

Implicit assumptions and their impact on photometric redshift PDF performance in the context of LSST

S.J. Schmidt¹, A.I. Malz^{2,3}, J.Y.H. Soo⁴, M. Brescia⁵, S. Cavaudi^{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

1 December 2018

ABSTRACT

In order to maximize scientific returns of current and upcoming galaxy surveys, the photometric redshift ($\text{photo-}z$) posterior distributions produced by redshift estimation codes must be accurate probability distribution functions (PDFs). However, the posteriors resulting from a number of current techniques are not, in general, consistent with each other, affected by implicit assumptions made by each code, and an optimal method for obtaining an accurate PDF estimate remains unclear. We present the results of an initial study of the Large Synoptic Survey Telescope Dark Energy Science Collaboration (LSST-DESC) evaluating twelve $\text{photo-}z$ algorithms using complete and representative training data and evaluate multiple metrics to test how accurately the posteriors represent probability distributions. We observe several trends, including systematic biases and an overall over/under-prediction in the broadness of the PDFs in many of the codes which may be symptomatic of implementation problems or problems in underlying algorithm design. 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), Hyper Suprime-Cam (HSC) Survey (Aihara et al. 2018a,b), 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 individual galaxy 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 LSST Dark Energy Science Collaboration (LSST-DESC) developed a separate Science Requirements Document (The LSST Dark Energy Science Collaboration et al. 2018), which forecasts the constraining power of five cosmological probes using somewhat conservative assumptions to define requirements on systematic errors for several measurements. These include even more stringent requirements on photometric redshift performance than those included in the LSST SRD, though most of the initial LSST-DESC requirements are defined in terms of tomographic bin populations rather than on individual object redshifts. The tremendous size of LSST’s galaxy catalogue 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. The inclusion of these populations introduce major physical degeneracies, for example the Lyman break/Balmer break degeneracy, that were not present in the populations covered in shallower 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). These issues are not unique to LSST, and are present in Stage III Dark Energy surveys; however, in order to meet the demanding error budgets of Stage IV projects such as LSST and LSST-DESC it will be necessary to fully characterize those degeneracies wherein multiple redshift solutions have comparable likelihood to per cent level accuracy.

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, 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 relative 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, there was a 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).

For cosmological measurements, certain science cases require redshift information on individual objects, e. g. identification of host galaxy redshift for supernova classification, or identifying potential cluster membership. Other science cases seem to need only ensemble redshift information; for instance many current cosmic shear techniques require only the overall redshift distribution $N(z)$ for tomographic redshift samples. However, even such cases require individual object redshift estimates for portions of the analysis, for example in determining galaxy intrinsic alignments in weak lensing samples. In addition, recent data-driven techniques employing hierarchical Bayesian or Gaussian Process methods have emerged that calibrate redshift distributions using individual $p(z)$ estimates (e. g. Sánchez & Bernstein 2018). These methods assume that the $p(z)$ for each galaxy is an accurate PDF, and such methods break down if this assumption is invalid. Thus, even methods that seem to need only ensemble $N(z)$ may actually require accurate $p(z)$ in order to meet stringent survey requirements. Large photometric surveys such as LSST must develop algorithms that simultaneously meet the needs of all science cases. In order to meet these ambitious goals for photo- z accuracy, every aspect of photo- z estimation will have to be optimized: the algorithms employed, both template and machine-learning based (both 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 catalogue compression schemes that balance information retention against limited storage resources.

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.

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 photomet-

ric bands. Rau et al. (2015) introduced a new method for improving redshift PDFs using an ordinal classification algorithm. DES compared several codes for point estimates and a subset with $p(z)$ information (Sánchez et al. 2014). A follow up paper examined summary statistics of photo- z interim posteriors for tomographically binned galaxy subsamples (Bonnett et al. 2016).

This paper is distinguished by its focus on 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 here are a major focus of the Photometric Redshift working group of the 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. In this initial paper we focus on evaluating the performance of photometric redshift codes and PDF-based performance metrics in the presence of complete and representative training sets. Specific implementation choices in each code will influence the resultant posterior distributions, for example choice of prior parameterization in template-based codes, the bandwidth size chosen for machine learning based codes, or even the output format chosen for storing the PDF. We have attempted to minimize the impact of many of these factors when comparing codes, for example by using the same template set for all template-based codes, and using a training set that is drawn from the same underlying population as the test sample, to create a controlled environment in which to compare the photo- z PDFs derived from each method. We explore a number of performance metrics in this paper that test whether the posterior estimates are actual PDFs. 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. Identifying and fixing failure modes within codes may aid us in reaching the stringent photo- z performance goals set out for LSST. We note that these initial tests are a necessary requirement for photo- z codes that will be used in cosmological analyses; however, meeting these requirements is only the first stage in the process, and can be thought of as an initial test under near perfect conditions to test for problems before further complexities are added in future analyses.

The outline of the paper is as follows: in § 2 we present the simulated data set; in § 3 we describe the current generation codes employed in the paper; in § 4 we discuss the interpretation of photo- z PDFs in terms of metrics of accuracy; in § 5 we show our results and compare the performance of the codes; in § 6 we offer our conclusions and discuss future extensions of this work.

2 THE SIMULATION AND MOCK GALAXY CATALOG

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

2.1 Buzzard-v1.0 simulation

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

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

³ Available at: http://lsst-desc.org/sites/default/files/DESC_SRM_V1_1.pdf

⁴ <https://github.com/vipasu/addseds>

⁵ <http://kcorrect.org>

4 LSST Dark Energy Science Collaboration

227 forced photometry, and the photo-z codes will treat them
228 as such. Only 2.0 per cent of our sample galaxies contain a
229 photometric band with a non-detection, the vast majority
230 of which are in the u -band.

