. J., Bershadsky, M. A., & Jangren, A. 2000, ApJ, 529,
1095 886
- 1096 Darg, D. W., Kaviraj, S., Lintott, C. J., et al. 2010, MNRAS,
1097 401, 1552
- 1098 de Vaucouleurs, G. 1959, Handbuch der Physik, 53, 275
- 1099 Dieleman, S., Willett, K. W., & Dambre, J. 2015, MNRAS, 450,
1100 1441
- 1101 Dressler, A. 1980, ApJ, 236, 351
- 1102 Elmegreen, B. G., Bournaud, F., & Elmegreen, D. M. 2008, ApJ,
1103 688, 67
- 1104 Elmegreen, B. G., Elmegreen, D. M., Sánchez Almeida, J., et al.
1105 2013, ApJ, 774, 86
- 1106 Freeman, P. E., Izbicki, R., Lee, A. B., et al. 2013, MNRAS, 434,
1107 282
- 1108 Galloway, M. A., Willett, K. W., Fortson, L. F., et al. 2015,
1109 MNRAS, 448, 3442
- 1110 Glasser, G. J. 1962, Journal of the American Statistical
1111 Association, 57, 648
- 1112 Griffith, R. L., Cooper, M. C., Newman, J. A., et al. 2012, ApJS,
1113 200, 9
- 1114 Holwerda, B. W., Muñoz-Mateos, J.-C., Comerón, S., et al. 2014,
1115 ApJ, 781, 12
- 1116 Hubble, E. P. 1936, The Realm of the Nebulae (Yale University
1117 Press)
- 1118 Huertas-Company, M., Rouan, D., Tasca, L., Soucail, G., & Le
1119 Fèvre, O. 2008, A&A, 478, 971
- 1120 Huertas-Company, M., Gravet, R., Cabrera-Vives, G., et al.
1121 2015, ApJS, 221, 8
- 1122 Kartaltepe, J. S., Mozena, M., Kocevski, D., et al. 2015, ApJS,
1123 221, 11
- 1124 Kauffmann, G., Heckman, T. M., White, S. D. M., et al. 2003,
1125 MNRAS, 341, 54
- 1126 Kormendy, J. 1977, ApJ, 217, 406
- 1127 Kormendy, J., & Kennicutt, Jr., R. C. 2004, ARA&A, 42, 603
- 1128 Land, K., Slosar, A., Lintott, C., et al. 2008, MNRAS, 388, 1686
- 1129 Laureijs, R., Amiaux, J., Arduini, S., et al. 2011, ArXiv e-prints,
1130 arXiv:1110.3193
- 1131 Lintott, C., Schawinski, K., Bamford, S., et al. 2011, MNRAS,
1132 410, 166
- 1133 Lintott, C. J., Schawinski, K., Slosar, A., et al. 2008, MNRAS,
1134 389, 1179
- 1135 Lotz, J. M., Primack, J., & Madau, P. 2004, AJ, 128, 163
- 1136 LSST Science Collaboration, Abell, P. A., Allison, J., et al. 2009,
1137 ArXiv e-prints, arXiv:0912.0201
- 1138 Marshall, P. J., Verma, A., More, A., et al. 2016, MNRAS, 455,
1139 1171
- 1140 Masters, K. L., Nichol, R. C., Hoyle, B., et al. 2011, MNRAS,
1141 411, 2026
- 1142 Meert, A., Vikram, V., & Bernardi, M. 2016, MNRAS, 455, 2440
