

An assessment of photometric redshift PDF performance in the context of LSST

S.J. Schmidt¹, A.I. Malz^{2,3}, J.Y.H. Soo⁴, M. Brescia⁵, S. Cavaoti^{5,6}, G. Longo⁶, I.A. Almosallam^{7,8}, M.L. Graham⁹, A.J. Connolly⁹, E. Nourbakhsh¹, J. Cohen-Tanugi¹⁰, H. Tranin¹⁰, P.E. Freeman¹¹, K. Iyer¹², J.B. Kalmbach¹³, E. Kovacs¹⁴, A.B. Lee¹¹, C. Morrison⁹, J. Newman¹⁵, E. Nuss¹⁰, T. Pospisil¹¹, M.J. Jarvis^{16,17}, R. Izbicki^{18,19}

(LSST Dark Energy Science Collaboration)

¹ Department of Physics, University of California, One Shields Ave., Davis, CA, 95616, USA

² Center for Cosmology and Particle Physics, New York University, 726 Broadway, New York, 10003, USA

³ Department of Physics, New York University, 726 Broadway, New York, 10003, USA

⁴ Department of Physics and Astronomy, University College London, Gower Street, London WC1E 6BT, UK

⁵ INAF-Capodimonte Observatory, Salita Moiariello 16, I-80131, Napoli, Italy

⁶ Department of Physics E. Pancini, University Federico II, via Cinthia 6, I-80126, Napoli, Italy

⁷ King Abdulaziz City for Science and Technology, Riyadh 11442, Saudi Arabia

⁸ Information Engineering, Parks Road, Oxford, OX1 3PJ, UK

⁹ Department of Astronomy, University of Washington, Box 351580, U.W., Seattle WA 98195, USA

¹⁰ Laboratoire Univers et Particules de Montpellier, Université de Montpellier, CNRS, Montpellier, France

¹¹ Department of Statistics & Data Science, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, USA

¹² Department of Physics and Astronomy, Rutgers, The State University of New Jersey, 136 Frelinghuysen Road, Piscataway, NJ 08854-8019 USA

¹³ Department of Physics, University of Washington, Box 351560, Seattle, WA 98195, USA

¹⁴ Argonne National Laboratory, Lemont, IL 60439, USA

¹⁵ Department of Physics and Astronomy and PITT PACC, University of Pittsburgh, 3941 O'Hara St., Pittsburgh, PA 15260, USA

¹⁶ Astrophysics, Department of Physics, University of Oxford, Denys Wilkinson Building, Keble Road, Oxford, OX1 3RH, UK

¹⁷ Department of Physics and Astronomy, University of the Western Cape, Bellville 7535, South Africa

¹⁸ Department of Statistics, Federal University of Sao Carlos, Sao Carlos, Brazil

¹⁹ External collaborator

30 October 2018

ABSTRACT

Photometric redshift (photo- z) probability distribution functions (PDFs) are a key data product of nearly all upcoming galaxy imaging surveys. However, the photo- z PDFs resulting from different techniques are not in general consistent with one another, and an optimal method for obtaining an accurate PDF remains unclear. We present the results of an initial study of the Large Synoptic Survey Telescope Dark Energy Science Collaboration (LSST-DESC), the first in a planned series of papers testing multiple photo- z codes on simulations of upcoming LSST galaxy photometry catalogues. This initial test evaluates photo- z algorithms in the presence of representative training data and in the absence of several common sources of systematic errors that affect the procedures by which photo- z PDFs are derived. The photo- z PDFs are evaluated using multiple metrics including the Kolmogorov-Smirnov statistic, Cramer-von Mises statistic, Anderson-Darling statistic, Kullback-Leibler divergence, $N(z)$ moments, quantile-quantile plots and probability integral transform. We observe several trends, including an overall over/under-prediction in the broadness of the PDFs for several of the codes. A careful accounting of all systematics discovered will be necessary for the codes employed in upcoming analyses in order to achieve unbiased cosmological measurements.

Key words: galaxies: distances and redshifts – galaxies: statistics – methods: statistical

2 LSST Dark Energy Science Collaboration

1 INTRODUCTION

Large-scale photometric galaxy surveys are entering a new era with currently or soon-to-be running Stage III and Stage IV dark energy experiments like the Dark Energy Survey (DES, Abbott et al. 2005), the Kilo-Degree Survey (KiDS, de Jong et al. 2013), Large Synoptic Survey Telescope (LSST, Abell et al. 2009), Euclid (Laureijs et al. 2011), and Wide-Field Infrared Survey Telescope (WFIRST, Green et al. 2012). The move to imaging based surveys, rather than spectroscopic based, for cosmological measurements makes proper understanding of photometric redshifts (“photo- z ’s”) of paramount importance, as cosmological distance measures for statistical samples are directly dependent on photo- z measurements.

The unprecedented sample size of LSST galaxies, expected to number several billion for the main cosmological sample, necessitates stringent constraints on photo- z accuracy if systematic errors are not to dominate the statistical errors. The LSST Science Requirements Document (SRD)¹ lists the photometric redshift goals for a magnitude limited sample with $i < 25$ as: root-mean-square error with a goal of $\sigma_z < 0.02(1+z)$; 3σ “catastrophic outlier” rate below 10%; bias below 0.003².

The tremendous size of LSST’s galaxy catalog will be enabled by its exceptional depth, pushing to fainter magnitudes and deeper imaging and including galaxies of lower luminosity and higher redshift than ever before. In addition to the contribution of low signal-to-noise photometry to the systematic error of photo- z s, these populations introduce major physical degeneracies, for example the Lyman break/Balmer break degeneracy, that were not present in the populations covered in previous large area surveys like the Sloan Digital Sky Survey (SDSS, York et al. 2000) and the Two Micron All Sky Survey (2MASS, Skrutskie et al. 2006). In order to meet LSST’s demanding error budget, it will be necessary to fully characterize those degeneracies wherein multiple redshift solutions have comparable likelihood.

There is often a desire to have a single valued “point-estimate” redshift for an individual galaxy. However, the complex, non-linear (and often non-unique) nature of the mapping between broad band fluxes and redshift means that a single value is unable to capture the full redshift information encoded in a galaxy’s magnitudes. For example, a common point-estimate for a template-based method is taking the highest likelihood solution as the point photo- z . A single valued redshift ignores degenerate redshift solutions of lower probability, potentially biasing photometric redshift estimates both for individual galaxies and ensemble distributions. Storing more information is necessary, most often photo- z codes output the redshift probability density function (PDF), also often referred to as $p(z)$, describing the relative likelihood as a function of redshift. Early template methods such as Fernández-Soto et al. (1999) converted rel-

ative χ^2 values of template spectra to likelihoods to estimate $p(z)$. Soon after, codes such as Benítez (2000) added a Bayesian prior and output a posterior probability distribution. While many early machine learning based algorithms focused on a point-estimate, Firth et al. (2003) used a neural net with 1000 realizations scattered within the photometric errors to estimate a $p(z)$. As more groups began to employ photometric redshifts in their cosmological analyses, realization that point-estimate photo- z ’s were inadequate for precision cosmology measurements (Mandelbaum et al. 2008). From around this point onward, most photo- z algorithms have attempted to implement some estimate of the overall redshift probability in their outputs, and some surveys began supplying a full $p(z)$ rather than a simple redshift point-estimate and error (e. g. de Jong et al. 2017).

There are numerous techniques for deriving photo- z PDFs from photometry, yet no one method has yet been established as clearly superior. Quantitative comparisons of photo- z methods have been made before. The Photo- z Accuracy And Testing (PHAT, Hildebrandt et al. 2010) effort focused on point estimates derived from many photometric bands. DES compared several codes for point estimates (Sánchez et al. 2014) and a summary statistic of photo- z interim posteriors for tomographically binned galaxy subsamples (Bonnett et al. 2016). This paper is distinguished by its inclusion of metrics of photo- z interim posteriors themselves and consideration of both classic and state-of-the-art photo- z algorithms, comparing the performance of several of the most widely employed codes as well as some that have been developed only recently on the basis of metrics appropriate for a probabilistic data product. The results presented in this work are a major focus of the Photometric Redshift working group of the LSST Dark Energy Science Collaboration (LSST-DESC). This work is laid out in the Science Roadmap (SRM)³ as one of the critical activities to be completed in preparation for dark energy science analysis on the first year LSST data. This is the first of multiple papers by the working group, which will grow in sophistication. In this initial paper we focus on evaluating the performance of photometric redshift codes and their ability to produce accurate PDFs in the presence of representative training sets. This can be thought of as an initial test under near perfect conditions, before further complexities are added in future papers. Comparing the relative performance of the codes enables us to evaluate whether each code is using information in an optimal way, and may reveal enhancements in some codes and deficiencies in others, either in the fundamental algorithm, or in specific implementation.

Certain science cases need redshift information on individual objects, e. g. identification of host galaxy redshift for supernova classification, or identifying potential cluster membership. Other science cases need only ensemble redshift information; for instance current weak lensing techniques require the overall redshift distribution $N(z)$ for tomographic redshift samples, but do not need single galaxy estimates. In the case of the multiple types of probes of cosmology enabled by the LSST cosmology sample of several billion galaxies, the number of redshift bins and their photo-

¹ available at <https://docushare.lsstcorp.org/docushare/dsweb/Get/LPM-17>

² Note that at the time the SRD was written, these goals were stated in terms of a photo- z point estimate for each galaxy, as was standard in many previous studies, while in this paper we emphasize the importance of using a full photo- z PDF.

³ Available at: http://lsst-desc.org/sites/default/files/DESC_SRMs_V1_1.pdf

113 z requirements vary with the specific probe; 2-point angular correlations benefit from many bins, while weak lensing 114 probes do not (due to the wide lensing kernel). Large photometric surveys such as LSST must develop algorithms that 115 meet the needs of all such science cases. In order to meet these ambitious goals for photo- z accuracy, every aspect of 116 photo- z estimation will have to be optimized: the algorithms employed, both template and machine-learning based (both 117 in design and implementation); the spectroscopic data used as a training set for machine learning algorithms or to estimate template sets and train Bayesian priors; and probabilistic catalog compression schemes that balance information 118 retention against limited storage resources.

119 Before moving forward, we must address how the best methods may be unique to the performance metrics and science 120 cases considered and what distinguishes photo- z PDFs of different methods from one another. Though photo- z 121 PDFs are often written simply as $p(z)$, the PDF itself must be an interim posterior distribution $p(z|d, I)$, the probability 122 of redshift conditioned on photometric data d that has actually been observed and the prior information I that guides 123 how a redshift is extracted from the photometry. If we run multiple photo- z codes on a single dataset, the photo- z interim 124 posteriors will not be identical because each code is based on assumptions in the form of an interim prior — 125 these assumptions form the premise for photo- z estimation as a whole and are the only way to introduce differences in estimates of what would otherwise be a shared photo- z posterior $p(z|d)$ regardless of the code used to obtain it. Though 126 explicit knowledge of the interim prior is necessary to use photo- z interim posteriors self-consistently in physical inference, the interim prior of a particular methodology is often 127 implicit and not necessarily shared among all galaxies in the catalog.

128 This paper therefore aims (1) to constrain the impact of the interim prior I by separating it into a component I_H representing 129 the method itself and a component I_D representing physical 130 information, such as a training set or SED template library and (2) to present a procedure for evaluating the 131 performance of photo- z codes in generic tests that may include many more systematics in the interim prior I . In order to 132 isolate the effects encapsulated by I_H of variation between 133 codes from issues with the training set or template library 134 encapsulated by I_D , we use an identical set of simulated 135 galaxies for every code and construct a template library 136 and training sample that are *complete and representative* 137 and shared among all codes; that is, our training sample for 138 machine learning codes is drawn from the same underlying 139 galaxy population as our test set, with no additional selections, and the SED library used for template-based codes is 140 the same as the one used to generate the photometric data. 141 We explore a number of performance metrics in this paper, 142 not to make a conclusion regarding the superiority or even 143 relative favorability of each code but to establish a method 144 for comparing photo- z PDFs derived by different methods. 145 These test conditions set the stage for addressing in a future 146 paper the crucial issue of incomplete and non-representative 147 prior information.

148 The outline of the paper is as follows: in § 2 we present 149 the simulated data set; in § 3 we describe the current generation 150 codes employed in the paper; in § 4 we discuss the interpretation 151 of photo- z PDFs in terms of metrics of accuracy;

152 in § 5 we show our results and compare the performance of 153 the codes; in § 6 we offer our conclusions and discuss future 154 extensions of this work.

2 THE SIMULATION AND MOCK GALAXY CATALOG

155 In order to test the current generation codes, we employ an 156 existing simulated galaxy catalogue. The simulation is completely 157 catalogue-based, with no image construction or mock 158 measurements made. We describe these in detail below.

2.1 Buzzard-v1.0 simulation

159 The BUZZARD-HIGHRES-V1.0 (De Rose et al., in prep; Wechsler et al., in prep) catalogue construction started with a 160 dark matter only simulation. This N-body simulation contained 161 2048^3 particles in a 400 Mpc h^{-1} box. A set of time 162 snapshots (with smoothing and interpolation between snapshots) were saved in order to construct a lightcone. Dark 163 matter halos were identified using the ROCKSTAR software 164 package (Behroozi et al. 2013). These dark matter halos were 165 populated with galaxies with a stellar mass and absolute r -band 166 magnitude in the SDSS system determined using a sub-halo abundance matching model constrained to match 167 both projected two-point galaxy clustering statistics and an 168 observed conditional stellar mass function (Reddick et al. 169 2013).

170 To assign an SED to each galaxy, the *Adding Density 171 Dependent Spectral Energy Distributions* (ADDSEDS, deRose in prep.)⁴ procedure was used. This consisted of 172 training an empirical relation between absolute r -band magnitude, local galaxy density, and SED using a sample of 173 $\sim 5 \times 10^5$ galaxies from the magnitude-limited Sloan Digital 174 Sky Survey Data Release 6 Value Added Galaxy Catalog (175 Blanton et al. 2005). Each SDSS spectrum is fit with a sum 176 of five SED components using the K-CORRECT v4.3? software 177 package⁵ (Blanton & Roweis 2007), thus each galaxy 178 SED is parameterized as five weights for the basis SEDs. 179 The distance to the spatial projected fifth-nearest neighbour 180 was used as a proxy for local density in the SDSS training 181 sample. For each simulated galaxy, a galaxy with similar 182 absolute r -band magnitude and local galaxy density was chosen 183 from the training set, and that training galaxy's SED 184 was assigned to the simulated galaxy. This process is done 185 in such a way as to preserve the colour-density relation of 186 galaxy environment. Given the SED, absolute r -band magnitude and redshift, we computed apparent magnitudes in 187 the six LSST filter passbands, $ugrizy$. We assigned magnitude 188 errors in the six bands using the simple model described 189 in Ivezić et al. (2008), assuming full 10-year depth observations 190 had been completed. The number of total 30-second 191 visits assumed when generating the photometric errors differs 192 slightly from the fiducial numbers assumed for LSST: 193 we assume 60 visits in u-band, 80 visits in g-band, 180 visits 194 in r-band, 180 visits in i-band, 160 visits in z-band, and 160 195

196 ⁴ <https://github.com/vipasu/addseds>

197 ⁵ <http://kcorrect.org>

4 LSST Dark Energy Science Collaboration

visits in y-band. In the course of simulating Gaussian photometric errors, we add noise to objects fluxes, and some of these noisy fluxes will become negative in one or more bands. We call such negative fluxes “non-detections” and signify them with a placeholder magnitude of 99.0 in the catalog. Thus, further mentions of “non-detections” refer to objects that would be “looked at but not seen” in multi-band forced photometry, and the photo-z codes will treat them as such.

2.1.1 Selection of training and test sets

The total catalogue covered 400 square degrees and contained 238 million galaxies to an apparent magnitude limit of $r=29$ and spanning the redshift range $0 < z \leq 8.7$. In order for statistical errors not to dominate, we need less than one million galaxies in our sample. Several studies claim that only a few tens of thousands of spectra are necessary to calibrate photo-z surveys to Stage IV requirements (e.g. Bernstein & Huterer (2010), Masters et al. (2017)). Therefore, we aim for a final number of training galaxies between 3×10^4 and 5×10^4 in our sample. In order to reduce our sample to a reasonable size, we limit our dataset to a subset of ~ 16.8 square degrees selected from five separate spatial regions of the simulation. Systematic problems with galaxy colors above $z > 2$ were observed, so the catalogue was limited to include only galaxies in the redshift range $0 < z \leq 2.0$. A random subset of the remaining galaxies was chosen, and placed at random into either a “training” set (10 per cent of the sample), for which the galaxies true redshifts will be supplied, or a “test” set (the remaining 90 per cent of the sample), for which each code will need to predict a redshift PDF for each galaxy. Finally, we restrict our analysis to a sample with an apparent magnitude limite $i < 25.3$, which give a signal-to-noise ~ 30 for most galaxies, a cut often referred to as the expected “LSST Gold Sample”. This magnitude cut results in a training set with 44 404 galaxies and a test set containing 399 356 galaxies. All subsequent results will evaluate this “gold sample” test set.