231 2.1.1 Selection of training and test sets

232 The total catalogue covered 400 square degrees and con-
233 tained 238 million galaxies to an apparent magnitude limit
234 of $r = 29$ and spanning the redshift range $0 < z \leq 8.7$. In order
235 for statistical errors not to dominate, we need less than one
236 million galaxies in our sample. Several studies claim that
237 only a few tens of thousands of spectra are necessary to cal-
238ibrate photo-z surveys to Stage IV requirements (e. g. [Bern-](#)
[stein & Huterer \(2010\)](#), [Masters et al. \(2017\)](#)). Therefore, we
240 aim for a final number of training galaxies between 3×10^4
241 and 5×10^4 in our sample. In order to reduce our sample to
242 a reasonable size, we limit our dataset to a subset of ~ 16.8
243 square degrees selected from five separate spatial regions
244 of the simulation. Systematic problems with galaxy colors
245 above $z > 2$ were observed, so the catalogue was limited to
246 include only galaxies in the redshift range $0 < z \leq 2.0$. A
247 random subset of the the remaining galaxies was chosen,
248 and placed at random into either a “training” set (10 per
249 cent of the sample), for which the galaxies true redshifts
250 will be supplied, or a “test” set (the remaining 90 per cent
251 of the sample), for which each code will need to predict a
252 redshift PDF for each galaxy. Finally, we restrict our analy-
253 sis to a sample with an apparent magnitude limite $i < 25.3$,
254 which give a signal-to-noise ~ 30 for most galaxies, a cut
255 often referred to as the expected “LSST Gold Sample”. This
256 magnitude cut results in a training set with 44 404 galaxies
257 and a test set containing 399 356 galaxies. All subsequent
258 results will evaluate this “gold sample” test set. In order
259 to blind results, initially redshifts were not revealed for the
260 “test” set, and were only supplied for the training sample
261 galaxies. This prevented code runners from tweaking results
262 and over-fitting to the specific test set.

263 2.1.2 Templates

264 As mentioned in Section 2.1, the SEDs in the Buzzard sim-
265 ulation are drawn from an empirical set of SEDs taken from
266 the SDSS DR6 NYU-VAGC, a sample of roughly $\sim 5 \times 10^5$
267 galaxies with spectra in SDSS. To determine a finite set of
268 templates to use with template fitting codes we take the five
269 SED weight coefficients for each of the galaxies in the SDSS
270 sample and run a simple K-means clustering algorithm on
271 this five dimensional space. Each dimension was normalized
272 such that it spanned an interval $[0, 1]$. The K-means clus-
273 ters partition the five-dimensional space of coefficients into
274 Voronoi cells, spanning the space of coefficients in a way
275 that properly reflects the underlying density in the coeffi-
276 cients. Thus, the resultant SEDs constructed using the cell
277 centers as weight coefficients will provide a reasonable span-
278 ning SED set. An ad-hoc number of $K = 100$ was chosen and
279 the 100 K-means centre positions are taken as the weights
280 for the k-CORRECT SED components to construct one hun-
281 dred template SEDs. These 100 templates were provided,
282 and the templates were used by both BPZ and LEPHARE;
283 however, because EAZY was designed and written to use

284 the same five basis templates employed by k-CORRECT when
285 constructing our mock galaxies, EAZY was run using linear
286 combinations of these five templates rather than using the
287 100 discrete templates. The ability to fit for linear combi-
288 nations of templates highlights an important implementation
289 difference between similar photo-z codes.

290 2.1.3 Limitations

291 For our initial investigation of photometric redshift codes,
292 we begin with a data set that is somewhat idealized, and
293 does not contain all of the complicating factors present in
294 real data. In several cases, the simplification is done with
295 a purpose, with potentially confounding effects excluded
296 in order to better isolate the differences between current-
297 generation photo-z codes, and their causes. We list several
298 of the simulations limitations in this section. As the sim-
299 ulation is catalogue-based, no image level effects, such as
300 photometric measurement effects, object blending, contami-
301 nation from sky background (Zodiacal light, scattered light,
302 etc...), lensing magnification, or Galactic reddening are in-
303 cluded. No stars are included in the catalogue, nor are the
304 effects of AGN. As all SEDs are constructed from only five
305 basis templates, properties of the galaxy population will be
306 restricted to follow linear combinations of the characteristics
307 of the five basis templates, so certain non-linear features, for
308 example the full range of emission line fluxes relative to the
309 continuum, will not be included in the model galaxy popu-
310 lation. Moreover, the linear combinations of templates are
311 modeled on the $\sim 5 \times 10^5$ SDSS galaxies discussed in Sec-
312 tion 2.1, and thus only galaxies that resemble those spec-
313 troscopically observed by the SDSS will be included in the
314 sample. No additional dust reddening intrinsic to the host
315 galaxy is included, the only approximation of dust extinc-
316 tion comes in the form of dust encoded in the five basis SEDs
317 via the training set used to create the basis templates. Sim-
318 ple linear combinations of these basis templates will, once
319 again, not explore the full range of realistic dust extinction
320 observed in galaxy populations. While these idealized con-
321 ditions limit the realism of our galaxy population, some are
322 also by design. We aim to test the photo-z codes at a very
323 basic level, and a simplified model assures that differences in
324 results seen between the codes are due to fundamental differ-
325 ences in their underlying assumptions and implementa-
326 tion details, rather than more nuanced implementation details.