- 1143 More, A., Verma, A., Marshall, P. J., et al. 2016, MNRAS, 455,
1144 1191
- 1145 Nair, P. B., & Abraham, R. G. 2010, ApJS, 186, 427
- 1146 Nakamura, O., Fukugita, M., Yasuda, N., et al. 2003, AJ, 125,
1147 1682
- 1148 Odewahn, S. C., Cohen, S. H., Windhorst, R. A., & Philip, N. S.
1149 2002, ApJ, 568, 539
- 1150 Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, Journal
1151 of Machine Learning Research, 12, 2825
- 1152 Peng, C. Y., Ho, L. C., Impey, C. D., & Rix, H.-W. 2002, AJ,
1153 124, 266
- 1154 Peng, Y.-j., Lilly, S. J., Kováč, K., et al. 2010, ApJ, 721, 193
- 1155 Peth, M. A., Lotz, J. M., Freeman, P. E., et al. 2016, MNRAS,
1156 458, 963
- 1157 Sandage, A. 1961, The Hubble atlas of galaxies
- 1158 Scarlata, C., Carollo, C. M., Lilly, S., et al. 2007, ApJS, 172, 406
- 1159 Schawinski, K., Urry, C. M., Simmons, B. D., et al. 2014,
1160 MNRAS, 440, 889
- 1161 Sersic, J. L. 1968, Atlas de galaxias australes
- 1162 Shen, S., Mo, H. J., White, S. D. M., et al. 2003, MNRAS, 343,
1163 978
- 1164 Sheth, K., Elmegreen, D. M., Elmegreen, B. G., et al. 2008, ApJ,
1165 675, 1141
- 1166 Simard, L., Mendel, J. T., Patton, D. R., Ellison, S. L., &
1167 McConnachie, A. W. 2011, ApJS, 196, 11
- 1168 Simmons, B. D., Melvin, T., Lintott, C., et al. 2014, MNRAS,
1169 445, 3466
- 1170 Simmons, B. D., Lintott, C., Willett, K. W., et al. 2017,
1171 MNRAS, 464, 4420
- 1172 Smethurst, R. J., Lintott, C. J., Simmons, B. D., et al. 2016,
1173 MNRAS, 463, 2986
- 1174 Snyder, G. F., Torrey, P., Lotz, J. M., et al. 2015, MNRAS, 454,
1175 1886
- 1176 Strateva, I., Ivezić, Ž., Knapp, G. R., et al. 2001, AJ, 122, 1861
- 1177 Taylor, M. B. 2005, in Astronomical Society of the Pacific
1178 Conference Series, Vol. 347, Astronomical Data Analysis
1179 Software and Systems XIV, ed. P. Shopbell, M. Britton, &
1180 R. Ebert, 29
- 1181 van den Bergh, S. 1976, ApJ, 206, 883
- 1182 Watanabe, M., Kodaira, K., & Okamura, S. 1985, ApJ, 292, 72
- 1183 Whitmore, B. C., Lucas, R. A., McElroy, D. B., et al. 1990, AJ,
1184 100, 1489
- 1185 Willett, K. W., Lintott, C. J., Bamford, S. P., et al. 2013,
1186 MNRAS, 435, 2835
- 1187 Willett, K. W., Galloway, M. A., Bamford, S. P., et al. 2017,
1188 MNRAS, 464, 4176