2.1.2 Templates

As mentioned in Section 2.1, the SEDs in the Buzzard simulation are drawn from an empirical set of SEDs taken from the SDSS DR6 NYU-VAGC, a sample of roughly $\sim 5 \times 10^5$ galaxies with spectra in SDSS. To determine a finite set of templates to use with template fitting codes we take the five SED weight coefficients for each of the galaxies in the SDSS sample and run a simple K-means clustering algorithm on this five dimensional space. Each dimension was normalized such that it spanned an interval $[0, 1]$. The K-means clusters partition the five-dimensional space of coefficients into Voronoi cells, spanning the space of coefficients in a way that properly reflects the underlying density in the coefficients. Thus, the resultant SEDs constructed using the cell centers as weight coefficients will provide a reasonable spanning SED set. An ad-hoc number of $K = 100$ was chosen and the 100 K-means centre positions are taken as the weights for the K-CORRECT SED components to construct one hundred template SEDs. These 100 templates were provided, and the templates were used by both BPZ and LEPHARE; however, because EAZY was designed and written to use

the same five basis templates employed by K-CORRECT when constructing our mock galaxies, EAZY was run using linear combinations of these five templates rather than using the 100 discrete templates. The ability to fit for linear combinations of templates highlights an important implementation difference between similar photo-z codes.

2.1.3 Limitations

For our initial investigation of photometric redshift codes, we begin with a data set that is somewhat idealized, and does not contain all of the complicating factors present in real data. In several cases, the simplification is done with a purpose, with potentially confounding effects excluded in order to better isolate the differences between current-generation photo-z codes, and their causes. We list several of the simulations limitations in this section. As the simulation is catalogue-based, no image level effects, such as photometric measurement effects, object blending, contamination from sky background (Zodiacal light, scattered light, etc...), lensing magnification, or Galactic reddening are included. No stars are included in the catalogue, nor are the effects of AGN. As all SEDs are constructed from only five basis templates, properties of the galaxy population will be restricted to follow linear combinations of the characteristics of the five basis templates, so certain non-linear features, for example the full range of emission line fluxes relative to the continuum, will not be included in the model galaxy population. Moreover, the linear combinations of templates are modeled on the $\sim 5 \times 10^5$ SDSS galaxies discussed in Section 2.1, and thus only galaxies that resemble those spectroscopically observed by the SDSS will be included in the sample. No additional dust reddening intrinsic to the host galaxy is included, the only approximation of dust extinction comes in the form of dust encoded in the five basis SEDs via the training set used to create the basis templates. Simple linear combinations of these basis templates will, once again, not explore the full range of realistic dust extinction observed in galaxy populations. While these idealized conditions limit the realism of our galaxy population, some are also by design. We aim to test the photo-z codes at a very basic level, and a simplified model assures that differences in results seen between the codes are due to fundamental differences in their underlying assumptions and implementation details, rather than more nuanced implementation details.

3 METHODS

Here we outline the photo-z PDF codes tested in this study. In total, eleven distinct codes are tested. This sample is not comprehensive, but does cover a broad range of current-generation codes. Both template-based and machine learning approaches are included and each are described separately in Secs. 3.1 and 3.2 respectively. The list of codes are summarized in Table. 1.

The questions that must be answered for each code are: what unique features are included in the specific implementation that influence the output $p(z)$. What form of validation was performed with the training data, how were photometric uncertainties employed in the analysis, how were

Table 1. List of photo- z codes featured in this study. ML here means machine learning.

Code	Type	Paper	Website
BPZ	template	Benítez (2000)	http://www.stsci.edu/~dcoe/BPZ/
EAZY	template	Brammer et al. (2008)	https://github.com/gbrammer/eazy-photoz
LePHARE	template	Arnouts et al. (1999)	http://www.cfht.hawaii.edu/~arnouts/lephare.html
ANNz2	ML	Sadeh et al. (2016)	https://github.com/IftachSadeh/ANNz2
DELIGHT	ML/template	Leistedt & Hogg (2017)	https://github.com/ixkael/Delight
FLEXZBOOST	ML	Izbicki & Lee (2017)	https://github.com/tospisic/flexcode; https://github.com/rizbicki/FlexCoDE
GPz	ML	Almosallam et al. (2016b)	https://github.com/OxfordML/GPz
METAPHOR	ML	Cavuoti et al. (2017a)	http://dame.dsfa.unina.it
CMNN	ML	Graham et al. (2018)	-
SKYNET	ML	Graff et al. (2014)	http://ccpforge.cse.rl.ac.uk/gf/project/skynet/
TPZ	ML	Carrasco Kind & Brunner (2013)	https://github.com/mgckind/MLZ
TRAINZ	N/A	See Section 3.3	

negative fluxes treated, what specific prior form was employed (for template based codes), or what specific machine learning architecture was used (for ML codes)?

3.1 Template-based Approaches

3.1.1 BPZ

BPZ⁶ (Bayesian Photometric Redshift, Benítez 2000) is a template-based photo- z code that compares the expected colors (C) calculated for a set of spectral energy distribution (SED) types/templates (T) to the observed colors to calculate the likelihood of observing colors at each redshift for each type, $p(C|z, T)$. The code employs an empirically determined Bayesian prior in apparent magnitude (m_0) and SED-type. Assuming that the SED-types are spanning and exclusive, we can determine the redshift posterior $p(z|C, m_0)$ by marginalizing over all SED-types with a simple sum (Eq. 3 from Benítez 2000):

$$p(z|C, m_0) \propto \sum_T p(z, T|m_0) p(C|z, T) \quad (1)$$

where the first term on the right-hand side is the Bayesian prior and the second term is the traditional likelihood. The prior is assumed to have the form: $p(z, T|m_0) = p(T|m_0) p(z|T, m_0)$, i.e. it parameterizes the prior as an evolving type fraction with apparent magnitude, combined with a prior on the expected redshift probability distribution as a function of both apparent magnitude and SED-type.

In this paper we use BPZ v 1.99.3. The template set employed here is the set of 100 discrete SEDs described in Section 2.1.2 To keep the number of free parameters to a manageable level the SEDs in the training set are sorted by the rest-frame $u-g$ colour and split into three “broad” SED classes, equivalent to the E, Sp and Im/SB types in Benítez (2000). We assume the same functional form for the Bayesian priors as used by Benítez (2000), and utilize the training-set galaxies with known SED-type, redshift, and apparent magnitude to determine the type fractions and the best fit for the eleven free parameters of the prior. For galaxies with negative flux in a measured band, the placeholder

value is replaced with an estimate one σ detection limit in that particular band, i. e. a value close to the estimated sky noise threshold. The type-marginalized $p(z)$ is generated by setting the parameter PROBS_LITE=TRUE in the BPZ parameter file.

3.1.2 EAZY

EAZY⁷ (Easy and Accurate Photometric Redshifts from Yale, Brammer et al. 2008) is a template-based photo- z code that includes several features that improve on the basic χ^2 fit used in many template codes. It can fit the observed photometry with SEDs created from a linear combination of a set of templates at each redshift, and the best-fit SED is found by simultaneously fitting one, two or all of the templates by minimizing χ^2 . The minimized $\chi^2(z)$ is then combined with an apparent magnitude prior to obtain the posterior redshift probability distribution, although some argue that this is not the mathematically correct way of calculating the posteriors. EAZY can also account for the uncertainties in the templates by adding an empirically derived template error in quadrature as a function of redshift to the flux errors.

In this paper we use the all-templates mode, which fits the photometric data with a linear combination of the five basis templates. We employed the 5 basis templates described in Section 2.1, and set the template error to zero since these same templates were used to produce the simulated catalog photometry. The likelihoods are calculated on a 200-point redshift grid spanning $0 \leq z \leq 2$, and include the application of a type-independent apparent magnitude prior estimated from the training data.

3.1.3 LePhare

LEPHARE⁸ (Photometric Analysis for Redshift Estimate, Arnouts et al. 1999; Ilbert et al. 2006) is a photo- z reconstruction code based on a χ^2 template-fitting procedure. The observed colors are matched with the colours predicted from

⁶ <http://www.stsci.edu/~dcoe/BPZ/>

⁷ <https://github.com/gbrammer/eazy-photoz>

⁸ <http://www.cfht.hawaii.edu/~arnouts/lephare.html>

6 LSST Dark Energy Science Collaboration

a set of spectral energy distribution (SED) which can be either synthetic or based on a semi-empirical approach. LEP-HARE has been used to produce the COSMOS2015 photo- z catalogue (Laigle et al. 2016).

Each SED is redshifted in steps of $\Delta z = 0.01$ and convolved with the simulated LSST filter transmission curves (accounting for instrument efficiency). The opacity of the inter-galactic medium has been set to zero as no additional reddening has been included in the Buzzard simulations. The computed photo- z is then the value that minimizes the merit function $\chi^2(z, T, A)$ from Arnouts et al. (1999):

$$\chi^2(z, T, A) = \sum_f^{N_f} \left(\frac{F_{\text{obs}}^f A \times F_{\text{pred}}^f(T, z)}{\sigma_{\text{obs}}^f} \right)^2 \quad (2)$$

where A is a normalization factor, $F_{\text{pred}}^f(T, z)$ is the flux predicted for a template T at redshift z . F_{obs}^f is the observed flux in a given band f and σ_{obs}^f the associated observational error. The index f refers to the considered band and N_f is the total number of filters.

In this paper we use LEPHARE v 2.2. The set of templates used for fitting the photo- z 's are the 100 discrete Buzzard SED templates as described in section 2.1.2. The full $p(z)$ corresponds to the likelihoods calculated at each point on our z -grid.

3.2 Training-based Codes

3.2.1 ANNz2

ANNz2⁹ (Sadeh et al. 2016) is a powerful package that has the ability to employ several machine learning algorithms, including artificial neural networks (ANN), boosted decision tree (BDT) and k-nearest neighbour (KNN). Using the Toolkit for Multivariate Data Analysis (TMVA) with ROOT¹⁰, it can run multiple machine learning algorithms for a single training and outputs photo- z 's based on a weighted average of their performances.

ANNz2 is capable of producing both photo- z point estimates and redshift posterior probability distributions $p(z)$, it could also conduct classifications and supports reweighting between samples. The PDFs are produced by propagating the intrinsic uncertainty on the input parameters and the uncertainty in the machine learning method to the expected photo- z solution, averaged over multiple runs weighted based on the performance of each run. ANNz2 presents its photo- z uncertainty different from many codes by using the KNN method: it estimates the photo- z bias between each object and a fixed number of nearest neighbours in parameter space, it then takes the 68th percentile width of the distribution of the bias. This is based on the implication that objects with similar photometric properties would have similar uncertainties, and therefore the photometric errors of the inputs are not propagated into the code.

In this study, ANNz2 v. 2.0.4 was used. The full PDF for each galaxy is also produced with a linear stepsize of $z = 0.01$ for $0 \leq z \leq 2$. A set of 5 ANNs with architecture 6 : 12 : 12 : 1 (6 *ugrizy* inputs, 2 hidden layers with

12 nodes each, and 1 output) with different random seeds are used during each training. Half of the training set is used as a validation set to prevent overtraining. All training objects are set to have detected magnitudes, however the non-detections ($\text{mag} = -99$) in the testing set are replaced with the mean of that particular band.

3.2.2 Color-Matched Nearest-Neighbours

The nearest-neighbours color-matching photometric redshift estimator (CMNN) is presented in (Graham et al. 2018, hereafter G18). This method uses a training set of galaxies with known redshifts that has equivalent or better photometry as the test set in terms of quality and filter coverage. For each galaxy in the test set we identify a color-matched subset of training galaxies, choose one (e.g. the nearest-neighbour or a random selection), and use its known redshift as the photo- z . This color-matched subset is identified by first calculating the Mahalanobis distance D_M in color-space between the test galaxy and all training-set galaxies: the difference between the test and a training set galaxy's color divided by the photometric error, summed over all colors (i.e., $u-g$, $g-r$, $r-i$, $i-z$, and $z-y$). Then, the threshold value for D_M that define a good color match is set by the percent point function (PPF): for example, for $N_{\text{dof}} = 5$, PPF = 95 per cent of all training galaxies consistent with the test galaxy will have $D_M < 11.07$ (where N_{dof} , the number of degrees of freedom, is the number of colors). For a given test galaxy, the $p(z)$ is the normalized distribution of the true catalogue redshifts of this color-matched subset of training galaxies, and the standard deviation of the color-matched subset is used as the photo- z uncertainty.

We have applied the nearest-neighbours color-matching photometric redshift estimator described in G18 to the simulated data. Compared to its application in G18, there are some minor differences in the application of this estimator to the Buzzard catalogue. First, we do not impose non-detections on galaxies with a magnitude fainter than the expected LSST 10-year limiting magnitude or bright enough to saturate with LSST: *all* of the photometry for all the galaxies in the test and training sets are used for this experiment. Second, as in G18 we do apply an initial cut in color to the training set before calculating the Mahalanobis distance in order to accelerate processing, and also use a magnitude pseudo-prior to improve photo- z estimates, but for both we have used different cut-off values that are appropriate for the Buzzard galaxies' colors and magnitudes. Third, we set different parameters for the identification of the color-matched subset of training galaxies and the selection of a photometric redshift estimate. In G18 we used a percent point function (PPF) value of 0.68 to identify the color-matched subset of training galaxies and used the redshift of nearest neighbour in color-space as the photo- z estimate. These choices work well when the desire is to obtain accurate photo- z estimates for most test-set galaxies, but does not return a robust $p(z)$ in all cases – especially for galaxies that are bright and/or have few matches in color-space. Since a robust estimate of $p(z)$ is desired for this work we make several changes to our implementation of the CMNN photo- z estimator. We continue to use a percent point function of PPF = 0.95 to generate the subset of color-matched training galaxies, but weight them by the inverse of their Mahalanobis distance.

⁹ <https://github.com/IftachSadeh/ANNZ>

¹⁰ <http://tmva.sourceforge.net/>

This weighting maintains some of the accuracy that was previously achieved by simply using the nearest neighbour in color-space. We then use the weights to create the $p(z)$ instead of having the redshift of each color-matched training-set galaxy count equally. To obtain a robust estimate of the $p(z)$ for galaxies with a small number of color-matched training set galaxies, when this number is less than 20 the nearest 20 neighbours in color-space are used instead, and we convolve the $p(z)$ with a Gaussian with a standard deviation of:

$$\sigma = \sigma_{\text{train}} \sqrt{(\text{PPF}_{20}/0.95)^2 - 1} \quad (3)$$

to appropriately broaden it so that the $p(z)$ for these test galaxies represents the enlarged PPF value associated with it. Overall, these three changes will yield poorer accuracy photo- z compared to those presented in G18, but they will all have significantly more robust estimates of the $p(z)$, particularly for the brightest test galaxies. This is sufficient for this work because, as described in G18, the goal of the CMNN photo- z estimator was never to provide the “best” (or even competitive) estimates in the first place, given its reliance on a deep training set, but rather to provide a means for direct comparisons between LSST photometric quality and photo- z estimates. With this work we show how the input parameters should be set in order to return robust $p(z)$ estimates in addition to point value estimates.

3.2.3 Delight

DELIGHT¹¹ (Leistedt & Hogg 2017) infers photo- z ’s by using a data-driven model of latent SEDs and a physical model of photometric fluxes as a function of redshift. Generally, machine learning methods rely on representative training data with similar band passes, while template based methods rely on a complete library of templates based on physical models constructed. DELIGHT is constructed in attempt to combine the advantages and eliminate the disadvantages of both template-based and machine learning algorithms: it constructs a large collection of latent SED templates (or physical flux-redshift models) from training data, with a template SED library as a guide to the learning of the model. The advantage of DELIGHT is that it neither needs representative training data in the same photometric bands, nor does it need detailed galaxy SED models to work.