3 METHODS

327 Here we outline the photo-z PDF codes tested in this study.
328 In total, eleven distinct codes are tested. This sample is not
329 comprehensive, codes were chosen based on the expertise
330 available within the group; however, those chosen do cover
331 a broad range of the current-generation methods used in the
332 field. Both template-based and machine learning approaches
333 are included and each are described separately in Secs. 3.1
334 and 3.2 respectively. The list of codes are summarized in Ta-
335 ble. 1. All code runners were asked to output redshift pos-
336 terior estimates on 200 linear-spaced bins between redshifts
337 0 and 2.

338 The questions that must be answered for each code are:

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/ixxael/Delight
FLEXZBOOST	ML	Izbicki & Lee (2017)	https://github.com/tospisici/flexcode; https://github.com/rizbicki/FlexCoDE
GPz	ML	Almosallam et al. (2016b)	https://github.com/OxfordML/GPz
METAPHOR	ML	Cavuoti et al. (2017)	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	

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

340 marginalizing over all SED-types with a simple sum (Eq. 3
341 from Benítez 2000):

$$342 p(z|C, m_0) \propto \sum_T p(z, T|m_0) p(C|z, T) \quad (2)$$

343 where the first term on the right-hand side is the Bayesian
344 prior and the second term is the traditional likelihood.
345 The prior is assumed to have the form: $p(z, T|m_0) =$
346 $p(T|m_0) p(z|T, m_0)$, i.e. it parameterizes the prior as an
347 evolving type fraction with apparent magnitude, combined
348 with a prior on the expected redshift probability distribution
349 as a function of both apparent magnitude and SED-type.

350 In this paper we use BPZ v 1.99.3. The template set
351 employed here is the set of 100 discrete SEDs described in
352 Section 2.1.2 To keep the number of free parameters to a
353 manageable level the SEDs in the training set are sorted
354 by the rest-frame $u-g$ colour and split into three “broad”
355 SED classes, equivalent to the E, Sp and Im/SB types in
356 Benítez (2000). We assume the same functional form for
357 the Bayesian priors as used by Benítez (2000), and utilize
358 the training-set galaxies with known SED-type, redshift, and
359 apparent magnitude to determine the type fractions and the
360 best fit for the eleven free parameters of the prior. For galaxies
361 that are not detected in a measured band, the placeholder
362 value is replaced with an estimate of the one σ detection
363 limit in that particular band, i.e. a value close to the es-
364 timated sky noise threshold. The type-marginalized $p(z)$ is
365 generated by setting the parameter PROBS_LITE=TRUE in the
366 BPZ parameter file.

3.1 Template-based Approaches

367 We test three publicly available and commonly used
368 template-based codes: BPZ, EAZY, and LEPHARE. All
369 three codes follow the standard procedure for template-
370 based redshift estimation: calculate model fluxes for a set
371 of template spectral energy distributions (SEDs) on a grid
372 of redshift values and calculating a χ^2 merit function using
373 the observed and model fluxes:

$$374 \chi^2(z, T, A) = \sum_i^{N_{\text{filt}}} \left(\frac{F_{\text{obs}}^i A \times F_{\text{pred}}^i(T, z)}{\sigma_{\text{obs}}^i} \right)^2 \quad (1)$$

375 where A is a normalization factor, $F_{\text{pred}}^i(T, z)$ is the flux
376 predicted for a template T at redshift z . F_{obs}^i is the observed
377 flux in a given band i and σ_{obs}^i is the observed flux error.
378 N_{tot} is the total number of filters, in our case the six *ugrizy*
379 LSST filters. Specific implementation details of each code,
380 e.g. prior form and implementation, are described below.

381 In this paper we use BPZ v 1.99.3. The template set
382 employed here is the set of 100 discrete SEDs described in
383 Section 2.1.2 To keep the number of free parameters to a
384 manageable level the SEDs in the training set are sorted
385 by the rest-frame $u-g$ colour and split into three “broad”
386 SED classes, equivalent to the E, Sp and Im/SB types in
387 Benítez (2000). We assume the same functional form for
388 the Bayesian priors as used by Benítez (2000), and utilize
389 the training-set galaxies with known SED-type, redshift, and
390 apparent magnitude to determine the type fractions and the
391 best fit for the eleven free parameters of the prior. For galaxies
392 that are not detected in a measured band, the placeholder
393 value is replaced with an estimate of the one σ detection
394 limit in that particular band, i.e. a value close to the es-
395 timated sky noise threshold. The type-marginalized $p(z)$ is
396 generated by setting the parameter PROBS_LITE=TRUE in the
397 BPZ parameter file.

3.1.1 BPZ

398 BPZ⁶ (Bayesian Photometric Redshift, Benítez 2000) is a
399 template-based photo-z code that compares the expected
400 colors (C) calculated for a set of spectral energy distribu-
401 tion (SED) types/templates (T) to the observed colors to
402 calculate the likelihood of observing colors at each redshift
403 for each type, $p(C|z, T)$. The likelihoods at each redshift are
404 related to the χ^2 in Equation 1 by the simple form: like-
405 lihood $\propto e^{-\chi^2/2}$. The code employs an empirically deter-
406 mined Bayesian prior in apparent magnitude (m_0) and SED-
407 type. Assuming that the SED-types are spanning and exclu-
408 sive, we can determine the redshift posterior $p(z|C, m_0)$ by