Replaced: ~~, reprocessed~~ replaced with: . We reprocess the top-level question of the GZ2 decision tree, on page ??, line ??.

Added: Though we simplify our analysis to examine a

List of Changes

binary classification scheme, and as, on page ??, line ??.

Added: As a proof of concept, we focus on the first question of the Galaxy Zoo decision tree in this paper., on page 2, line 97.

Replaced: Our replaced with: We demonstrate that our, on page 2, line 99.

Added: , well-suited for the top-level question of the GZ2 decision tree, as discussed below, on page 2, line 132.

Added: ⁷, on page 3, line 171.

Added: choice, on page 3, line 180.

Added: By combining “star or artifact” responses with “features or disk” responses, we obtain a binary task., on page 3, line 185.

Replaced: ~~By combining the “star or artifact” vote fraction, f_{artifact} , with the “features or disk” vote fraction, f_{features} we obtain a binary response. Here, a vote fraction is simply the fraction of volunteers who voted for a particular response. We define a label for each GZ2 subject as the majority vote fraction; that is, if $f_{\text{features}} + f_{\text{artifact}} > f_{\text{smooth}}$, the galaxy is labelled ‘Featured’, otherwise it is labelled ‘Not’.~~ replaced with: In order to compare our classification output with GZ2 we assign each subject a descriptive label. GZ2 classifications are comprised of volunteer vote fractions (f_{response}) for each response to every task in the decision tree, where vote fractions are derived from the fraction of volunteers who voted for a particular response (more on this below). GZ2 classifications are thus continuous. A common technique is to place a threshold on these vote fractions to select samples with an emphasis on purity or completeness, depending on the science case. For our current analysis we choose a threshold of 0.5, that is, if $f_{\text{features}} + f_{\text{artifact}} > f_{\text{smooth}}$, the galaxy is labelled ‘Featured’, otherwise it is labelled ‘Not’. We note that this threshold is not much different from the suggested 0.430 threshold in Willett et al. (2013) that produces a well-sampled subset of ‘Featured’ galaxies. Though naive, we will demonstrate throughout this paper that this threshold produces adequate results, though a more sophisticated mechanism will be explored in a future publication., on page 3, line 187.

Added: by combining the “star or artifact” with “features or disk” responses, on page 4, line 217.

Replaced: ~~assigns every subject three types of volunteer vote fractions:~~ replaced with: publishes three types of

vote fractions for each subject:, on page 4, line 219.

Deleted: ~~, a task we perform as well.~~, on page 4, line 225.

Replaced: ~~However, because our mechanism is entirely different from GZ2,~~ replaced with: The SWAP algorithm (described below) also has a mechanism to weight volunteer votes, however, the two methods are in stark contrast. In order to ensure equal footing,, on page 4, line 226.

Added: that have received no post-processing whatsoever., on page 4, line 232.

Deleted: ~~actually~~, on page 4, line 259.

Added: (denoted N), on page 4, line 261.

Added: where X is the true classification of the subject and “ X ” is the classification made by the volunteer upon viewing the subject. Thus, on page 4, line 264.

Added: Therefore the confusion matrix for a single volunteer goes as, on page 4, line 270.

Added: where probabilities are normalised such that $P("A" | A) = 1 - P("N" | A)$., on page 4, line 273.

Replaced: ~~derive~~ replaced with: compute, on page 4, line 277.

Replaced: ~~each subject’s probability is continually updated~~ replaced with: each subject’s posterior probability is updated after every volunteer classification, on page 4, line 281.

Replaced: ~~Probability thresholds can be set such that subjects crossing a threshold are highly likely to exhibit the feature of interest or the absence thereof. These subjects are then considered retired.~~ replaced with: Upper and lower probability thresholds can be set such that when a subject’s posterior crosses the upper threshold it is highly likely to exhibit feature A , while if it crosses the lower threshold it is highly likely that feature A is absent. Subjects whose posteriors cross either of these thresholds are considered retired., on page 4, line 285.

Replaced: ~~they are adept at correctly identifying both ‘Featured’ and ‘Not’ subjects.~~ replaced with: they correctly identify both ‘Featured’ and ‘Not’ subjects more than 50% of the time., on page ??, line ??.

Replaced: ~~they are adept at correctly identifying both ‘Featured’ and ‘Not’ subjects.~~ replaced with: they correctly identify both ‘Featured’ and ‘Not’ subjects more than 50% of the time., on page 5.

Replaced: ~~We describe how we engineer the GZ2 data to mimic the Space Warps system.~~ replaced with: Though we cannot retroactively train GZ2 volunteers in such a manner, we develop a gold standard sample in lieu of simulated data, and arrange the classification order of gold standard data in order to mimic the Space Warps

⁷ A visualization of this decision tree can be found at https://data.galaxyzoo.org/gz_trees/gz_trees.html

system., on page 5, line 306.