This conceptually novel approach is done by using Gaussian processes operating in flux-redshift space. The posterior distribution on the redshift of a target galaxy is obtained via a pairwise comparison with training galaxies,

$$p(z|\hat{\mathbf{F}}) \approx \sum_i p(\hat{\mathbf{F}}|z, t_i) p(z|t_i) p(t_i), \quad (4)$$

where $p(z|t_i)p(t_i)$ captures prior information about the redshift distributions and abundances of the galaxies, with t_i denoting the galaxy template; while $p(\hat{\mathbf{F}}|z, t_i)$ is the poste-

rior of noisy flux $\hat{\mathbf{F}}$ at redshift z . For each training-target pair, $p(\hat{\mathbf{F}}|z, t_i)$ is evaluated as follows:

$$p(\hat{\mathbf{F}}|z, t_i) = \int p(\hat{\mathbf{F}}|\mathbf{F}) p(\mathbf{F}|z, z_i, \hat{\mathbf{F}}_i) d\mathbf{F}, \quad (5)$$

where $p(\hat{\mathbf{F}}|\mathbf{F})$ is the likelihood function, it compares the noisy real flux $\hat{\mathbf{F}}$ with the noiseless flux \mathbf{F} obtained from the linear combination of template models, carefully constructed to account for model uncertainties and different normalization of the same SED; while $p(\mathbf{F}|z, z_i, \hat{\mathbf{F}}_i)$ is the prediction of flux at a different redshift z with respect to the training object with redshift z_i and flux $\hat{\mathbf{F}}_i$. Eq. 5 is essentially the probability that the training and the target galaxies having the same SED but at a different redshift. The flux prediction $p(\mathbf{F}|z, z_i, \hat{\mathbf{F}}_i)$ of the training galaxy at redshift z is modeled via a Gaussian process,

$$F_b \sim \mathcal{GP}(\mu^F, k^F), \quad (6)$$

with mean function μ^F and kernel k^F , both imposed to capture expected correlations resulting from the known underlying physics (i.e., fluxes resulting from observing SEDs through filter response, and the SEDs being redshifted). The reader should refer to Leistedt & Hogg (2017) for further details.

In this study, all 100 ordered Buzzard templates, as described in Section 2.1.2, were used in DELIGHT, and the Gaussian process was trained with a subset of 50 000 galaxies. Photometric uncertainties from the inputs are propagated into the code, while non-detections for each band are set to the mean of the respective bands. Default settings of DELIGHT were used, with the exception that the PDF bins were set to be linear instead of logarithmic, with 200 equally-spaced bins between $0.0 < z < 2.0$. In this study a flat prior is assumed.

3.2.4 FlexZBoost

FLEXZBOOST¹² (Izbicki & Lee 2017) is a particular realization of FlexCode, which is a general-purpose methodology for converting any conditional mean point estimator of z to a conditional density estimator $f(z|\mathbf{x})$, where \mathbf{x} here represents our photometric covariates and errors.¹³ The key idea is to expand the unknown function $f(z|\mathbf{x})$ in an orthonormal basis $\{\phi_i(z)\}_i$:

$$f(z|\mathbf{x}) = \sum_i \beta_i(\mathbf{x}) \phi_i(z). \quad (7)$$

By the orthogonality property, the expansion coefficients are just conditional means

$$\beta_i(\mathbf{x}) = \mathbb{E}[\phi_i(z)|\mathbf{x}] \equiv \int f(z|\mathbf{x}) \phi_i(z) dz. \quad (8)$$

¹² <https://github.com/tospispi/flexcode>; <https://github.com/rizbicki/FlexCoDE>

¹³ Instead of $p(z)$, we use the notation $f(z|\mathbf{x})$ to explicitly show the dependence on \mathbf{x} .

¹¹ <https://github.com/ixkael/Delight>

8 LSST Dark Energy Science Collaboration

These coefficients can easily be estimated from data by regression.

In this paper, we use XGBOOST (Chen & Guestrin 2016) for the regression part as these techniques scale well for massive data; it should however be noted that FLEXCODE-RF (also on GitHub), based on Random Forests, generally performs better for smaller data sets. As our basis, we choose a standard Fourier basis. There are two tuning parameters in our $p(z)$ estimate: (i) the number of terms, I , in the series expansion in Eq. 7, and (ii) an exponent α that we use to sharpen the computed density estimates $\hat{f}(z|\mathbf{x})$, according to $\hat{f}(z|\mathbf{x}) \propto \hat{f}(z|\mathbf{x})^\alpha$. We split the “train data” into a training set (85%) and a validation set (15%), and choose both I and α in an automated way by minimizing the weighted L_2 -loss function (Eq. 5 in Izbicki & Lee 2017) on the validation set.

Although FlexCode offers a *lossless compression* of the photo- z estimates (in this study, one can reconstruct $\hat{f}(z|\mathbf{x})$ exactly at any resolution from estimates of the first 35 coefficients, Eq. 8, for a Fourier basis $\{\phi_i(z)\}_i$), we discretize our final estimates into 200 bins linearly spaced in $0 < z < 2$ for easy comparison with other algorithms. Using a higher resolution may yield better results (with no added cost in storage).

3.2.5 GPz

GPz¹⁴ (Almosallam et al. 2016a,b) is a sparse Gaussian process based code, a fast and a scalable approximation of full Gaussian Processes (Rasmussen & Williams 2006), with the added feature of being able to produce input-dependent variance estimations (heteroscedastic noise). The model assumes that the probability of the output y , the redshift, given the input x , the photometry, is $p(y|x) = \mathcal{N}(y|\mu(x), \sigma(x)^2)$. The mean function, $\mu(x)$, and the variance function $\sigma(x)^2$ are both linear combinations of basis functions that take the following form:

$$f(x) = \sum_{i=1}^m \phi_i(x) w_i, \quad (9)$$

where $\{\phi_i(x)\}_{i=1}^m$ and $\{w_i\}_{i=1}^m$ are sets of m basis functions and their associated weights respectively. Basis function models (BFM), for specific classes of basis functions such as the sigmoid or the squared exponential, have the advantage of being universal approximators, i.e. there exist a function of that form that can approximate any function, with mild assumptions, to any desired degree of accuracy. The details on how to learn the parameters of the model and the hyper-parameters of the basis functions are described in Almosallam et al. (2016b).

A unique feature in GPz, is that the variance estimate is composed of two terms each quantifying a different source of uncertainty. One term (the model uncertainty) reflects how much of the uncertainty is due to lack of training samples at the location of interest, whereas the second term (the noise uncertainty) reflects how much of the uncertainty is caused from observing many noisy samples at that location. Thus, the predictive variance can determine whether we need more

representative samples or more precise samples for any particular location in the input space. GPz can also emphasize the importance of some samples as weights. This weight can be for example $|z_{\text{spec}} - z_{\text{phot}}|/(1 + z_{\text{spec}})$ to target the desired objective of minimizing the normalized redshift error or as a function of their probability in the test set relative to the training set in order to pressure the model to better fit samples that are rare in the training set but are expected to be abundant during testing.

The data is prepared for GPz by taking the log of the magnitude errors, decorrelating the data set using PCA and imputing the missing values using a simple linear model that estimates the missing variables given the observed ones. The log transformation helps to smooth the long tail distribution of the magnitude errors, which is more stable numerically and makes the optimization process unconstrained. The missing values are imputed by computing the mean of the training set μ and its covariance Σ , then we use the following equation to estimate the missing values from the observed ones

$$x_u = \mu_u + \Sigma_{uo}\Sigma_{oo}^{-1}(x_o - \mu_o), \quad (10)$$

where the subscript o in x_o indexes the *observed* part of the input x , whereas the subscript u indexes the *unobserved* set (similarly for μ and Σ). This is the optimal expected value of the unobserved variables given the observed ones if the distribution is jointly Gaussian, note that if the variables are independent, i.e. $\Sigma_{uo} = 0$, this will reduce to a simple average predictor.

We use the Variable Covariance (VC) option in GPz with 200 basis functions after we note that there is no significant increase in the performance on the validation set (using 80%-20% training-validation split) and with no cost-sensitive learning applied.

3.2.6 METAPhOR

METAPhOR (Machine-learning Estimation Tool for Accurate Photometric Redshifts, Cavuoti et al. 2017a) is a pipeline designed to provide photo- z ’s point estimates and a reliable PDF for machine learning (ML) based techniques. It includes pre- and post-processing phases, hosting a photo- z prediction engine based on the Multi Layer Perceptron with Quasi Newton Algorithm (MLPQNA), already validated on photo- z ’s in several cases (de Jong et al. 2017; Cavuoti et al. 2017b, 2015; Brescia et al. 2014, 2013; Biviano et al. 2013). Due to its plug-in based modular nature, METAPhOR can be easily replaced by any other photo- z prediction kernel, regardless its implementation, by taking the I/O interface compliance as unique constrain.

At a higher level, the pipeline mainly consists of three modules: (i) *data pre-processing*, including a catalogue cross-matching sub-module (based on the tool C3, Riccio et al. 2017), a sub-module for photometric evaluation and error estimation of the multi-band catalogue used as Knowledge Base (KB), and a sub-module dedicated to the perturbation of the photometric KB, propaedeutic to the PDF estimation; (ii) *photo- z prediction*, which is the training/validation/test phase, producing the photo- z ’s point estimates, based on a pre-selected ML method; (iii) *PDF estimation*, specifically designed to calculate the PDF of the photo- z estimation

¹⁴ <https://github.com/OxfordML/GPz>

errors. The last module includes also a post-processing tool, providing some statistics on the produced point estimates and PDFs.

The photometry perturbation law is based on the formula $m_{ij} = m_{ij} + \alpha_i F_{ij} * u_{\mu=0,\sigma=1}$, where α_i is a user selected multiplicative constant (useful in case of multi-survey photometry), $u_{\mu=0,\sigma=1}$ is a random value from the standard normal distribution and F_{ij} is a bimodal function (a constant function + polynomial fitting of the mean magnitude errors on the binned bands), heuristically tuned in such a way that the constant component is the threshold under which the polynomial function is considered too low to provide a significant noise contribution to the photometry perturbation.

As introduced, the photo- z point estimate prediction engine of METAPHOR is based on the MLPQNA model, whose photo- z regression training error, used by the quasi Newton learning rule, is based on the least square error and Tikhonov L_2 -norm regularization (Hofmann & Mathé 2018).

As main prerogative, METAPHOR is able to provide a PDF for ML methods by taking into account the photometric errors provided with data, by running N trainings on the same training set, or M trainings on M different random extractions from the KB. The different test sets, used to produce the PDF, are thus obtained by introducing a proper perturbation, parametrized from the photometric error distribution in each band, on the photometric data populating the original test set (Brescia et al. 2018).

For the present work since it was required to produce a redshift (and a PDF) for each object of the test set we decided to apply a hierarchical kNN to fill the missing detection, it goes without saying that for such points the reliability of PDFs and point estimation is lower. No cross validation has been used.

3.2.7 SkyNet

SKYNET¹⁵ (Graff et al. 2014) is a publicly available neural network software, based on a 2nd order conjugate gradient optimization scheme (see Graff et al. 2014, for further details). It has been used efficiently for redshift PDF estimates (Sánchez et al. 2014; Bonnett 2015; Bonnett et al. 2016).

The neural network is configured as a standard multi-layer perceptron with three hidden layers and one input layer with 12 nodes (the 6 magnitudes and their errors). The classifier is laid out such that the hidden layers have 20:40:40 nodes each, all rectified linear units, and the output layer has 200 nodes (corresponding to 200 bins for the PDF) activated with a “softmax” function so that they automatically sum to 1.

To avoid over-fitting, a 30 per cent fraction of the training set is used as validation, and the training is stopped as soon as the error rate begins to increase in the validation set. The weights are randomly initialized based on normal sampling. The error function is a standard chi-square function for the regressor, and a cross-entropy function for the classifier. Finally, the data are all whitened before processing, with magnitudes pegged to (45,45,40,35,42,42) and their

errors pegged to (20,20,10,5,15,15) for *ugrizy* filters, respectively.

3.2.8 TPZ

TPZ¹⁶ (Trees for Photo- z , Carrasco Kind & Brunner 2013; Carrasco Kind & Brunner 2014) is a parallel machine learning algorithm that generates photometric redshift PDFs using prediction trees and random forest techniques. The code recursively splits the input data (i. e. the training sample), into two branches, one after another, until a terminal leaf is created that meets a termination criterion (e. g. a minimum leaf size or a variance threshold). Bootstrap samples from the training data and associated errors are used to build a set of prediction trees. In order to minimize correlation between the trees, the data is divided in such a way that the highest information gain among the random subsample of features is obtained at every point. The regions in each terminal leaf node corresponds to a specific subsample of the entire data that possesses similar properties.

The training data is examined before running TPZ. Since TPZ does not handle non-detections (magnitudes flagged as 99.0), we replace these values with an approximation of the 1σ detection threshold, i. e. a signal to noise ratio of 1 in terms of magnitude uncertainty using the equation $dm = 2.5 \log(1 + N/S)$ where $dm \sim 0.7526$ mag for $N/S = 1$. That is, for each band, we replace the non-detection with the magnitude corresponding to the error of 0.7526 from the error model forecasted for 10-year LSST data. The Out-of-Bag (Breiman et al. 1984; Carrasco Kind & Brunner 2013) cross-validation technique is used within TPZ to evaluate its predictive validity and determine the relative importance of the different input attributes. We employed this information to calibrate our algorithm.

In the present work, the LSST magnitudes u , g , r , i and colors $u-g$, $g-r$, $r-i$, $i-z$, $z-y$ and their associated errors are used in the process of growing 100 trees with a minimum leaf size of 5 (the z and y magnitudes did not show significant correlation with the redshift in our cross-validation, so we did not use them when constructing our trees). We partitioned our redshift space into 100 bins from $z = 0.005$ to $z = 2.0$ and smoothed each individual PDF with a smoothing scale of twice the bin size.

3.3 Simple Ensemble Estimator

In addition to the main photo- z algorithms described above we also include a very simple method. For TRAINZ, as we will we call this simple estimator, we well define $p(z)$ as simply:

$$p(z) = \frac{1}{N_{train}} \sum_{i=1}^{N_{train}} z_{train} \quad (11)$$

That is, we simply set the redshift PDF of every galaxy equal to the normalized $N(z)$ of the training sample. As the training sample is drawn from the same underlying distribution as the test sample, modulo small deviations due to sample size, the quantiles of the training and test distributions

¹⁵ <http://ccpforge.cse.rl.ac.uk/gf/project/skynet/>

¹⁶ <https://github.com/mgckind/MLZ>

834 should be identical. This is a wildly unrealistic estimator, as
 835 it assigns all galaxies, no matter their apparent magnitude,
 836 colour, or true redshift, the same redshift PDF, and is thus
 837 uninformative at the level of individual object redshifts, but
 838 is designed to perform very well for the ensemble of all ob-
 839 jects. We will discuss this method and cautions relative to
 840 metrics in Section 5.3.

841 4 METRICS FOR QUANTIFYING PDF 842 COMPARISONS

843 The overloaded “ $p(z)$ ” is a widespread abuse of notation;
 844 we would like the outputs of photo- z PDF codes to be in-
 845 terpretable as probabilities. Obviously photo- z PDFs must
 846 not take negative values and must integrate to unity over
 847 the range of possible redshifts. Additionally, an estimator
 848 derived by method H for the photo- z PDF of galaxy i must
 849 be understood as a posterior probability distribution

$$850 \hat{p}_j(z_i) = p(z|d_i, I_D, I_H), \quad (12)$$

851 conditioned not only on the photometric data d_i for that
 852 galaxy but also on parameters encompassing a number of
 853 things that will differ depending on the method H used to
 854 produce it, namely the assumptions I_H necessary for the
 855 method to be valid and any inputs I_D it takes as prior infor-
 856 mation, such as a template library or training set. Because of
 857 this, direct comparison of photo- z PDFs produced by differ-
 858 ent methods is in some sense impossible; even if they share
 859 the same prior information I_D , by definition they cannot
 860 be conditioned on the same assumptions I_H , otherwise they
 861 would not be distinct methods at all.

862 In this study, we isolate the differences in prior infor-
 863 mation specific to each method by using a single training set
 864 I_D^{ML} for all machine learning-based codes and a single tem-
 865 ple library I_D^T for all template-based codes, and these sets
 866 of prior information are carefully constructed to be represen-
 867 tative and complete, we have $I_D^{ML} \equiv I_D^T$ for every method
 868 H . Thus, we are saying

$$869 \frac{\hat{p}_{i,H}(z)}{\hat{p}_{i,H'}(z)} \approx \frac{p(z|d_i, I_H)}{p(z|d_i, I_{H'})}, \quad (13)$$

870 meaning that we assume comparisons of $\hat{p}_{i,H}(z)$ isolate the
 871 effect of the method used to obtain the estimator, which
 872 should make examination of differences caused by specifics
 873 of the method implementations easier to isolate.

874 As mentioned previously, there are cosmology probes
 875 that require knowledge of individual galaxy $p(z)$ and others
 876 that require only knowledge of the ensemble redshift distri-
 877 bution, $N(z)$. Due to the paucity of principled techniques
 878 for using and validating photo- z PDFs, there are few alter-
 879 natives to the common practice of reducing photo- z PDFs
 880 to point estimates. Though this practice should not be en-
 881 couraged, we also calculate traditional metrics based on the
 882 most common point estimators derived from photo- z PDFs.
 883 Those seeking to establish a connection to traditional ways
 884 of thinking about redshift estimation may consult the Ap-
 885 pendix for these results.