409 3.1.2 EAZY

410 EAZY⁷ (Easy and Accurate Photometric Redshifts from
411 Yale, Brammer et al. 2008) is a template-based photo-z
412 code that includes several features that extend the basic
413 χ^2 fit used in many template codes. The code can fit the
414 observed photometry with SEDs created from a linear com-
415 bination of a set of templates at each redshift, and the best-
416 fit SED is found by simultaneously fitting one, two or all
417 of the templates by minimizing χ^2 . The minimized $\chi^2(z)$ is

⁶ <http://www.stsci.edu/~dcoe/BPZ/>

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

6 LSST Dark Energy Science Collaboration

then combined with an apparent magnitude prior to obtain the posterior redshift probability distribution. On examination of the source code for EAZY, it appears that rather than marginalizing across all templates in the χ^2 calculation, EAZY takes only the minimum value of χ^2 at each redshift. This improper marginalization does not lead to the correct posterior distribution, an implementation issue that will need to be addressed in the future. 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 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 a set of spectral energy distribution (SED) which can be either synthetic or based on a semi-empirical approach.

Each SED is convolved with the simulated LSST filter transmission curves (accounting for instrument efficiency). The computed photo- z is then the value that minimizes the merit function $\chi^2(z, T, A)$ from Arnouts et al. (1999), and given in Equation 1.

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, and the full $p(z)$ corresponds to the likelihoods calculated at each point on our z -grid.

3.2 Training-based Codes

The training-based codes use a variety of algorithms in order to estimate $p(z)$, specifics of each implementation are described in the subsections. Some aspects of data treatment were left to the individual code runners, for example, whether/how to split the available data with known redshifts into separate training and validation sets. Another key difference is the treatment of non-detections in one or more bands. Some codes choose to ignore a band, others replace the value with either an estimate for the detection limit, the mean of other values in the training set, or another default value. There are varying conventions among training-based codes for treatment of non-detections, and no one prescription dominates in the photo- z literature. The specific choices for each code affect the results, and contribute to the implicit prior influencing their output. However, we remind the reader that only 2.0 per cent of our sample has non-detections, almost exclusively in the u-band, and thus should not dominate the code performance differences.

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

3.2.1 ANNz2

ANNz2⁹ (Sadeh et al. 2016) is a software 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 either run a single machine learning algorithm and return results, or it can run multiple algorithms simultaneously and output photo- z 's as a weighted combination of the different algorithms. In this study, only ANNs were employed. The redshift PDFs are produced by running an ensemble set of ANNs, each with different random seeds used in initialization of input parameters for training. Uncertainties for each method are estimated from a KNN-uncertainty estimator (Oyaizu et al. 2008). The final PDF can either be the “best” of the candidate PDFs, or a weighted average of the PDFs based on their error estimates for each of the ensemble members.

In this study, ANNz2 v. 2.0.4 was used. 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 Colour-Matched Nearest-Neighbours

The nearest-neighbours colour-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 colour-matched subset of training galaxies. This subset is identified by first calculating the Mahalanobis distance D_M in colour-space between the test galaxy and all training-set galaxies:

$$D_M = \sum_{\text{colours}}^{N_{\text{colours}}} \frac{(c_{\text{train}} - c_{\text{test}})^2}{(\delta c_{\text{test}})^2} \quad (3)$$

where c colour, δc_{test} is the measurement error on the colour, and N_{colours} is the total number of colors (i.e., in our case $u-g$, $g-r$, $r-i$, $i-z$, and $z-y$). Then, we choose a threshold value for D_M that defines the colour match set based on a set value of the percent point function (PPF): for example, for $N_{\text{dof}} = 5$, choosing a vPPF = 95 per cent of all training galaxies consistent with the test galaxy will have $D_M < 11.07$ (where N_{dof} is the number of degrees of freedom, in this case the number of colours). If a galaxy had a non-detection in a band, that band was dropped and N_{dof} was reduced by one in the colour-matching space. For a given test galaxy, the $p(z)$ is the normalized distribution of the true catalogue redshifts of this colour-matched subset of training galaxies.

We have applied the nearest-neighbours colour-matching photometric redshift estimator described in G18

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

¹⁰ <http://tmva.sourceforge.net/>

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 colour 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’ colours and magnitudes. Third, we set different parameters for the identification of the colour-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 colour-matched subset of training galaxies and used the redshift of nearest neighbour in colour-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 an accurate $p(z)$ in all cases – especially for galaxies that are bright and/or have few matches in colour-space. Since an accurate 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 colour-matched training galaxies, but weight them by the inverse of their Mahalanobis distance. This weighting maintains some of the accuracy that was previously achieved by simply using the nearest neighbour in colour-space. We then use the weights to create the $p(z)$ instead of having the redshift of each colour-matched training-set galaxy count equally. To obtain a robust estimate of the $p(z)$ for galaxies with a small number of colour-matched training set galaxies, when this number is less than 20 the nearest 20 neighbours in colour-space are used instead, and we convolve the $p(z)$ with a Gaussian with a standard deviation of:

$$\sigma = \sigma_{\text{train}} \sqrt{(\text{PPF}_{20}/0.95)^2 - 1} \quad (4)$$

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 less precise photo- z estimates compared to those presented in G18, but they will all have significantly more accurate 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 accurate $p(z)$ 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. Delight models the underlying latent SEDs as a linear combination of a set of pre-defined template SEDs, plus zero mean Gaussian processes with factorized kernels. 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,

$$p(z|\hat{\mathbf{F}}) \approx \sum_i p(\hat{\mathbf{F}}|z, t_i) p(z|t_i) p(t_i), \quad (5)$$

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 posterior 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:

$$p(\hat{\mathbf{F}}|z, t_i) = \int p(\hat{\mathbf{F}}|\mathbf{F}) p(\mathbf{F}|z, z_i, \hat{\mathbf{F}}_i) d\mathbf{F}, \quad (6)$$

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. 6 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,

$$F_b \sim \mathcal{GP} \left(\mu^F, k^F \right), \quad (7)$$

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 using the provided training sample. 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. The default settings of DELIGHT were used, with the exception that the PDF bins were set to be linearly-spaced rather than logarithmic. In this study a flat prior in magnitude/type is assumed.