Replaced: ~~Expert classifications were obtained~~ replaced with: We must generate expert labels for these galaxies that are consistent with the labels we defined for GZ2 classifications. These are obtained, on page 5, line 319.

Added: This ensures that our expert labels are defined in exactly the same manner as the labels for all SDSS galaxies in the GZ2 sample., on page 6, line 328.

Replaced: ~~are generally good at correctly identifying both ‘Featured’ and ‘Not’ subjects.~~ replaced with: correctly identify both ‘Featured’ and ‘Not’ subjects more than 50% of the time., on page 6, line 362.

Added: Figure 3 (adapted from Figure 5 of Marshall et al. 2016) demonstrates how subject posterior probabilities are updated with each classification. The arrow in the top panel denotes the prior probability, $p_0 = 0.5$. With each classification, that prior is updated into a posterior probability thus creating a trajectory through probability space for each subject. The yellow and orange lines show the trajectories of a random sample of ‘Featured’ and ‘Not’ subjects from our gold standard sample, while the black lines show the trajectories of a random sample of GZ2 subjects that were not part of the gold standard sample. The vertical yellow and orange dashed lines correspond to the retirement thresholds, t_F and t_N . The lower panel shows the full distribution of GZ2 subject posteriors at the end of our simulation. An overwhelming majority of subjects cross one of these retirement thresholds., on page 7, line 372.

Replaced: ~~The original GZ2 project took approximately one year to classify as many subjects, representing a factor of four increase in the classification rate.~~ replaced with: Here we are considering simply the number of classifications logged each day and not the length of time spent on a single classification. Under the assumption that collapsing the GZ2 decision tree to a single question would not decrease the number of classifications collected each day during the GZ2 project, processing volunteer classifications through SWAP presents a promising increase in classification efficiency., on page 7, line 421.

Added: We discuss these discrepancies in the next section., on page 7, line 436.

Added: We explore these subjects in redshift, magnitude, physical size, and concentration. We find no correlation with any of these variables, suggesting there are

no physical reasons why SWAP’s label disagrees with GZ2_{raw}., on page 7, line 470.

Replaced: ~~We find the majority of these disagreements are due to uncertainties in the GZ2_{raw} label.~~ replaced with: Instead we consider the errors associated with the GZ2 vote fraction, which can be estimated as binomial. Let success be a response of “smooth” and failure be any other response. The 68% confidence interval on a subject with $f_{\text{smooth}} = 0.5$ is then (0.42, 0.57) assuming 40 classifications, each with a probability of 0.5., on page 7, line 475.

Replaced: ~~incorrectly labelled~~ replaced with: disagreements between SWAP and GZ2 are for, on page 8, line 489.

Replaced: ~~indicating that the GZ2 raw vote fractions are simply too uncertain to provide high quality labels.~~ replaced with: It is thus unsurprising that SWAP and GZ2 disagree most within the errors of GZ2 vote fractions., on page 8, line 492.

Added: compared to expert classifications, on page 8, line 509.

Replaced: ~~The training score is the accuracy of the trained machine applied to its own training sample. The cross-validation score is the accuracy of the machine computed during the cross-validation process.~~ replaced with: These curves show the mean accuracy computed during cross-validation and on the training sample, where the shaded regions denote the standard deviation., on page ??, line ??.

Replaced: ~~The training score is the accuracy of the trained machine applied to its own training sample. The cross-validation score is the accuracy of the machine computed during the cross-validation process.~~ replaced with: These curves show the mean accuracy computed during cross-validation and on the training sample, where the shaded regions denote the standard deviation., on page 9.

Replaced: ~~GZX thus provides an order of magnitude increase in the rate of classification over the traditional crowd-sourced approach.~~ replaced with: Though our analysis considers the only the top-level task of GZ2’s decision tree, GZX suggests a tantalizing potential to increase the classification rate by an order of magnitude over the traditional crowd-sourced approach., on page 10, line 722.