886 There are a number of metrics that can be used to test
 887 the accuracy of a photo- z interim posterior as an estimator

888 of a true photo- z posterior if it is known. Even for simulated
 889 data, the true photo- z PDF is in general not accessible un-
 890 less the redshifts are in fact drawn from the true photo- z
 891 PDFs, a mock catalogue generation procedure that has not
 892 yet appeared in the literature. Furthermore, only limited ap-
 893 plications of photo- z PDFs that could be used as the basis
 894 for a metric have been presented in the literature. The most
 895 popular application by far is the calculation of the overall
 896 redshift distribution $N(z)$, the true value of which is known
 897 for the BUZZARD simulation and will be denoted as $N'(z)$.
 898 Though alternatives exist (Malz & Hogg prep), stacking ac-
 899 cording to

$$900 \hat{N}^H(z) \approx \frac{1}{N_{tot}} \sum_i^{N_{tot}} \hat{p}_i^H(z) \quad (14)$$

901 is the most widely accepted method for estimating the red-
 902 shift distribution from photo- z PDFs. If we assume that
 903 the response of estimators of $N(z)$ is uniform across all ap-
 904 proaches H , then we may interpret metrics on the accuracy
 905 of $\hat{N}(z)$ obtained in this way. We must note, however, that
 906 this is a poor assumption in general. Under the setup of this
 907 paper, the true redshift distribution $N'(z) = p(z|I_D)$ (i.e.
 908 because our training data is representative, the interim prior
 909 is the truth). In this ideal case, the method that would give
 910 the best approximation to $N'(z)$ would be one that neglects
 911 all the information contained in the photometry $\{d_i\}_{N_{tot}}$
 912 and gives every galaxy the same photo- z PDF $\hat{p}_i(z) = N'(z)$
 913 for all i . This is the exact estimator, TRAINZ, that we have
 914 described in Section 3.3, and which will serve as a point of
 915 reference for the other codes.

916 The exact implementation of the stacked estimator
 $\hat{N}^H(z)$ will depend on the parametrization of the photo- z
 917 PDFs, which may differ across codes and can affect the pre-
 918 cision of the estimator (Malz et al. 2018); even considering a
 919 single method under the same parametrization, say a piece-
 920 wise constant function over bins or a set of samples from
 921 the posterior, an estimator using $2N$ bins or samples will
 922 trivially be more precise than an estimator using N bins or
 923 samples. In order to minimize the effects of such choices,
 924 we asked those running all eleven codes to output $p(z)$ para-
 925 meterized with a generous ≈ 200 piecewise constant bins
 926 spanning $0 < z < 2$. The piecewise constant format is chosen
 927 because of its established presence in the literature, and the
 928 choice of 200 bins was motivated by the approximate number
 929 of columns expected to be available for storage of $p(z)$ for
 930 the final LSST Project tables.¹⁷ All $p(z)$ catalogues are pro-
 931 cessed using the QP software package (Malz et al. 2018)¹⁸
 932 for manipulating and calculating metrics of 1-dimensional
 933 PDFs. We will discuss the choice of $p(z)$ parameterization
 934 further in Section 5.

¹⁷ See, e. g. the LSST Data Products Definition Document, available at: <https://ls.st/dpdd>

¹⁸ available at: <http://github.com/aimalz/qp>

936 **4.1 Metrics of an ensemble of photo-z interim
937 posteriors**

938 **4.1.1 Probability integral transform (PIT)**

939 The probability integral transform (PIT) (Polsterer et al.
940 2016) is defined for each individual galaxy as:

$$941 \quad \text{PIT} = \int_{-\infty}^{z_{\text{true}}} p(z) dz. \quad (15)$$

943 The distribution of PIT values quantifies the behavior of the
944 *ensemble* of photo- z PDFs, enabling us to evaluate whether
945 the $p(z)$ is, on average, accurate: The PIT value is the Cumu-
946 lative Distribution Function (CDF) of the $p(z)$ evaluated at
947 the true redshift. A catalogue of photo- z PDFs that are ac-
948 curate should have a flat PIT histogram (i.e., the individual
949 PIT values as samples from each CDF should match a Uni-
950 form(0,1) distribution if the CDFs are accurate). Specific
951 deviations from flatness indicate inaccuracy: overly broad
952 photo- z PDFs would manifest as underrepresentation of the
953 lowest and highest PIT values, whereas overly narrow photo-
954 z PDFs would manifest as over-representation of the lowest
955 and highest PIT values. High frequency at only PIT ≈ 0
956 and PIT ≈ 1 indicates the presence of catastrophic outliers
957 with highly inaccurate photo- z PDFs where the true red-
958 shift is outside of the support of $p(z)$. Tanaka et al. (2017)
959 use the histogram of PIT values as a diagnostic indicator of
960 overall code performance, while Freeman et al. (2017) inde-
961 pendently define the PIT and demonstrate how its individ-
962 ual values may be used both to perform hypothesis testing
963 (via, e.g., the KS, CvM, and AD tests; see below) and to
964 construct quantile-quantile plots.

965 **4.1.2 Quantile-quantile (QQ) plot**

966 The quantile-quantile (QQ) plot is a graphical method for
967 comparing two distributions, where the quantiles of one dis-
968 tribution are plotted against the quantiles of the other distri-
969 bution (A quantile being defined by partitioning a distribu-
970 tion into consecutive intervals containing equal amounts of
971 probability, or equal numbers of objects in each interval). In
972 this paper we show the quantiles of the PIT values compared
973 to the quantiles of the Uniform distribution that we expect
974 the PIT values to match if $p(z)$ is an accurate probability dis-
975 tribution for all objects. The QQ plot provides an easy way
976 to qualitatively assess the differences in various properties
977 such as the moments of an estimating distribution relative
978 to a true distribution. In this paper, QQ plots are used for
979 two purposes: (1) for comparing $N(z)$ from photo- z PDFs
980 (estimated using Eq. 14) with the true $N(z)$, i.e. comparing
981 the estimated distribution of redshifts with the true redshift
982 distribution, and (2) for assessing the overall consistency of
983 an ensemble of photo- z PDFs with their true redshifts on
984 a population level, where the distribution of the PIT values
985 (see previous section) is compared to a uniform distribution
986 between 0 and 1. The QQ plot contains very similar infor-
987 mation to that shown in the PIT histogram plot, we include
988 both forms, as visually they each convey the information in
989 a somewhat distinct manner.

990 **4.1.3 Conditional density estimation loss**

991 With the conditional density estimation loss (CDE loss) we
992 can compare how well different methods estimate individual
993 PDFs for photometric covariates \mathbf{x} rather than looking only
994 at the ensemble distribution. As in Section 3.2.4, we use the
995 notation $f(z|\mathbf{x})$ instead of $p(z)$ to explicitly show the
996 dependence on \mathbf{x} .

997 The CDE loss is defined as:

$$998 \quad L(f, \hat{f}) = \int \int (f(z | \mathbf{x}) - \hat{f}(z | \mathbf{x}))^2 dz dP(\mathbf{x}) \quad (16)$$

999 This loss is the CDE equivalent of the RMSE in regression.
1000 To estimate this loss we rewrite the loss as

$$1001 \quad \mathbb{E}_{\mathbf{X}} \left[\int \hat{f}(z | \mathbf{X})^2 dz \right] - 2\mathbb{E}_{\mathbf{X}, Z} \left[\hat{f}(Z | \mathbf{X}) \right] + K_f, \quad (17)$$

1002 where the first expectation is with respect to the marginal
1003 distribution of the covariates \mathbf{X} , the second expectation is
1004 with respect to the joint distribution of \mathbf{X} and Z , and K_f is a
1005 constant depending only upon the true conditional densities
1006 $f(z | \mathbf{x})$. For each method we can estimate these expecta-
1007 tions as empirical expectations on the test or validation data
1008 (Eq. 7 in Izbicki et al. 2017) without knowledge of the true
1009 densities.

1010 **4.2 Metrics over estimated probability
1011 distributions**

1012 In tandem with the QQ and PIT metrics introduced above,
1013 we additionally compute the following metrics comparing
1014 the empirical CDF of a distribution to the true or expected
1015 distribution. These metrics give a more quantitative mea-
1016 sure of the departure from ideal than the more visual PIT
1017 histogram and QQ plot. We compute metrics comparing the
1018 CDF of PIT values to a the CDF of a Uniform distribution,
1019 and also compute the CDF of the true redshift distribution
1020 $N'(z)$ compared the $\hat{N}(z)$ distribution derived from sum-
1021 ming the $p(z)$ as described in Eq. 14.

1022 **4.2.1 Root-mean-square error (RMSE)**

1023 We employ the familiar root-mean-square error:

$$1024 \quad \text{RMSE} = \sqrt{\int_{-\infty}^{\infty} (\hat{f}(z) - f'(z))^2 dz}, \quad (18)$$

1025 Though this metric does not account for the fact that the
1026 redshift distribution function is, in fact, a probability dis-
1027 tribution, it can still be interpreted as a measure of the in-
1028 tegrated difference between the estimated distribution and
1029 the true distribution, and it can be used to quantify the
1030 otherwise qualitative metrics.

1031 **4.2.2 Kolmogorov-Smirnov (KS) and related statistics**

1032 The *Kolmogorov-Smirnov* statistic N_{KS} is the maximum dif-
1033 ference between $F_{\text{phot}}(z)$ and $F_{\text{spec}}(z)$, the CDFs of the
1034 photo- z and spectroscopic redshift respectively:

$$1035 \quad N_{\text{KS}} = \max_z (|F_{\text{phot}}(z) - F_{\text{spec}}(z)|). \quad (19)$$

The KS test quantifies the similarity between two distributions, independent of binning. A lower N_{KS} value corresponds to more similar distributions.

We also consider two variants of the KS statistic: the Cramer-von Mises (CvM) and Anderson-Darling (AD) statistics. The CvM statistic is similar to the KS statistic as it is also computed from the distance between the measured CDF and the ideal CDF, but instead of the maximum distance, the CvM statistic calculates the average of the distance squared:

$$\omega^2 = \int_{-\infty}^{+\infty} (F_{\text{meas.}}(x) - F_{\text{ideal}}(x))^2 dF_{\text{ideal}} \quad (20)$$

The AD statistic is a weighted version of the CvM statistic, making it more sensitive to the tails of the distribution:

$$A^2 = n \int_{-\infty}^{+\infty} \frac{(F_{\text{meas.}}(x) - F_{\text{ideal}}(x))^2}{F_{\text{ideal}}(x)(1 - F_{\text{ideal}}(x))} dF_{\text{ideal}} \quad (21)$$

where n is the sample size.

4.2.3 Moments

For the $\hat{N}(z)$ distributions we additionally calculate the first three moments of the estimated redshift distribution for each code and compare them to the moments of the true redshift distribution $N'(z)$. The m th moment of a distribution is defined as

$$\langle z^m \rangle = \int_{-\infty}^{\infty} z^m N(z) dz. \quad (22)$$

Here, we use the moments of the stacked estimator of the redshift distribution function as the basis for a metric. The closer the moments of $\hat{N}(z)$ for a photo-z PDF method are to the moments of the true redshift distribution function $N'(z)$, the better the photo-z PDF method.

5 RESULTS

5.1 Ensembles of photo-z interim posteriors

Fig. 1 Shows the $p(z)$ produced by each of our eleven photo-z codes for four example galaxies which exemplify some prominent cases that arise when estimating photo-z PDFs: a narrow, unimodal redshift solution, a broader unimodal solution, a bimodal distribution, and a complex, multimodal distribution. The red vertical line represents the true redshift of the individual galaxy, and the blue curve represents the redshift probability. Several features are obvious even in these illustrative examples. ANNz2, METAPHOR, NN, and SKYNET all show an excess of small-scale features, which appear to be print-through of the underlying training set galaxies. GPZ (in its current implementation), on the other hand, always produces a single Gaussian, which broadens to cover the multi-modal redshift solutions seen in other codes.

As stated in Section 4, $p(z)$ is parameterized as ≈ 200 piecewise constant bins covering $0 < z < 2$ for all eleven codes, giving a grid size of roughly $\delta z = 0.01$ for each code. A piecewise constant grid was a natural choice for some

photo-z codes, for instance most template-based codes compute likelihoods on a fixed grid. In contrast, FlexZBoost, for example, can return estimates on any grid without compression errors as its a basis expansion method where only the expansion coefficients need to be stored. Codes with a native output format other than the shared piecewise constant binning scheme (or one that can be losslessly converted to it) may suffer from loss of information when converting to it, which could artificially favor some codes over others.

Furthermore, the fidelity of photo-z interim posteriors in this format varies with the quality of the photometry. For faint galaxies, this redshift resolution is sufficient to capture the shape of $p(z)$ for the majority of the test sample, where photometric errors on the faint galaxies lead to somewhat broad peaks in the redshift posterior. However, as can be seen in e. g. the top left panel of Fig. 1, for bright galaxies with narrow $p(z)$ the grid spacing of $\delta z = 0.01$ is not sufficient to resolve the peak. This is consistent with the results described in Malz et al. (2018), who find that quantiles (and, to a lesser degree, samples) often outperform gridded $p(z)$, particularly for bright objects and in the presence of harsher storage constraints. With a full 200 numbers to capture the information of each photo-z PDF, any parametrization will perform adequately, but other storage parametrizations and limits on storage resources may be considered in future work. We will discuss this further in Section 6.

Fig. 2 shows both the quantile-quantile plots (red) and the histogram of PIT values (blue) summarizing the results from each photo-z code. The red line shows the measured quantiles, while the black diagonal represents the ideal QQ values if the distribution were perfectly reproduced. A second panel below the main panel for each code shows the difference between Q_{data} and Q_{theory} , i. e. the departure from the diagonal, for clarity. Biases and trends in whether the average width of the $p(z)$ values being over/under-predicted are evident. An overall bias where the predicted redshift is systematically low manifests as the measured QQ value falling above the diagonal, as is the case for BPZ and EAZY, while a systematic overprediction shows up as the measured QQ value falling below the diagonal, as seen in TPZ. In terms of PIT histograms, a systematic underprediction of redshift corresponds to fewer PIT values at $PIT < 0.5$ and more at $PIT > 0.5$, while a systematic overprediction will show the opposite.

Examination of the PIT histograms and QQ plots shows that there are fairly generic issues with the width of $p(z)$ uncertainties: DELIGHT, NN, SKYNET and TPZ all show a PIT histogram with an dearth of low values and an excess of high values, signs that, on average, their $p(z)$ are more broad than the true distribution of redshifts. METAPHOR shows the opposite trend, indicating the $p(z)$ are more narrow than the distributions given by the true redshifts. In all of these code cases there is a free parameter or bandwidth that can be used to tune uncertainties. The sensitivity of multiple codes to this bandwidth choice emphasizes the fact that great care must be taken in setting user-defined parameters in photo-z codes, even in the presence of representative training/validation data. for FLEXZBOOST the “sharpening” parameter (described in Section 3.2.4) plays a key role in improving the results, resulting in a QQ plot that is very nearly diagonal. A similar sharpening procedure could be beneficial for several codes. Interestingly, the three