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

8 LSST Dark Energy Science Collaboration

628 3.2.4 FlexZBoost

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

$$636 f(z|\mathbf{x}) = \sum_i \beta_i(\mathbf{x})\phi_i(z). \quad (8)$$

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

$$639 \beta_i(\mathbf{x}) = \mathbb{E}[\phi_i(z)|\mathbf{x}] \equiv \int f(z|\mathbf{x})\phi_i(z)dz. \quad (9)$$

640 These coefficients can easily be estimated from data by regression since $\mathbb{E}[\phi_i(z)|\mathbf{x}]$ is the regression of $\phi_i(z)$ (a transformation of Z) on X .

641 In this paper, we use XGBOOST (Chen & Guestrin 2016)
 642 for the regression part; it should however be noted that
 643 FLEXCODE-RF (also on GitHub), based on Random Forests,
 644 generally performs better for smaller data sets. As our basis,
 645 we choose a standard Fourier basis. There are two tuning
 646 parameters in our $p(z)$ estimate: (i) the number of terms, I ,
 647 in the series expansion in Eq. 8, and (ii) an exponent α that
 648 we use to sharpen the computed density estimates $\hat{f}(z|\mathbf{x})$,
 649 according to $\hat{f}(z|\mathbf{x}) \propto \hat{f}(z|\mathbf{x})^\alpha$. We reserve 15% of the training
 650 set data as a validation set, and choose both I and α
 651 in an automated way by minimizing the weighted L_2 -loss
 652 function (Eq. 5 in Izbicki & Lee 2017) on the validation set.
 653 While the native storage format for FLEXCODE encodes the
 654 PDF using the coefficients shown in Equation 9, to match
 655 the output format requested of other codes we discretize our
 656 final estimates into 200 bins linearly spaced in $0 < z < 2$.

659 3.2.5 GPz

660 GPz¹⁴ (Almosallam et al. 2016a,b) is a sparse Gaussian process based code, 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:

$$670 f(x) = \sum_{i=1}^m \phi_i(x)w_i, \quad (10)$$

671 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

¹² <https://github.com/tospispis/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} .

¹⁴ <https://github.com/OxfordML/GPz>

674 such as the sigmoid or the squared exponential, have the
 675 advantage of being universal approximators, i.e. there exist
 676 a function of that form that can approximate any function,
 677 with mild assumptions, to any desired degree of accuracy.
 678 The details on how to learn the parameters of the model and
 679 the hyper-parameters of the basis functions are described in
 680 Almosallam et al. (2016b).

681 A unique feature in GPz, is that the variance estimate is
 682 composed of two terms each quantifying a different source of
 683 uncertainty. One term (the model uncertainty) reflects how
 684 much of the uncertainty is due to lack of training samples at
 685 the location of interest, whereas the second term (the noise
 686 uncertainty) reflects how much of the uncertainty is caused
 687 from observing many noisy samples at that location. Thus,
 688 the predictive variance can determine whether we need more
 689 representative samples or more precise samples for any par-
 690 ticular location in the input space. GPz can also emphasize
 691 the importance of some samples as weights. This weight can
 692 be for example $|z_{\text{spec}} - z_{\text{phot}}|/(1 + z_{\text{spec}})$ to target the de-
 693 sired objective of minimizing the normalized redshift error
 694 or as a function of their probability in the test set relative
 695 to the training set in order to pressure the model to better
 696 fit samples that are rare in the training set but are expected
 697 to be abundant during testing.

698 The data is prepared for GPz by taking the log of the
 699 magnitude errors, decorrelating the data set using PCA and
 700 imputing any missing magnitude values using a simple lin-
 701 ear model that estimates the missing magnitudes given the
 702 observed ones. The log transformation helps to smooth the
 703 long tail distribution of the magnitude errors, which is more
 704 stable numerically and makes the optimization process un-
 705 constrained. The missing values are imputed by computing
 706 the mean of the training set μ and its covariance Σ , then
 707 we use the following equation to estimate the missing values
 708 from the observed ones

$$709 x_u = \mu_u + \Sigma_{uo}\Sigma_{oo}^{-1}(x_o - \mu_o), \quad (11)$$

710 where the subscript o in x_o indexes the *observed* part of
 711 the input x , whereas the subscript u indexes the *unobserved*
 712 set (similarly for μ and Σ). This is the optimal expected
 713 value of the unobserved variables given the observed ones
 714 if the distribution is jointly Gaussian, note that if the vari-
 715 ables are independent, i.e. $\Sigma_{uo} = 0$, this will reduce to a
 716 simple average predictor. We use the Variable Covariance
 717 (VC) option in GPz with 200 basis functions after we note
 718 that there is no significant increase in the performance on
 719 the validation set (using 80%-20% training-validation split)
 720 and with no cost-sensitive learning applied.

721 3.2.6 METAPhOR

722 METAPhOR (Machine-learning Estimation Tool for Accu-
 723 rate Photometric Redshifts, Cavuoti et al. 2017) is a pipeline
 724 designed to provide photo-z point estimates and a reliable
 725 PDF for machine learning (ML) based techniques. It in-
 726 cludes pre- and post-processing phases, hosting a photo-z
 727 prediction engine based on the Multi Layer Perceptron with
 728 Quasi Newton Algorithm (MLPQNA).