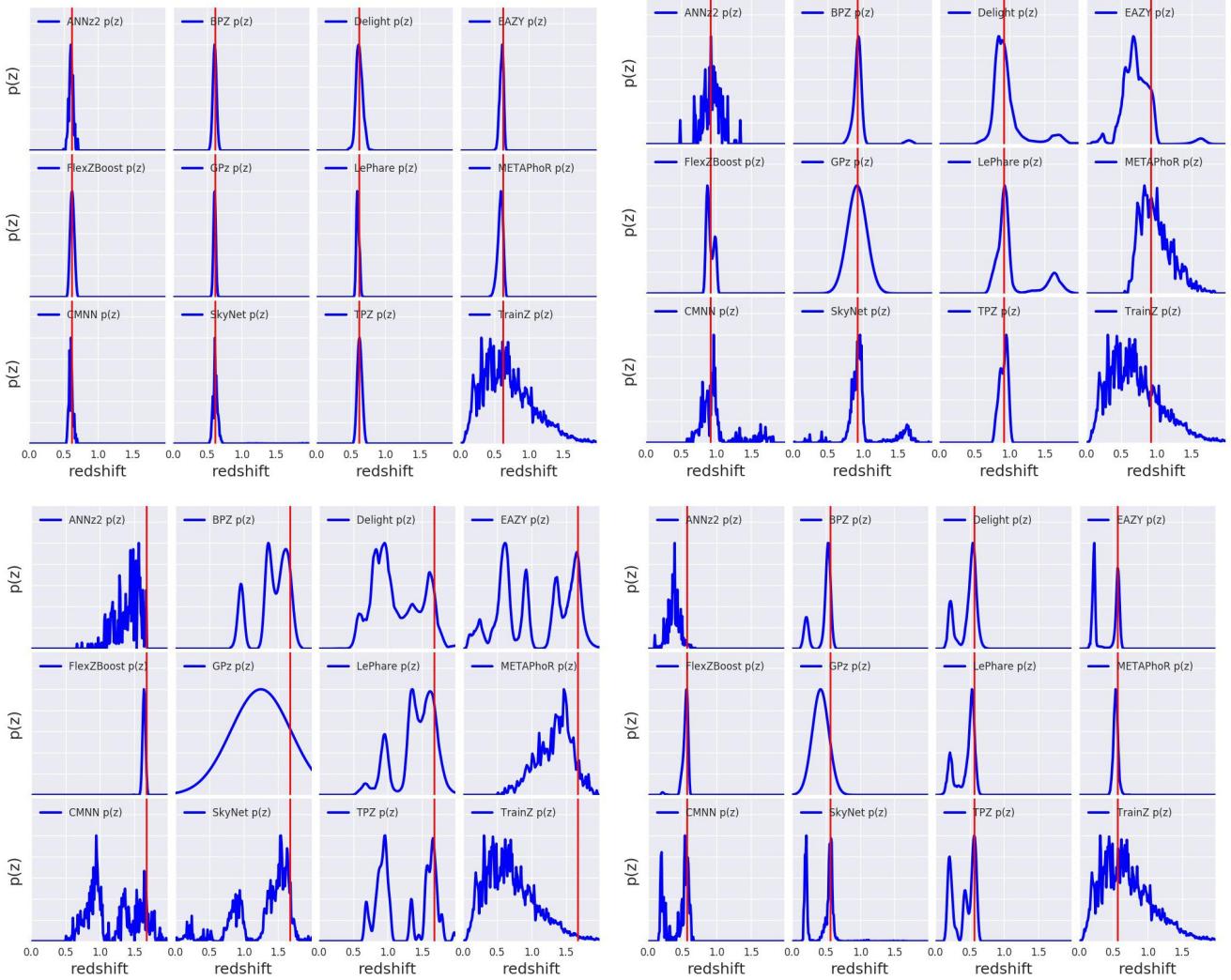


Figure 1. Four illustrative examples of individual $p(z)$ distributions produced by the codes. The red vertical line represents the true redshift. Examples are chosen with common features seen in PDFs: tight unimodal $p(z)$ (upper left), broad unimodal $p(z)$ (upper right), bimodal $p(z)$ (lower right), and complex/multimodal $p(z)$ (lower left). Codes show varying amounts of small-scale structure in their reconstruction of the posterior distribution. We see varying responses from the codes in the presence of color degeneracies and photometric errors, resulting in narrow and broad unimodal, bimodal, and multi-modal $p(z)$ curves.

purely template-based codes, BPZ, EAZY, and LEPHARE, show relatively well behaved $p(z)$ statistics (albeit with some bias), which may indicate that the likelihood estimation with representative templates is accurately capturing the uncertainties on individual redshifts.

The ideal PIT histogram would follow the black dashed line, representing a uniform distribution of PIT values, equivalent to the diagonal line in the QQ plot. Overly broad $p(z)$ values show up as an excess of PIT values near 0.5 and a dearth of values at the edges, while overly narrow $p(z)$ will have an excess at the edges and will be missing values at the centre. Another feature evident in the PIT histograms is the number of ‘catastrophic outlier’ values where the true redshift falls outside of the non-zero support of $p(z)$, corresponding to $\text{PIT} = 0.0$ or 1.0 is more apparent than in the QQ plots. Following Kodra & Newman (in prep.) we define f_O as the fraction of objects with $\text{PIT} < 0.0001$ or $\text{PIT} > 0.9999$. Table 2 lists these fractions for each of

the codes. For a proper Uniform distribution we expect a value of 0.0002. Several codes show a marked excess, with ANNz2, FLEXZBOOST, LEPHARE, AND METAPHOR with $f_O > 0.02$, indicating a sizeable number of catastrophic redshift solutions where the true redshift is not covered by the extent of $p(z)$. For METAPHOR this may be partially due to an overall underprediction of the $p(z)$ width, however this is not the case for the other codes. LEPHARE is a particular outlier with nearly 5 per cent of objects outside of $p(z)$ support. Further study will be necessary to determine what is causing these misclassifications for LEPHARE. As expected, and by design, TRAINZ has the proper fraction of outliers for the f_O statistic.

Fig. 3 shows comparative metric values for the quantitative Kolmogorov-Smirnoff (KS), Cramer-Von Mises (CvM), and Anderson Darling (AD) test statistics for each of the codes based on comparing the distribution of their PIT values to the expected uniform distribution over the interval

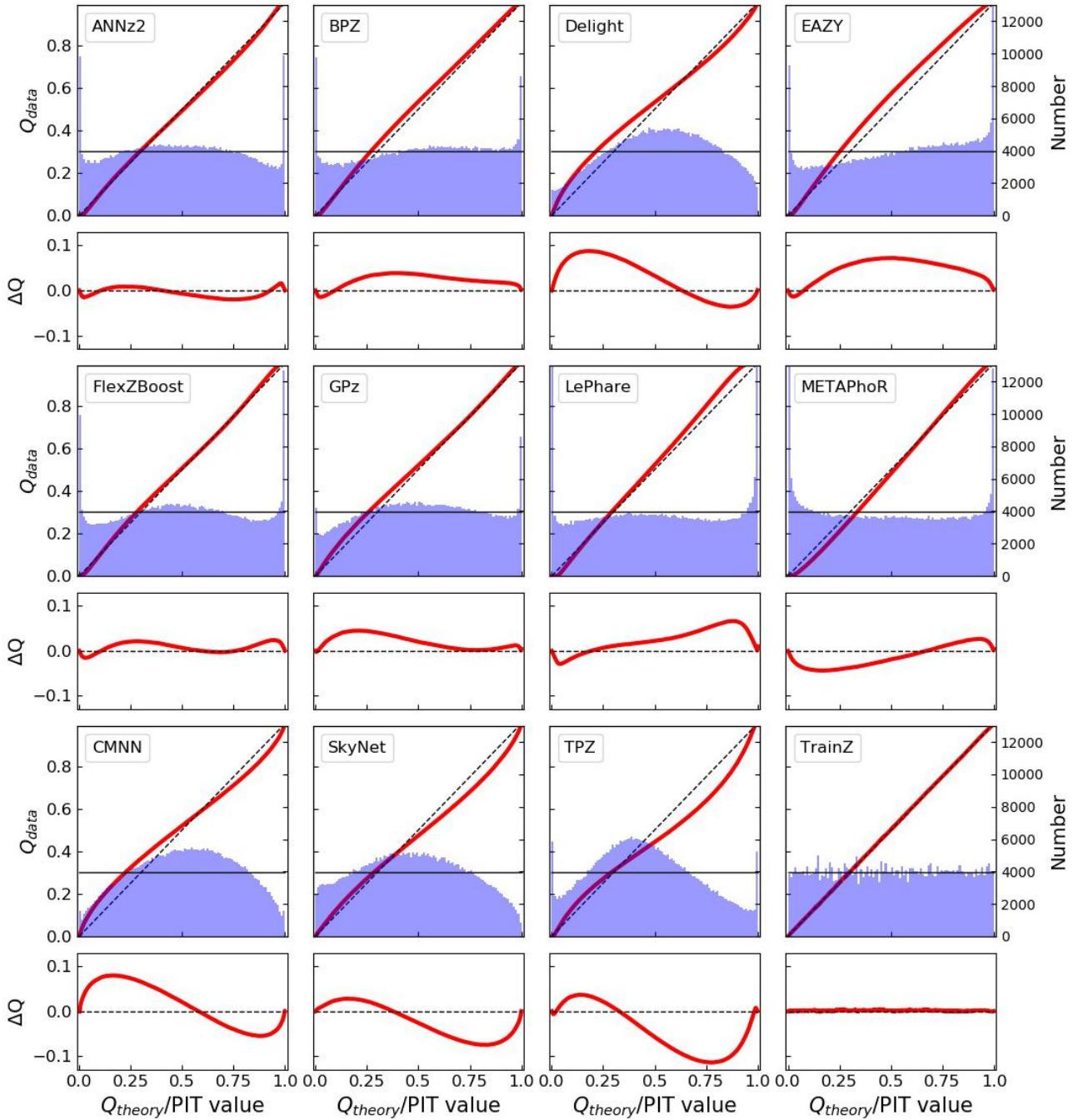


Figure 2. Summary plots for all eleven photo-z codes illustrating performance for the interim posterior statistics. The top panel of each pair shows both the Quantile-Quantile (QQ) plot (red) and the histogram of PIT values (blue). The desired behavior is a QQ plot that matches the diagonal dashed line, and a PIT histogram that matches a uniform distribution matching the thin horizontal black line. The bottom panel of each pair shows the difference between the QQ quantile and the diagonal, illustrating departure from the desired performance. Histograms with an overabundance of PIT values at the centre of the distribution indicate $p(z)$ distributions that are overly broad, while an excess of values at the extrema indicate $p(z)$ distributions that are overly narrow. Values of PIT=0 and PIT=1 indicate “catastrophic failures” where the true redshift is completely outside the support of $p(z)$. Asymmetric features are indicative of systematic bias in the redshift predictions. A variety of behaviors are evident, and specific details are discussed in the text.

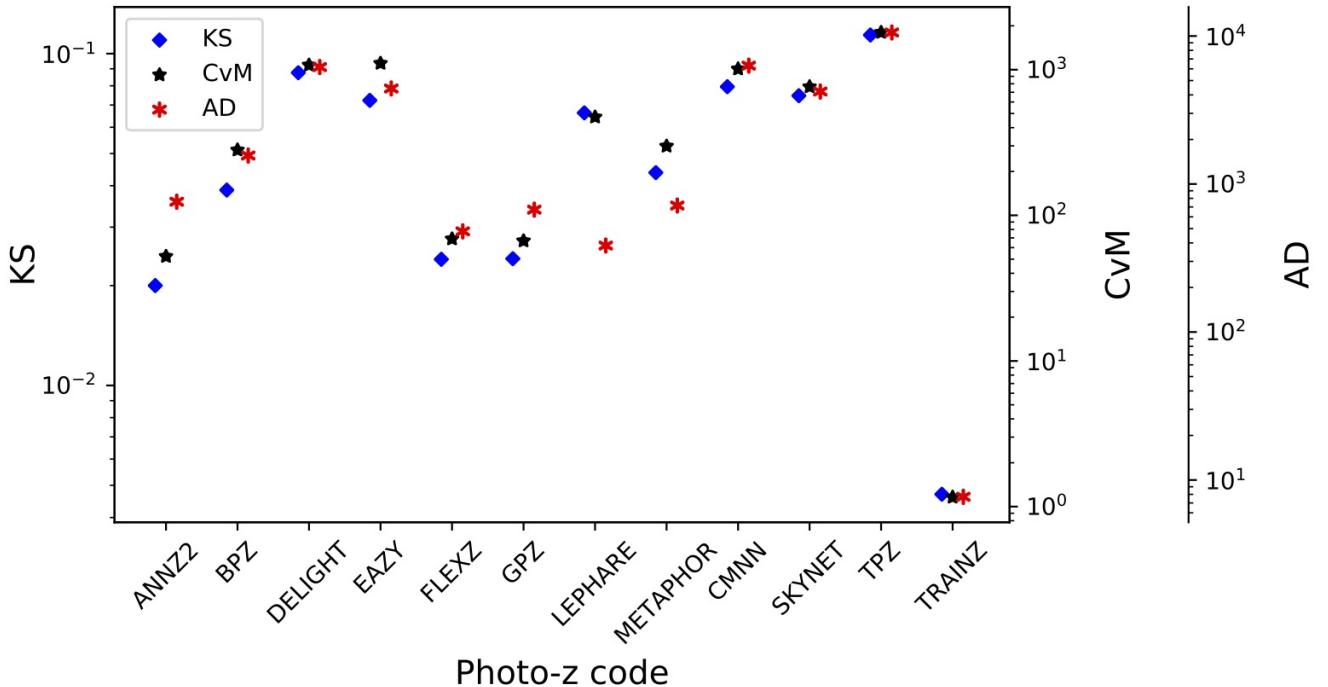


Figure 3. A visual representation of the Kolmogorov-Smirnov (KS, blue diamond), Cramer-von Mises (CvM, black star), and Anderson-Darling (AD, red asterisk) statistics for the PIT distributions. The statistics are often highly correlated, though the AD statistic truncates the extrema of the distribution and can have disparate values compared to KS and CvM.

Table 2. The fraction of “catastrophic outlier” PIT values. We expect a value of 0.0002 for a proper Uniform distribution. An excess over this small value indicates true redshifts that fall outside the non-zero support of the $p(z)$.

Photo-z Code	fraction $\text{PIT} < 10^{-4}$ or > 0.9999
ANNz2	0.0265
BPZ	0.0192
DELIGHT	0.0006
EAZY	0.0154
FLEXZBOOST	0.0202
GPZ	0.0058
LEPHARE	0.0486
METAPHOR	0.0229
CMNN	0.0034
SKYNET	0.0001
TPZ	0.0130
TRAINZ	0.0002

a large number of catastrophic outliers, resulting in higher KS and CvM scores.

Given the near-perfect training data, examining the individual codes for explanations for departures from the expected behaviour will be instructive in avoiding similar problems in future tests. ANNz2 performs quite well in $p(z)$ based metrics. In the specific implementation employed in this paper, the final $p(z)$ is a weighted average of five neural-nets. During the training process ANNz2 compares the percentiles of the redshift training sample against the CDFs of the $p(z)$ sample. Distributions that more closely match are given extra weight, and the final weights are designed to produce accurate percentiles. Given that our metrics are focused on the percentile distributions, it is unsurprising that ANNz2 performs well in the given metrics. The discreteness in the individual $p(z)$ estimated by ANNz2 can be attributed to the fact that the code was run as a classifier, assigning weights to discrete bins of redshift. While multiple bins may receive weight, the bins themselves will still be discretized, and no additional smoothing was performed. Overall, FLEXZBOOST and ANNz2 show the best ensemble agreement in their distribution of PIT values.

[0,1]. The individual values of the statistic are not as important as the comparative score between the different codes. The AD test statistic diverges for values that include the extrema, and thus is calculated by excluding the edges of the distribution. We calculate the AD statistic over the range of PIT values $v = [0.01, 0.99]$. ANNz2 and FLEXZBOOST score very well for the PIT metrics. METAPHOR and LEPHARE score very well in the PIT AD statistic, but both have

5.2 Metrics of the stacked estimator of the redshift distribution

Fig. 4 shows the stacked $\hat{N}(z)$ distribution compared to the true redshift distribution $N'(z)$ for all tested codes. The red line indicates the summed $p(z)$ for each code, while the blue line shows the true redshift distribution smoothed via kernel density estimation (KDE), with a bandwidth chosen via

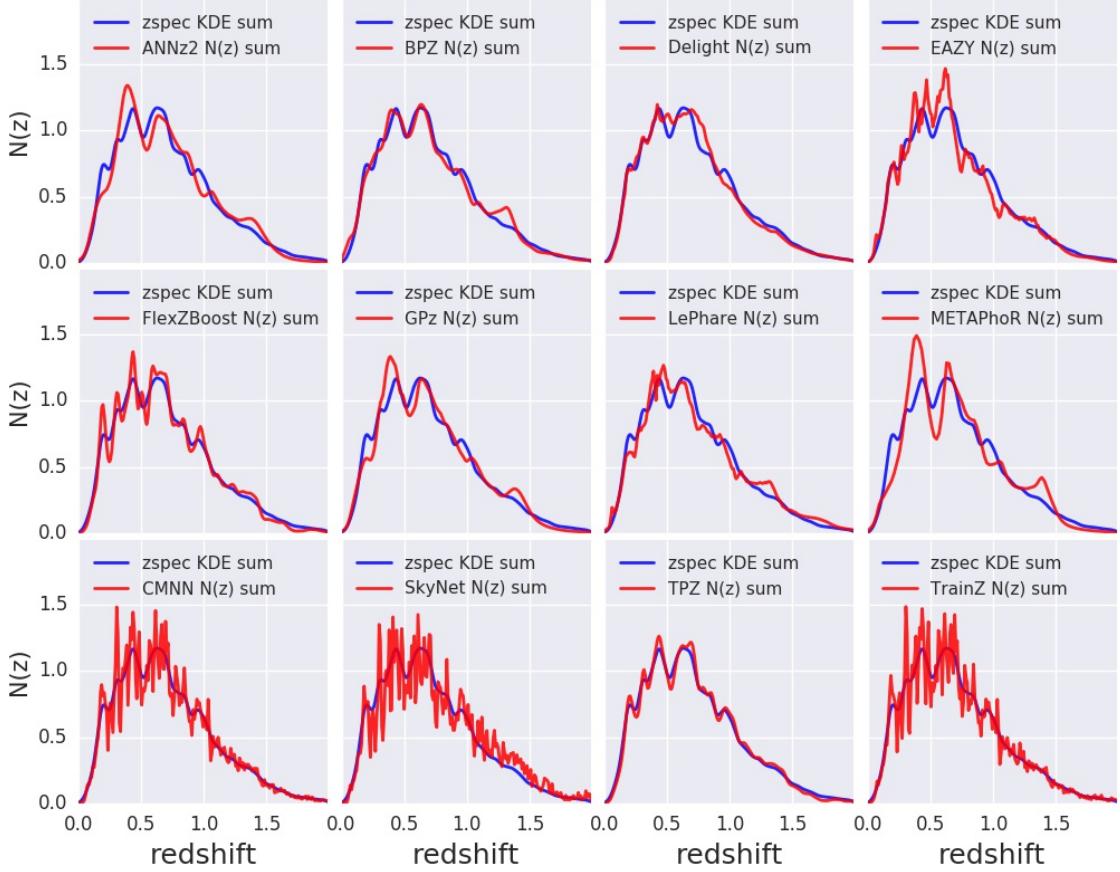


Figure 4. The stacked $p(z)$ produced by each photo-z code ($\hat{N}(z)$, red) compared to the spectroscopic redshift distribution ($N'(z)$, blue). Varying levels of small-scale structure are seen in the codes. $N'(z)$ is smoothed using a single bandwidth chosen via Scott's rule for all codes.