729 METAPhOR includes data modules for pre-processing,
 730 photo-z estimation, and PDF estimation, and post-
 731 processing. The pre-processing includes a model for pertur-

bation of the photometry that is employed in calculating the PDF of the photo- z estimation errors. The photometric perturbation is defined as: $m_{ij} = m_{ij} + \alpha_i F_{ij} * u_{\mu=0,\sigma=1}$, 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 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 replace the missing detections with values based on their neighbors. 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).

The neural network is configured as a standard multilayer 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. 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.

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 $dm \sim 0.7526 \text{ mag}$ 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 colours $u-g, g-r, r-i, i-z, z-y$ 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 200 bins 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 as a pathological example. For TRAINZ, as we will we call this simple estimator, we well define $p(z)$ as simply:

$$p(z) = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} z_{\text{train}} \quad (12)$$

That is, we simply set the redshift PDF of every galaxy equal to the normalized $N(z)$ of the training sample. This estimator is essentially a k nearest-neighbour estimator with k equal to the number of galaxies in 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 should be identical, modulo fluctuations due to

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

¹⁶ <https://github.com/mgckind/MLZ>

finite sample size. This is a wildly unrealistic estimator, as it assigns all galaxies, no matter their apparent magnitude, colour, or true redshift, the same redshift PDF, and is thus uninformative at the level of individual object redshifts, but is designed to perform very well for the ensemble of all objects. If the training set was not representative, this estimator would produce biased results, and any attempts to break up the sample into tomographic bins will fail, as every galaxy has an identical $p(z)$. We will discuss this method and cautions relative to metrics in Section 5.3.

4 METRICS FOR QUANTIFYING PDF COMPARISONS

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

$$\hat{p}_j(z_i) = p(z|d_i, I_D, I_H), \quad (13)$$

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

In this study, we isolate the differences in prior information specific to each method by using a single training set I_D^{ML} for all machine learning-based codes and a single template library I_D^T for all template-based codes, and these sets of prior information are carefully constructed to be representative and complete, we have $I_D^{ML} \equiv I_D^T$ for every method H . Thus, we are saying

$$\frac{\hat{p}_{i,H}(z)}{\hat{p}_{i,H'}(z)} \approx \frac{p(z|d_i, I_H)}{p(z|d_i, I_{H'})}, \quad (14)$$

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

As mentioned previously, there are cosmology probes that require knowledge of individual galaxy’s photo-z PDF, for example galaxy intrinsic alignment studies, strong lensing foreground shear prediction, and photometric supernova classification, and others that require only knowledge of the ensemble redshift distribution, $N(z)$. Due to the paucity of principled techniques for using and validating individual galaxy photo-z PDFs, there have been few alternatives to the common practice of reducing photo-z PDFs to point estimates when evaluating and comparing photo-z code performance in the literature. Though this practice should not be encouraged, we also calculate traditional metrics based on the most common point estimators derived from photo-z

PDFs. Those seeking to establish a connection to traditional ways of thinking about redshift estimation may consult the Appendix for these results.

There are a number of metrics that can be used to test the accuracy of a photo-z posterior as an estimator of a true photo-z posterior if it is known. Even for simulated data, the true photo-z PDF is in general not accessible unless the photometry is in fact drawn from the true photo-z likelihoods, a mock catalogue generation procedure that has not yet appeared in the literature. Furthermore, only limited applications of photo-z PDFs that could be used as the basis for a metric have been presented in the literature.

The most popular application of photo-z PDFs by far is the estimation of the overall redshift distribution $N(z)$, the true value of which is known for the BUZZARD simulation and will be denoted as $N'(z)$. Though alternatives exist (Malz & Hogg prep), “stacking” according to

$$\hat{N}^H(z) \equiv \frac{1}{N_{tot}} \sum_i^{N_{tot}} \hat{p}_i^H(z) \quad (15)$$

is the most widely accepted method for obtaining the stacked estimator $\hat{N}^H(z)$ of the redshift distribution from photo-z PDFs derived by a method H . Though we do not endorse the use of the stacked estimator of the redshift distribution, we use it under the assumption that the response of our metrics of $\hat{N}^H(z)$ will be analogous to the same metrics applied to a principled estimator of the redshift distribution. We must note, however, that this is a poor assumption in general.

Returning to the prior of photo-z PDFs, the true redshift distribution satisfies the tautology $N'(z) = p(z|I_D)$, because our training data is representative, I_D represents a prior that is also equal to the truth. In this ideal case of representative training data, the method that would give the best approximation to $N'(z)$ would be one that neglects all the information contained in the photometry $\{d_i\}_{N_{tot}}$ and gives every galaxy the same photo-z PDF $\hat{p}_i(z) = N'(z)$ for all i . (In fact, including any information from the photometry would only add noise to the optimal result of returning the prior for every galaxy.) This is the exact estimator, TRAINZ, that we have described in Section 3.3, and which will serve as an experimental control.

The exact implementation of the stacked estimator $\hat{N}^H(z)$ should depend on the parametrization of the photo-z PDFs, which may differ across codes and can affect the precision of the estimator (Malz et al. 2018); even considering a single method under the same parametrization, say a piecewise constant function over bins or a set of samples from the posterior, an estimator using $2N$ bins or samples will trivially be more precise than an estimator using N bins or samples.

In order to minimize the effects of the choice of parameterization, we asked those running all twelve codes to output photo-z PDFs parameterized with ≈ 200 piecewise constant bins spanning $0 < z < 2$. The piecewise constant format is chosen because of its established presence in the literature, and the choice of 200 bins was motivated by the approximate number of columns expected to be available for storage of

photo- z PDFs for the final LSST Project tables.¹⁷ All the photo- z PDF catalogs are processed using the `qp` software package (Malz et al. 2018)¹⁸ for manipulating and calculating metrics of univariate PDFs. We will discuss the choice of photo- z PDF parameterization further in Section 5.

4.1 Metrics of an ensemble of photo- z posteriors

Because the photo- z PDFs in the LSST catalog will be used for many applications, some of which require accuracy of each individual catalog entry, we consider several metrics that get at the population-level performance of the photo- z PDFs as distinct from a summary statistic thereof.

4.1.1 Probability integral transform (PIT)

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

$$\text{PIT} = \int_{-\infty}^{z_{\text{true}}} p(z) dz. \quad (16)$$

The distribution of PIT values quantifies the behavior of the ensemble of photo- z PDFs, enabling us to evaluate whether the population of photo- z PDFs is, on average, accurate. The PIT value for each galaxy is the Cumulative Distribution Function (CDF) of its photo- z PDF evaluated at its true redshift. A catalogue of photo- z PDFs that are accurate should have a flat PIT histogram (i.e., the individual PIT values as samples from each CDF should match a Uniform(0,1) distribution if the CDFs are accurate). Specific deviations from flatness indicate inaccuracy: overly broad photo- z PDFs would manifest as underrepresentation of the lowest and highest PIT values, whereas overly narrow photo- z PDFs would manifest as over-representation of the lowest and highest PIT values. High frequency at only PIT ≈ 0 and PIT ≈ 1 indicates the presence of catastrophic outliers with highly inaccurate photo- z PDFs where the true redshift of a galaxy is outside of the support of its photo- z PDF. Tanaka et al. (2018) use the histogram of PIT values as a diagnostic indicator of overall code performance, while Freeman et al. (2017) independently define the PIT and demonstrate how its individual values may be used both to perform hypothesis testing (via, e.g., the KS, CvM, and AD tests; see below) and to construct quantile-quantile plots.

4.1.2 Quantile-quantile (QQ) plot

A quantile is defined by partitioning a distribution into consecutive intervals containing equal amounts of probability, or equal numbers of objects in each interval in the case of a distribution of objects. The quantile-quantile (QQ) plot serves as a graphical visualization for comparing two distributions, where the quantiles of one distribution are plotted against the quantiles of the other distribution. The QQ plot provides an easy way to qualitatively assess the differences in various properties such as the moments of an estimating distribution relative to a true distribution.

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

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

In this paper, QQ plots are used for two purposes: (1) for comparing $N(z)$ from photo- z PDFs (estimated using Eq. 15) with the true $N(z)$, and (2) for assessing the overall consistency of an ensemble of photo- z PDFs with their true redshifts on a population level, where the distribution of the PIT values (see previous section) is compared to a uniform distribution between 0 and 1. Though the QQ plot contains very similar information to that shown in the PIT histogram plot, due to being the PIT being the derivative of the QQ, we include both forms for completeness and enhanced visual interpretability.

4.1.3 Conditional density estimation loss

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

The CDE loss is defined as

$$L(f, \hat{f}) \equiv \int \int (f(z | \mathbf{x}) - \hat{f}(z | \mathbf{x}))^2 dz dP(\mathbf{x}). \quad (17)$$

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

$$L(f, \hat{f}) = \mathbb{E}_{\mathbf{X}} \left[\int \hat{f}(z | \mathbf{X})^2 dz \right] - 2\mathbb{E}_{\mathbf{X}, Z} \left[\hat{f}(Z | \mathbf{X}) \right] + K_f, \quad (18)$$

where upper-case letters denote random variables and lower-case the observed variables. The first expectation is with respect to the marginal distribution of the covariates \mathbf{X} , the second expectation is with respect to the joint distribution of \mathbf{X} and Z , and K_f is a constant depending only upon the true conditional densities $f(z | \mathbf{x})$. For each method we can estimate these expectations as empirical expectations on the test or validation data (Eq. 7 in Izbicki et al. 2017) without knowledge of the true densities.

4.2 Metrics over estimated probability distributions

In tandem with the QQ and PIT metrics introduced above, we additionally compute the following metrics comparing the empirical CDF of a distribution to the true or expected distribution. These metrics give a more quantitative measure of the departure from ideal than the more visual PIT histogram and QQ plot. We compute metrics comparing the CDF of PIT values to the CDF of a Uniform distribution, and also compute the CDF of the true redshift distribution $N'(z)$ compared the $\hat{N}(z)$ distribution derived from summing the photo- z PDFs as described in Eq. 15.

4.2.1 Root-mean-square error (RMSE)

We employ the familiar root-mean-square error

$$\text{RMSE} \equiv \sqrt{\int_{-\infty}^{\infty} (\hat{f}(z) - f'(z))^2 dz}. \quad (19)$$

Though this metric does not account for the fact that the redshift distribution function is, in fact, a probability distribution, it can still be interpreted as a measure of the integrated difference between the estimated distribution and the true distribution.

4.2.2 Kolmogorov-Smirnov (KS) and related statistics

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

$$N_{\text{KS}} \equiv \max_z (|F_{\text{phot}}(z) - F_{\text{spec}}(z)|). \quad (20)$$

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

$$\omega^2 \equiv \int_{-\infty}^{+\infty} (F_{\text{meas.}}(x) - F_{\text{ideal}}(x))^2 dF_{\text{ideal}} \quad (21)$$

is the average of the distance squared.