Scott's rule (Scott 1992). While Scott's rule is used to display $N'(z)$ in the figure, all quantitative statistics are computed via the empirical CDF, and are thus unaffected by bandwidth/smoothing choice. Several of the codes show an excess at $z \sim 1.4$, particularly the template-based codes BPZ, EAZY, and LEPHARE. This is likely due to the 4000 angstrom break passing through the gap between the z and y filters. This feature is one of the most prominent in individual galaxy $p(z)$, and is readily seen in the point-estimate plots shown in Fig. A1 and described in the Appendix. Several of the machine learning based codes appear to be overtrained, adding excess galaxy probability to the redshift peaks and missing probability in the troughs. Given that our training data is drawn from the same galaxy population as the test set, and our data has prominent peaks in $N'(z)$, perhaps it is not unexpected that such overtraining occurs. A more extensive training/validation set might allow for a better choice of smoothing parameters in individual codes that would avoid such overtraining.

As with the $p(z)$ values in Figure 2, different levels of substructure are obvious for the different codes. While

Scott's rule provides a relatively good general smoothing scale to represent the true $N'(z)$, there are smaller scale fluctuations: while FLEXZBOOST and CMNN appear somewhat discrepant in Fig. 4, they are actually the two most accurate in terms of their quantitative measurements. Interestingly, while ANNz2 shows an abundance of small scale structure in individual $p(z)$ measurements (see Fig. 1), the summed $\hat{N}(z)$ is rather smooth, where the small scale features average out. This is not the case for the two other codes that show an abundance of substructure in their individual $p(z)$: both CMNN and SKYNET show small scale features both in $p(z)$ and $\hat{N}(z)$. For CMNN the $p(z)$ are simply a weighted histogram of all spectroscopic training galaxies in nearby colour space with no smoothing applied, so the substructure is not unexpected. The PIT histogram and shape of the QQ plot in Figure 2 show that CMNN is producing $p(z)$ that are overly broad, additional smoothing of the $p(z)$ would exacerbate this problem. While the $\hat{N}(z)$ plot shows more small scale features than other codes, these features are actually representative of real structure in the true $N'(z)$, as evidenced by the very good metric scores for

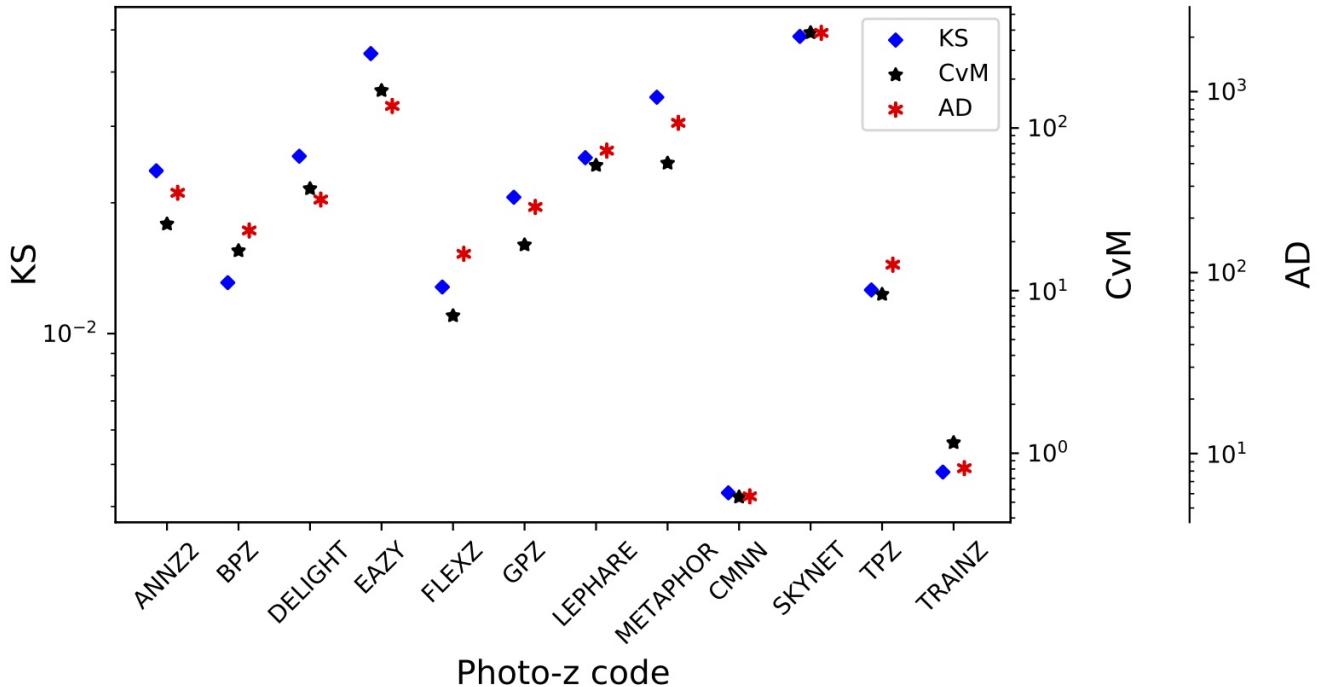


Figure 5. A visual representation of the Kolmogorov-Smirnoff (KS, blue diamond), Cramer-von Mises (CvM, black star), and Anderson-Darling (AD, red asterisk) statistics for the $\hat{N}(z)$ distributions. The statistics are correlated, the codes with the lowest KS statistics tend to have the lowest CvM and AD statistics. CMNN performs markedly better than the others in reconstructing the overall $N(z)$ distribution, while SKYNET scores poorly due to an overall bias in its redshift predictions.

CMNN. SKYNET $p(z)$ were also not smoothed: while previous implementations of the code such as Sánchez et al. (2014) and Bonnett (2015) (see Appendix C.3) implement a “sliding bin” smoothing, no such procedure was used in this study. In addition to excess substructure, SKYNET shows an obvious redshift bias, evident both visually in Figure 4 and in the first moment of $N(z)$ listed in Table 5, where it is clearly an outlier. SKYNET employed a method where a random sample of training galaxies was chosen, but there was no test that the subset was completely representative of the overall redshift distribution. Also unlike Bonnett (2015), no effort was made to add extra weight to more rare low and high redshift galaxies. Either of these decisions could be the cause of the bias seen in our results. Future runs of SKYNET will explore these implementation choices and their effects.

Figure 5 shows the quantitative Kolmogorov-Smirnoff (KS), Cramer-Von Mises (CvM), and Anderson Darling (AD) test statistics for each of the codes for the $\hat{N}(z)$ based measures. FLEXZBOOST, CMNN, and TPZ outperform the other codes in the $\hat{N}(z)$ metrics. It is unsurprising that CMNN scores well, as with a near perfectly representative training set means that choosing neighbouring points in color/magnitude space should lead to excellent agreement in the final $\hat{N}(z)$ estimate. TPZ performed quite poorly in $p(z)$ statistics, but results in a good fit to the overall $N(z)$. This is somewhat surprising, as performance was optimized for accurate $p(z)$, not $\hat{N}(z)$. During the validation stage for TPZ, there was a trade off between the width of the $p(z)$ when adjusting a smoothing parameter and overall redshift bias. The optimal result in the PIT metrics, as illustrated in

the shape of the QQ plot, does contain some level of bias as well as a slight underprediction of mean $p(z)$ width, which translates to poor metric scores. This is something that will be looked into for TPZ in the future.

It is also of note that all three template-based codes show an excess in their stacked $p(z)$ at $z \sim 1.3 - 1.4$. This redshift range corresponds to the wavelengths where the 4000 Angstrom break is passing between the borders of the z and y filters. This strong break entering the gap between the two reddest filters can cause problems with redshift estimation of individual galaxies, as can be seen in the point-estimate photo- z ’s shown in Figure A1. This is not unique to this dataset, it is a common occurrence in photo- z estimation. The fact that similar excesses appear in Figure 4 for ANNz2 and METAPHOR shows that the effect is not limited to template-based codes. However, the lack of such a feature in the other codes shows that it is possible to eliminate the degeneracies. Further study on this issue may provide a solution for codes that suffer from this shortcoming.

Table 3 shows the CDE loss statistic for each photo- z code. Once again FLEXZBOOST and CMNN score very well for the stacked $\hat{N}(z)$ metrics, as do GPZ and TPZ. The CDE loss measures how well individual PDFs are estimated, and codes with a low CDE loss tend to have good $\hat{N}(z)$ estimates (though the reverse is not necessarily true). FLEXZBOOST is optimized to minimize CDE loss which may explain why the method has good ensemble metrics as well. Note from Table 3 that both FLEXZBOOST and CMNN have low CDE losses. Empirically, we have found that PIT RMSE is not as

1321 closely correlated to CDE loss as it is to the $N(z)$ statistics.
 1322 As CDE loss is a better measure of individual redshift performance,
 1323 rather than ensemble distribution performance, this statistic is a
 1324 better indicator of which codes will be most likely to perform well for
 1325 science cases where single objects are employed.

1326 Table 4 gives the root-mean-square-error (RMSE) statistics for both the PIT and $N(z)$ estimators. The PIT value calculates the RMSE between the quantiles shown in the QQ plot in Figure 2 and the diagonal, while the $N(z)$ calculates the RMSE between the cumulative distribution of the stacked $\hat{N}(z)$ and the true redshift distribution $N'(z)$.

1327 Table 5 lists the first three moments of the stacked $\hat{N}(z)$ distribution, including the moments of the “truth” distribution for comparison. Several codes are able to reproduce the mean and variance of the distribution to less than a per cent, while several codes do not, which may be a cause for concern, given that mean and variance of the redshift distribution are key properties in cosmological analyses. We note that this stated goal of the study as defined for participants was to accurately reproduce $p(z)$, the “stacking” of the probability distributions to estimate $\hat{N}(z)$ was not the focus as stated to the participants. This explains why some of the best-performing empirical codes in terms of $p(z)$ measures (e. g. FLEXZBOOST) do not do as well at reproducing $\hat{N}(z)$ moments. Had we defined a different parameter to optimize, in this case overall accuracy of $\hat{N}(z)$ rather than individual $p(z)$, would result in improved performance in a particular metric. That is, optimizing photo-z performance for one metric does not automatically give optimal performance for other metrics. As previously stated, there are a variety of scientific use cases for photo-z’s in large upcoming surveys, and care must be taken in how the metrics used to optimize catalog photometric redshifts are defined as well as in how they are used. In addition, very few scientific use cases will employ the overall $\hat{N}(z)$ with no cuts, as we explore in this paper. We discuss more realistic tomographic bin selections that will be explored in a follow-up paper in Section 6.1.

1359 5.3 Interpretation of metrics

1360 Samples from accurate photo-z posteriors should reproduce 1385 the space of $p(z, data)$. However, it is difficult to test this re- 1386 construction given our data set, as the galaxy distributions 1387 arise from mock objects pasted on to an underlying dark 1388 matter halo catalogue with properties designed to match 1389 empirical relations, rather than being drawn from statisti- 1390 cal distributions in redshift. In previous sections we have 1391 mentioned that optimizing for a specific metric does not 1392 guarantee good performance on other metrics, nor is there 1393 any guarantee that good performance by our metrics cor- 1394 responds to accurate photo-z posteriors. In other words, we 1395 can construct photo-z estimators that provide good coverage 1396 in many of our tests, but which have very little predictive 1397 power.

1374 The TRAINZ estimator, which assigns every galaxy a 1399 $p(z)$ equal to $N(z)$ of the training set as described in Sec- 1400 tion 3.3, is introduced as a “null test” to demonstrate this 1401 point via *reductio ad absurdum*. TRAINZ outperforms all 1402 codes on the PIT-based metrics, and all but one code on 1403 the $N(z)$ based statistics. Because our training set is per- 1404

Table 3. CDE loss statistic for each photo-z code.

Photo-z Code	CDE Loss
ANNZ2	-6.88
BPZ	-7.82
DELIGHT	-8.33
EAZY	-7.07
FLEXZBOOST	-10.60
GPz	-9.93
LEPHARE	-1.66
METAPHoR	-6.28
CMNN	-10.43
SKYNET	-7.89
TPZ	-9.55
TRAINZ	-0.83

Table 4. Root-Mean-Square-Error (RMSE) statistics for the eleven photo-z codes for both PIT and $\hat{N}(z)$ distributions.

Photo-z Code	PIT RMSE	$N(z)$ RMSE
ANNZ2	0.019	0.0054
BPZ	0.032	0.0050
DELIGHT	0.111	0.0056
EAZY	0.054	0.0102
FLEXZBOOST	0.021	0.0022
GPz	0.027	0.0042
LEPHARE	0.028	0.0062
METAPHoR	0.064	0.0081
CMNN	0.108	0.0009
SKYNET	0.054	0.0144
TPZ	0.082	0.0031
TRAINZ	0.0025	0.0013

flectly representative of the test set, $N(z)$ should be identical for both sets down to statistical noise.

The CDE loss and point estimate metrics, however, successfully identify problems with TRAINZ. As shown in Appendix A, TRAINZ has identical *ZPEAK* and *ZWEIGHT* values for every galaxy, and thus the photo-zs are constant as a function of spec-zs, i.e. a horizontal line at the mode and mean of the training set distribution respectively. The explicit dependence on the *individual posteriors in the calculation of the CDE loss, described in Section 4.1.3, distinguishes this metric from the other p(z) metrics that test the overall ensemble of p(z) distributions. With a representative training set, TRAINZ will score well on the ensemble metrics, but fails miserably for metrics tied to individual redshifts. We note that many of the ensemble-based metrics are prominent in the photo-z literature despite their inability to identify problems such as those exemplified by TRAINZ.*

In summary, context is crucial to interpreting metrics and defending against the likes of TRAINZ. The best photo-z method is the one that most effectively achieves our science goals, not the one that performs best on a metric that does not accurately reflect those goals. In the absence of clear goals or the information necessary for a principled metric definition, we must think carefully before choosing a single metric

Table 5. Moments of the stacked $\hat{N}(z)$ distribution

Stacked $n(z)$ Moments			
	1st Moment	2nd Moment	3rd Moment
TRUTH	0.701	0.630	0.671
Photo-z Code	1st Moment	2nd Moment	3rd Moment
ANNz2	0.702	0.625	0.653
BPZ	0.699	0.629	0.671
DELIGHT	0.692	0.609	0.638
EAZY	0.681	0.595	0.619
FLEXZBOOST	0.694	0.610	0.631
GPz	0.696	0.615	0.639
LEPHARE	0.718	0.668	0.741
METAPHO.R	0.705	0.628	0.657
CMNN	0.701	0.628	0.667
SKYNET	0.743	0.708	0.797
TPZ	0.700	0.619	0.643
TRAINZ	0.699	0.627	0.666

6 DISCUSSION

In this paper we presented results evaluating the photometric redshift PDF computation for eleven photo-z codes. As discussed in Section 4 the $p(z)$ should accurately reflect the relative likelihood as a function of redshift for each galaxy. All codes were provided a set of representative training data and tested on an idealized set of model galaxies with high signal-to-noise and photometry with no confounding effects due to blending, instrumental effects, the night sky, etc... included. The goal was not to determine a “best” photo-z code: in many ways, this was a baseline test of a “best case scenario” to predict the expected photo-z performance if a stage IV dark energy survey was to obtain complete training samples and perfectly calibrated their multi-band photometry. Given these idealized conditions, any deficiencies observed in a photo-z code’s performance should be a cause for concern, and may be evidence in a problem with either/both of the specific code implementation or the underlying algorithm. In order to meet the stringent LSST requirements on photo-z performance, identifying and correcting such problems is an important first step before tackling more realistic data in future challenges. Most of the codes tested performed well, however, several did not meet the stringent goals that have been laid out for LSST photometric redshift performance. This is a cause for concern, given the idealized conditions, and the individual code responses will be studied in detail moving forward. One obvious trend in several of the codes tested was an overall over or underprediction of the widths of $p(z)$, as evidenced by the QQ plots and PIT histograms shown in Fig. 2. A more careful tuning of bandwidth or smoothing during the validation process appears to be necessary for many of the machine learning based codes in order to improve the accuracy of $p(z)$. For narrow peaked $p(z)$ the parameterization of the PDF as evaluated on a fixed redshift grid could also have contributed to some overestimates of $p(z)$ width simply due to the finite resolution. After evaluating results such as those presented in Malz et al. (2018), in future analyses we plan to switch from a fixed grid to quantile-based storage of $p(z)$ in order to more efficiently and accurately store redshift PDF results.

Another important factor to keep in mind when examining the results presented in this paper is the fact that they are at some level dependent on the metrics that we aim to optimize: in this case code participants were asked to submit their optimal measures of an accurate $p(z)$, so participants used the training/validation data to optimize their codes accordingly. Had we, instead, asked for an optimal $N(z)$ the resulting metrics would be different for most, if not all, of the codes, as they would optimize toward a different goal. Specific metric choice can affect which codes are among the “best” codes. As stated earlier, there are cosmological science cases that require either individual galaxy photo-z measures, or ensemble $\hat{N}(z)$ measures. We must be aware of that the optimal method for one is not necessarily optimal for the other, and in fact several photo-z algorithms may be necessary in the final cosmological analysis in order to satisfy the requirements of all science use cases. The example of the simple TRAINZ estimator described in Section 5.3 shows a simple model with a $p(z)$ that is unrealistic for individual objects can still score very well on many of our metrics. It is important to look at all metrics, and keep in mind what information each metric conveys. We re-emphasize that the dataset tested was quite idealized, and discuss enhancements that will be added in future simulations to test photo-z codes on increasingly realistic conditions in the following section.

6.1 Future work

The work presented in this paper is only the first step in characterizing current photo-z codes and moving toward an improved photometric redshift estimator. This initial paper explored code performance in idealized conditions with perfect catalog-based photometry and representative training data. As mentioned in Section 5.2 for the stacked $N(z)$ metrics we examined only the entire galaxy population with no selections in either photo-z “quality” or redshift. The cosmological analyses for weak lensing and large scale structure based measures plan to break galaxy samples into tomographic redshift bins, using photo-z $p(z)$ to infer the redshift distribution for each bin. The specific selection used to determine these bins, both algorithmically and the specific bin boundaries, could induce biases due to indirect selections inherent in the photo-z or other bin selection parameters. The effects of tomographic bin selection will be explored in a dedicated future paper. [are there any references for this? I remember Gary Bernstein talking about this at a photo-z workshop in Japan, but I don’t know that it was published. I believe Michael Troxel has discussed this as well.] We also plan to propagate the uncertainties measured in a set of fiducial tomographic redshift bins in order to estimate impact on cosmological parameter estimation.

In future papers we will add more and more complexity to our simulated data in order to test photo-z algorithms in increasingly realistic conditions. The most pressing concern is the impact of incomplete spectroscopic training samples. As discussed extensively in Newman et al. (2015) a representative set of spectroscopically confirmed galaxies spanning the full range of both redshift and apparent magnitude is necessary as a training set to characterize the mapping from broad-band fluxes to photometric redshifts. However, due to a combination of factors due to both the galaxy SEDs and lim-

iterations of spectrographic instruments, redshift samples are known to be systematically incomplete, where certain galaxy types and redshift intervals fail to yield a redshift even at the longest integration times on current and near-future instruments. The more representative the training data, the better the performance of photo-z algorithms will be. Current and upcoming surveys are putting in significant effort into obtaining these training samples (e.g. [Masters et al. 2017](#)), however we still expect significant incompleteness for LSST-like samples, particularly at faint magnitudes. One major focus of an upcoming LSST Dark Energy Science Collaboration Photo-z Working Group data challenge is to produce a realistically incomplete training set of spectroscopic galaxies, modeling the performance of spectrographs, emission-line properties, and expected signal-to-noise to determine which galaxies will fail to yield a secure redshift. In addition to outright redshift failures we will model the inclusion of a small number of falsely identified secure redshifts where misidentified emission lines or noise spikes cause an incorrect redshift solution to be marked as a high quality identification. Even sub-per cent level contamination by false redshifts can impact photo-z solutions at levels comparable to the stringent requirements of some LSST science cases. We expect different systematics to occur in different photo-z codes in response to training on incomplete data, particularly some of the machine learning methods. The response of the codes will inform future directions of code development.

This initial paper explored a data set that was constructed at the catalog level, with no inclusion of the complications that come from measuring photometry from images. Future data challenges will move to catalogs constructed from mock images, including effects that will have great impact on photo-z measurements. Object blending will be a major area of investigation, as the mixing of flux from multiple objects and the resultant change in measured colours is predicted to affect a large fraction of LSST galaxies ([Dawson et al. 2016](#)), and will be one of the major contributing systematics for photo-z's. Inclusion of differing observing conditions (seeing, clouds, variations in filter curves, Galactic dust, ...), as well as models for instrumental and system effects, sky masks, will all impact object photometry, and will be explored in the upcoming data challenge and their impacts described in upcoming papers. All underlying SEDs were parameterized as a weighted combination of five basis SEDs, with no additional accounting for host galaxy dust obscuration beyond what was encoded in the basis templates. This, in effect, limited the simulation to a very simple model of internal obscuration. Future simulations will include a more complicated and realistic treatment of host galaxy dust.

The underlying simulation used in this work was based on a light-cone constructed to a maximum redshift of $z = 2$. LSST imaging after 10 years of observations will include a significant number of $z > 2$ galaxies in expected cosmology samples, and their inclusion does have potential significant implications for photo-z measures: the high redshift galaxies lie at fainter apparent magnitudes and can have anomalous colours due to evolution of stellar populations and the shift to rest-frame magnitudes probing UV features of the underlying SED. More importantly, one of the most common “catastrophic outlier” degeneracies observed in deep photometric samples occurs when the Lyman break is mistaken for the Balmer break, leading to multiple redshift solutions

at $z \sim 0.2 - 0.3$ and $z \sim 2 - 3$ ([Massarotti et al. 2001](#)). This degeneracy, along with other potential degeneracies, are currently not covered by the limited redshift range of this initial paper, which could mean that we are not probing the full range of potential extreme outlier populations and how our photo-z estimators respond to them. Extending simulations to include the high-redshift galaxy population will be a priority in future data challenges.

7 CONCLUSION

In this study we have not accounted for the presence of Active Galactic Nuclei (AGN) contributions to galaxy fluxes. In some cases, AGN will be easily identified from the colors and morphologies, i.e. the case of the brightest quasars where the nuclear activity outshines the host galaxy, and numerous studies have utilized color selection to create large samples of quasars (e.g. [Richards et al. 2006](#); [Maddox et al. 2008](#); [Richards et al. 2015](#)). In current deep fields, similar in depth to what we expect from LSST, variability information and multi-wavelength data have been critical to not only identify AGN dominated galaxies, but also obtain more accurate photometric redshifts (e.g. [Salvato et al. 2011](#)).

In addition to AGN dominated galaxies, those with lower levels of nuclear activity present a more insidious problem, where AGN features may not be apparent, but the colors and other host galaxy properties are perturbed relative to galaxies with an inactive nucleus. In such cases, the presence of the AGN may induce a bias if the template SEDs or empirical datasets do not include low-level AGN counterparts. For LSST, we will need to identify and obtain accurate photometric redshifts of all types of AGN for a range of science goals, whether it is to eliminate such objects from cosmology experiments, or to use them with confidence, all the way through to understanding galaxy evolution and the role that AGN may play in influencing galaxy properties over cosmic time.

A promising route to classifying and obtaining accurate photometric redshifts for the AGN population is by combining machine learning with template-fitting techniques, as has recently been demonstrated by [Duncan et al. \(2018\)](#) for radio-selected AGN. This is because AGN are relatively easy to obtain spectroscopic redshifts for over all redshifts due to the strong emission lines that they exhibit, allowing very good training sets for machine learning algorithms to use. Whereas for those galaxies where the AGN is sub-dominant the galaxy templates are still adequate for obtaining reasonable photometric redshifts.

In addition to these improvements, the DESC Photo-z group plans to look at all potential methods to combine the results from multiple photo-z codes to improve $p(z)$ accuracy, similar to the work presented in [Dahlen et al. \(2013\)](#); [Carrasco Kind & Brunner \(2014\)](#); [Duncan et al. \(2018\)](#). Taking advantage of multiple algorithms that use observables in slightly different ways has shown promise, however we must be very conscious of whether a potential combination properly treats the covariance between the methods, given that they are estimating quantities based on the same underlying observables. Several science cases wish to estimate physical quantities along with redshift, for example galaxy stellar mass and star formation rate. Proper joint estimation of

1625	redshift and physical quantities requires an in depth under-	1685	P.E. Freeman: Contributed to choice of CDE metrics and
1626	standing of galaxy evolution, and progress on accurate bivari-	1686	to implementation of FlexZBoost
1627	ate redshift probability distributions will go hand in hand with	1687	K. Iyer: assisted in writing metric functions used to evaluate
1628	progress on understanding galaxies themselves. Parameter-	1688	codes
1629	ization and storage of a complex 2-dimensional probability	1689	J.B. Kalmbach: Worked on preparing the figures for the
1630	surface for potentially billions of galaxies (or even subsets	1690	paper.
1631	of hundreds of thousands of particular interest) pose a po-	1691	E. Kovacs: Ran simulations, discussed data format and
1632	tential challenge. These issues will be examined in another	1692	properties for SEDs, dust, and ELG corrections
1633	future paper.	1693	A.B. Lee: Co-developed FlexZBoost and the CDE loss statis-
1634	tic, wrote text on the work, and supervised the development	1694	of FlexZBoost software packages
1635	Finally, while this paper and future papers discussed	1694	C. Morrison: Managerial support; Discussions with authors
1636	above focus on photometric redshift codes and estimating	1695	regarding metrics and style; Some coding contribution to
1637	accurate $p(z)$ from training data, we plan a separate, but	1696	metric computation.
1638	complementary, project to examine calibration of the resul-	1697	J. Newman: Contributions to overall strategy, design of
1639	tant redshifts via spatial cross-correlations (Newman 2008),	1698	metrics, and supervision of work done by Rongpu Zhou
1640	which will be explored in a separate series of future papers.	1699	E. Nuss: contributed to running code, analysis discussion,
1641	The overarching plan describing everything laid out in this	1700	and editing, reviewing the paper
1642	section is described in more detail in the LSST DESC Sci-	1701	T. Pospisil: Co-developed FlexZBoost software and CDE
1643	ence Roadmap (see Footnote in Section 1). These plans will	1702	loss calculation code
1644	require significant effort, but they are necessary if we are to	1703	M.J. Jarvis: Contributed text on AGN to Discussion section
1645	make optimal use of the LSST data for astrophysical and	1704	and portions of GPz work
	cosmological analyses.	1705	R. Izbicki: Co-developed FlexZBoost and the CDE loss
		1706	statistic, and wrote software for FlexZBoost
		1707	
		1708	
1646	Acknowledgments		
1647	Author contributions are listed below.	1709	The authors would like to thank their LSST-DESC pub-
1648	S.J. Schmidt: Led the project. (conceptualization, data	1710	lication review committee.
1649	curation, formal analysis, investigation, methodology,	1711	AIM is advised by David W. Hogg and was supported by
1650	project administration, resources, software, supervision,	1712	National Science Foundation grant AST-1517237.
1651	visualization, writing – original draft, writing – review &	1713	The DESC acknowledges ongoing support from the In-
1652	editing)	1714	stitut National de Physique Nucléaire et de Physique des
1653	A.I. Malz: Contributed to choice of metrics, implementation	1715	Particules in France; the Science & Technology Facilities
1654	in code, and writing. (conceptualization, methodology,	1716	Council in the United Kingdom; and the Department of En-
1655	project administration, resources, software, visualization,	1717	ergy, the National Science Foundation, and the LSST Cor-
1656	writing – original draft, writing – review & editing)	1718	poration in the United States. DESC uses resources of the
1657	J.Y.H. Soo: Ran ANNz2 and Delight, updated abstract,	1719	IN2P3 Computing Center (CC-IN2P3-Lyon/Villeurbanne -
1658	edited sections 1 through 6, added tables in Methods	1720	France) funded by the Centre National de la Recherche Sci-
1659	and Results, updated references.bib and added references	1721	entifique; the National Energy Research Scientific Comput-
1660	throughout the paper	1722	ing Center, a DOE Office of Science User Facility supported
1661	M. Brescia: main ideator of METAPHOR and of MLPQNA;	1723	by the Office of Science of the U.S. Department of Energy
1662	modification of METAPHOR pipeline to fit the LSST data	1724	under Contract No. DE-AC02-05CH11231; STFC DiRAC
1663	structure and requirements	1725	HPC Facilities, funded by UK BIS National E-infrastructure
1664	S. Cavaudi: Contributed to choice and test of metrics, ran	1726	capital grants; and the UK particle physics grid, supported
1665	METAPHOR, minor text editing	1727	by the GridPP Collaboration. This work was performed in
1666	G. Longo: Scientific advise, test and validation of the	1728	part under DOE Contract DE-AC02-76SF00515.
1667	modified METAPHOR pipeline, text of the METAPHOR	1729	
1668	section		
1669	I.A. Almosallam: vetted the early versions of the data set		
1670	and ran many photo-z codes on it, applied GPz to the final		
1671	version and wrote the GPz subsection	1730	
1672	M.L. Graham: Ran the colour-matched nearest-neighbours	1731	
1673	photo-z code on the Buzzard catalog and wrote the relevant	1732	
1674	piece of Section 2; participated in discussions of the analy-	1733	
1675	sis.	1734	
1676	A.J. Connolly: Developed the colour-matched nearest-	1735	
1677	neighbours photo-z code; participated in discussions of the	1736	
1678	analysis.	1737	
1679	E. Nourbakhsh: Ran and optimized TPZ code on the	1738	
1680	Buzzard catalog and wrote a subsection of Section 2 for that	1739	
1681	J. Cohen-Tanugi: contributed to running code, analysis	1740	
1682	discussion, and editing, reviewing the paper	1741	
1683	H. Tranin: contributed to providing SkyNet results and	1742	
1684	writing the relevant section	1743	

APPENDIX A: POINT ESTIMATE PHOTOMETRIC REDSHIFTS

While we do not recommend the use of single point estimates of redshift for most science applications, plots of the point estimates can be a useful qualitative diagnostic of photo-z code performance, i. e. examining point photo-z vs. spec-z plots visually can give a quick impression of some common trends in different codes. Computing point estimate statistics may also be useful for more direct comparisons with previous photo-z evaluations. If a point-estimate is preferred for a specific science case, it is fairly simple to compute the mean, mode, or some other simple estimator from each $p(z)$, so these point estimates can be easily derived from the stored $p(z)$.

There are several common point estimators of photo-z posteriors employed by different codes, e.g. the mode, mean, median of the $p(z)$ distribution. In addition, many of the machine learning based estimators can be set up to return a single redshift solution. For example, SkyNet can be configured to run as a regressor that returns a single float rather than a classifier that returns a 200-bin $p(z)$ estimate. The single value returned by a machine learning based code may not correspond to a particular measure such as the mode or mean, and so to avoid interpretation of results that might be introduced by variations in choice of specific point-estimate implementation per code, we discard the code-specific point estimates. We instead calculate point estimates more uniformly across the codes directly from the $p(z)$ using two measures, z_{PEAK} and z_{WEIGHT} . z_{PEAK} is simply the maximum value attained for each galaxy $p(z)$, the mode of the probability distribution. z_{WEIGHT} is defined similarly to how it is defined in Dahlen et al. (2013), as the weighted mean of the redshift over the main peak of $p(z)$ containing the z_{PEAK} value. The main peak is defined by subtracting $0.05 \times z_{\text{PEAK}}$ from $p(z)$ and identifying the roots to isolate the peak containing z_{PEAK} , z_{WEIGHT} is defined as the weighted mean redshift within this peak. We restrict to a single peak in order to avoid confusion from bimodal and multimodal $p(z)$ such as those shown in bottom panels of Figure 1. For example, for a bimodal probability distribution a weighted mean calculated over both peaks would fall between the peaks, at a redshift where the probability is minimal. Restricting the weighting to a single peak ensures that the point estimate will fall in the region of maximum redshift probability.

- $\text{RMS scatter} < 0.02(1+z)$
- $\text{bias} < 0.003$
- $\text{catastrophic outlier rate} < 10\%$

These definitions are similar, but not exactly the same, as the σ_{IQR} and median bias calculated here, but are similar enough for qualitative comparisons to the LSST goals.

Fig. A1 shows the point estimates for both z_{PEAK} and z_{WEIGHT} . Point density is shown with mixed contours to emphasize that most of the galaxies do fall close to the $z_{\text{phot}} = z_{\text{spec}}$ line, while blue points show differing characteristics of the outlier populations. The red dashed lines indicated the cutoff for catastrophic outliers, defined as: $\max(0.06, 3\sigma_{\text{IQR}})$. As with the full $p(z)$ results, a variety of behaviours are evident in the different codes. Table A1 lists the scatter, bias, and catastrophic outlier fractions for the codes. The performance of the codes for point metrics is highly correlated with performance on $p(z)$ based tests, which is to be expected, given that the point-estimates were derived from the $p(z)$. Some discretization is evident in z_{PEAK} , particularly for SKYNET, due to the finite grid spacing of the reported $p(z)$. These discreteness effects are mitigated by the weighting of z_{WEIGHT} , resulting in a smoother distribution of redshift estimates. Several features perpendicular to the main $z_{\text{phot}} = z_{\text{spec}}$ line are evident. These features are due to the 4000 angstrom break passing through the gaps between adjacent LSST filters. These features are most prominent in template-based codes, but appear to some degree in all codes tested.

In even the best performing codes, there are visible occupied regions away from the $z_{\text{phot}} = z_{\text{spec}}$ line, corresponding to degenerate redshift solutions for certain LSST magnitudes and colors. While use of the full information available via $p(z)$ mitigates their impact, a full understanding of the outlier population is critical for LSST science, particularly in tomographic applications

Finally, we note that all eleven codes are at or near the goals for point-estimates as outlined in the LSST Science Requirements Document¹⁹ and Graham et al. (2018). This is to be expected, given that the requirements were designed such that a point estimate photo-z would meet these requirements for perfect training data to a depth of $i < 25$. But, it is still an encouraging sign, given an updated mock galaxy simulation and the expanded set of photo-z codes tested.

A1 Point Estimate Metrics

We calculate the commonly used point estimate metrics of the overall photo-z scatter (σ_z , the standard deviation of the photo-z residuals), bias, and ‘‘catastrophic outlier rate’’. Specifically, we calculate the metrics as follows: we define e_z as

$$e_z = \frac{z_p - z_s}{1 + z_s} \quad (\text{A1})$$

where z_p is the point estimate and z_s is the true redshift. In practice, because the standard deviation calculation is quite sensitive to the outliers, we define the photo-z scatter, σ in terms of the Interquartile Range (IQR), the difference between the 75th and 25th percentiles of the e_z distribution. In order to match the usual meaning of a 1σ interval, we scale the IQR and define $\sigma_{\text{IQR}} = \text{IQR}/1.349$, as there is a factor of 1.349 difference between the IQR and the standard deviation of a Normal distribution. While many other studies define the bias based on the mean offset between true and estimated redshift, in this study we define the bias as the median value of e_z for the sample. We use median as it is, once again, less sensitive to outliers than the mean. The catastrophic outlier fraction is defined as the fraction of galaxies with e_z greater than the larger of $3\sigma_{\text{IQR}}$ or 0.06, i.e. 3σ outliers with a floor of $\sigma_{\text{IQR}}=0.02$. For reference, the goals stated in Section 3.8 of the LSST Science Book (Abell et al. 2009) for photo-z performance in these metrics, assuming perfect training knowledge (as we are testing in this paper) are:

REFERENCES

- Abbott T., et al., 2005, preprint ([arXiv:astro-ph/0510346](http://arxiv.org/abs/astro-ph/0510346))
 Abell P. A., et al., 2009, preprint ([arXiv:0912.0201](http://arxiv.org/abs/0912.0201)),
 Almosallam I. A., Lindsay S. N., Jarvis M. J., Roberts S. J., 2016a, *MNRAS*, **455**, 2387
 Almosallam I. A., Jarvis M. J., Roberts S. J., 2016b, *MNRAS*, **462**, 726
 Arnouts S., Cristiani S., Moscardini L., Matarrese S., Lucchin F., Fontana A., Giallongo E., 1999, *MNRAS*, **310**, 540
 Behroozi P. S., Wechsler R. H., Wu H.-Y., 2013, *ApJ*, **762**, 109
 Benítez N., 2000, *ApJ*, **536**, 571
 Bernstein G., Huterer D., 2010, *MNRAS*, **401**, 1399
 Biviano A., et al., 2013, *A&A*, **558**, A1
 Blanton M. R., Roweis S., 2007, *AJ*, **133**, 734

¹⁹ available at: <http://ls.st/srd>

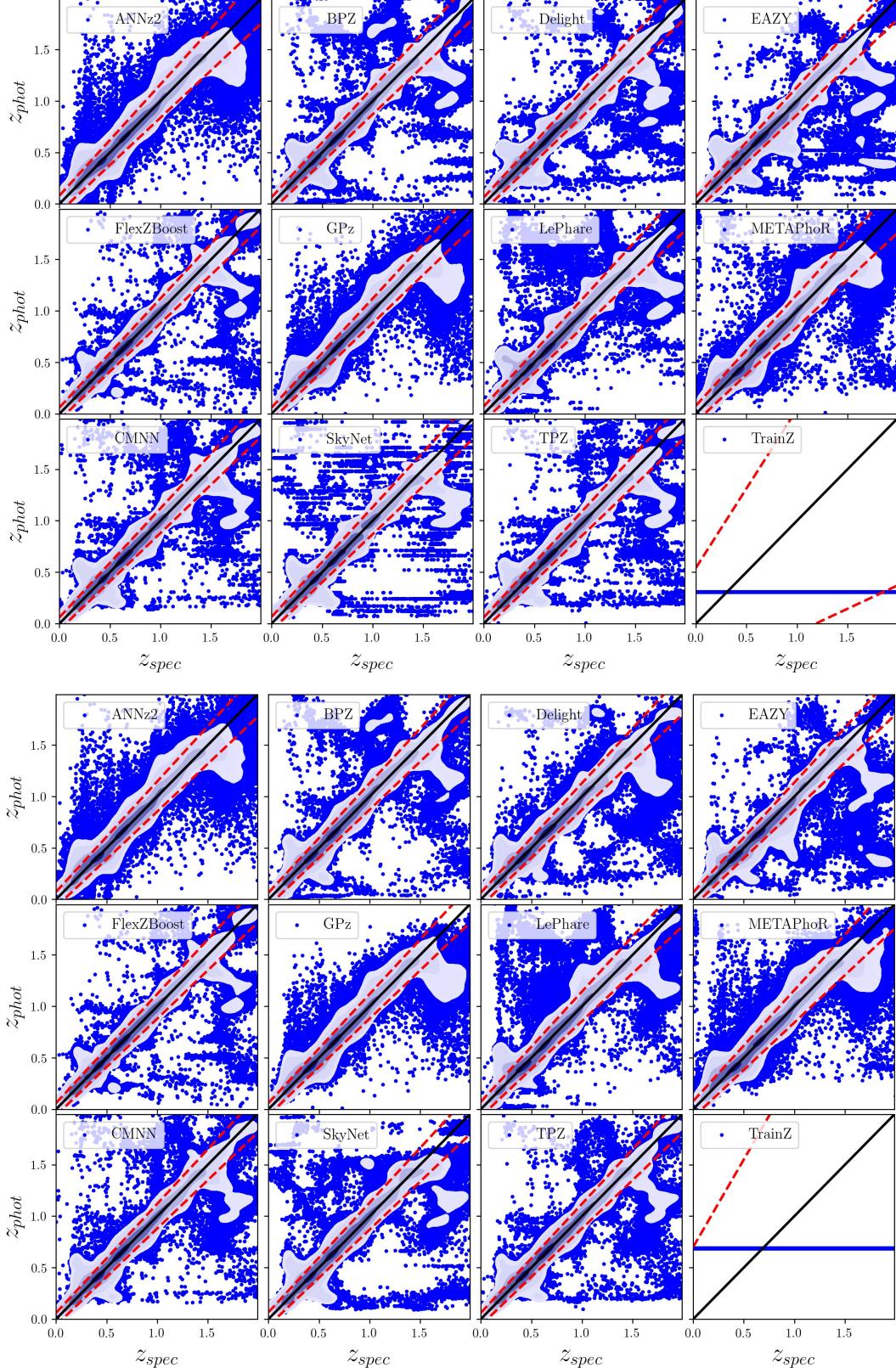


Figure A1. Point estimate photo-z's derived from the posteriors. Top panel shows z_{PEAK} , while bottom panel shows z_{WEIGHT} . Point estimate density is represented with fixed density contours, while outliers at lower density are represented by blue points. While use of point-estimate photo-z's is not recommended, they do make for useful comparative and visual diagnostics. In the lower-right panel of each plot, the TRAINZ estimator results in identical photo-z estimates at the mode and mean of the training set $N'(z)$ distribution for all galaxies.

Table A1. Point estimate statistics

Photo-z Code	ZPEAK			ZW EIGHT		
	$\frac{\sigma_{IQR}}{(1+z)}$	median	outlier fraction	$\frac{\sigma_{IQR}}{(1+z)}$	median	outlier fraction
ANNz2	0.0270	0.00063	0.044	0.0244	0.000307	0.047
BPZ	0.0215	-0.00175	0.035	0.0215	-0.002005	0.032
DELIGHT	0.0212	-0.00185	0.038	0.0216	-0.002158	0.038
EAZY	0.0225	-0.00218	0.034	0.0226	-0.003765	0.029
FLEXZBOOST	0.0154	-0.00027	0.020	0.0148	-0.000211	0.017
GPz	0.0197	-0.00000	0.052	0.0195	0.000113	0.051
LEPHARE	0.0236	-0.00161	0.058	0.0239	-0.002007	0.056
METAPHOR	0.0264	0.00000	0.037	0.0262	0.001333	0.048
CMNN	0.0184	-0.00132	0.035	0.0170	-0.001049	0.034
SKYNET	0.0219	-0.00167	0.036	0.0218	0.000174	0.037
TPZ	0.0161	0.00309	0.033	0.0166	0.003048	0.031
TRAINZ	0.1808	-0.2086	0.000	0.2335	0.022135	0.000

- 1858 Blanton M. R., et al., 2005, *AJ*, **129**, 2562
 1859 Bonnett C., 2015, *MNRAS*, 449, 1043
 1860 Bonnett C., et al., 2016, *Phys. Rev. D*, **94**, 042005
 1861 Brammer G. B., van Dokkum P. G., Coppi P., 2008, *ApJ*, **686**, 1503
 1862 Breiman L., Friedman J. H., Olshen R. A., Stone C. J., 1984, *Classification and Regression Trees, Statistics/Probability Series*. Wadsworth Publishing Company, Belmont, California, U.S.A
 1863 Bresia M., Cavuoti S., D'Abrusco R., Longo G., Mercurio A., 2013, *ApJ*, 772
 1864 Bresia M., Cavuoti S., Longo G., De Stefano V., 2014, *A&A*, 568
 1865 Bresia M., Cavuoti S., Amaro V., Riccio G., Angora G., Velucci C., Longo G., 2018, preprint, ([arXiv:1802.07683](https://arxiv.org/abs/1802.07683))
 1866 Carrasco Kind M., Brunner R. J., 2013, *MNRAS*, 432, 1483
 1867 Carrasco Kind M., Brunner R. J., 2014, *MNRAS*, 442, 3380
 1868 Carruoti S., Bresia M., De Stefano V., Longo G., 2015, *Exp. Astron.*, 39, 45
 1869 Cavuoti S., Amaro V., Bresia M., Vellucci C., Tortora C., Longo G., 2017a, *MNRAS*, 465, 1959
 1870 Cavuoti S., et al., 2017b, *MNRAS*, 466, 2039
 1871 Chen T., Guestrin C., 2016, in *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '16. ACM, New York, NY, USA, pp 785–794, doi:10.1145/2939672.2939785, http://doi.acm.org/10.1145/2939672.2939785*
 1872 Dahlen T., et al., 2013, *ApJ*, 775, 93
 1873 Dawson W. A., Schneider M. D., Tyson J. A., Jee M. J., 2016, *ApJ*, 816, 11
 1874 Duncan K. J., Jarvis M. J., Brown M. J. I., Röttgering H. J. A., 2018, *Monthly Notices of the Royal Astronomical Society*, p. 940
 1875 Fernández-Soto A., Lanzetta K. M., Yahil A., 1999, *ApJ*, 513, 34
 1876 Firth A. E., Lahav O., Somerville R. S., 2003, *MNRAS*, 339, 1195
 1877 Freeman P. E., Izbicki R., Lee A. B., 2017, *MNRAS*, 468, 4556
 1878 Graff P., Feroz F., Hobson M. P., Lasenby A., 2014, *MNRAS*, 441, 1741
 1879 Graham M. L., Connolly A. J., Ivezić Ž., Schmidt S. J., Jones R. L., Jurić M., Daniel S. F., Yoachim P., 2018, *AJ*, 155, 1
 1880 Green J., et al., 2012, preprint ([arXiv:1208.4012](https://arxiv.org/abs/1208.4012)),
 1881 Hildebrandt H., et al., 2010, *A&A*, 523, A31
 1882 Hofmann B., Mathé P., 2018, *Inverse Problems*, 34, 015007
 1883 Ibert O., et al., 2006, *A&A*, 457, 841
 1884 Ivezic Ž., et al., 2008, preprint ([arXiv:0805.2366](https://arxiv.org/abs/0805.2366)),
 1885 Izbicki R., Lee A. B., 2017, *Electron. J. Statist.*, 11, 2800
 1886 Izbicki R., Lee A. B., Freeman P. E., 2017, *Ann. Appl. Stat.*, 11, 698
 1887 Laigle C., et al., 2016, *ApJS*, 224, 24
 1888 Laureijs R., et al., 2011, preprint (1110.3193),
 1889 Leistedt B., Hogg D. W., 2017, *ApJ*, 838, 5
 1890 Maddox N., Hewett P. C., Warren S. J., Croom S. M., 2008, *MNRAS*, 386, 1605
 1891 Malz A., Hogg D., in prep., CHIPPR, chippr
 1892 Malz A., Marshall P., DeRose J., Graham M., Schmidt S., Wechsler R., 2018, *AJ*, Accepted,
 1893 Mandelbaum R., et al., 2008, *MNRAS*, 386, 781
 1894 Massarotti M., Iovino A., Buzzoni A., 2001, *A&A*, 368, 74
 1895 Masters D. C., Stern D. K., Cohen J. G., Capak P. L., Rhodes J. D., Castander F. J., Paltani S., 2017, *ApJ*, 841, 111
 1896 Newman J. A., 2008, *ApJ*, 684, 88
 1897 Newman J. A., et al., 2015, *Astroparticle Physics*, 63, 81
 1898 Polsterer K. L., D'Isanto A., Giesecke F., 2016, preprint ([arXiv:1608.08016](https://arxiv.org/abs/1608.08016)),
 1899 Rasmussen C., Williams C., 2006, *Gaussian Processes for Machine Learning. Adaptative computation and machine learning series*, MIT Press, Cambridge, MA
 1900 Reddick R. M., Wechsler R. H., Tinker J. L., Behroozi P. S., 2013, *ApJ*, 771, 30
 1901 Riccio G., Bresia M., Cavuoti S., Mercurio A., di Giorgio A., Molinari S., 2017, *PASP*, 129
 1902 Richards G. T., et al., 2006, *ApJS*, 166, 470
 1903 Richards G. T., et al., 2015, *ApJS*, 219, 39
 1904 Sadeh I., Abdalla F. B., Lahav O., 2016, *PASP*, 128, 104502
 1905 Salvato M., et al., 2011, *ApJ*, 742, 61
 1906 Sánchez C., et al., 2014, *MNRAS*, 445, 1482
 1907 Scott D. W., 1992, *Multivariate Density Estimation. Theory, Practice, and Visualization*. Wiley
 1908 Skrutskie M. F., et al., 2006, *AJ*, 131, 1163
 1909 Tanaka M., et al., 2017, preprint ([arXiv:1704.05988](https://arxiv.org/abs/1704.05988)),
 1910 York D. G., et al., 2000, *AJ*, 120, 1579
 1911 de Jong J. T. A., Verdoes Kleijn G. A., Kuijken K. H., Valentijn E. A., 2013, *Exp. Astron.*, 35, 25
 1912 de Jong J. T. A., et al., 2017, *A&A*, 604, A134