The AD statistic

$$A^2 \equiv N_{\text{tot}} \int_{-\infty}^{+\infty} \frac{(F_{\text{meas.}}(x) - F_{\text{ideal}}(x))^2}{F_{\text{ideal}}(x)(1 - F_{\text{ideal}}(x))} dF_{\text{ideal}} \quad (22)$$

is a weighted version of the CvM statistic, making it more sensitive to the tails of the distribution, where N_{tot} is the sample size.

4.2.3 Moments

We additionally calculate the first three moments of the estimated redshift distribution $\hat{N}^H(z)$ 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

$$\langle z^m \rangle \equiv \int_{-\infty}^{\infty} z^m N(z) dz. \quad (23)$$

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 twelve 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. For example, in CMNN the $p(z)$ are a simply a weighted histogram of all spectroscopic training galaxies in nearby colour space with no smoothing applied, so the substructure is due to the finite number of neighbours, and is not unexpected. 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 twelve codes, giving a grid size of $\delta z = 0.01$ 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 in a limited number of cases, for example bright galaxies with very narrow $p(z)$ where the true peak falls between grid points. We will discuss PDF storage in Section 6.

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.

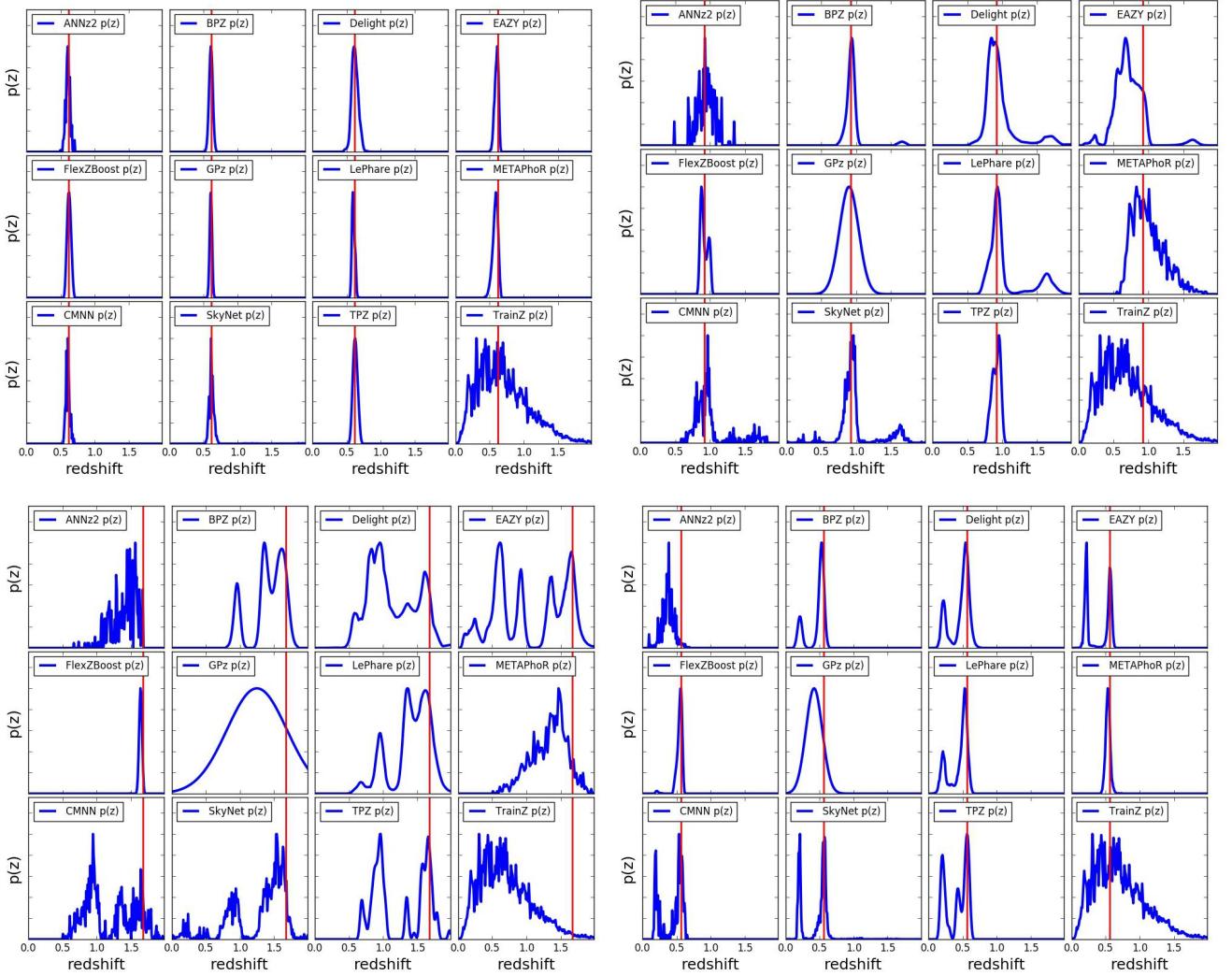


Figure 1. Four illustrative examples of individual $p(z)$ distributions produced by the codes. The red vertical line represents the true redshift. Examples are chosen with common features seen in PDFs: tight unimodal $p(z)$ (upper left), broad unimodal $p(z)$ (upper right), bimodal $p(z)$ (lower right), 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.

Examination of the PIT histograms and QQ plots shows that there are fairly generic issues with the width of $p(z)$ uncertainties: DELIGHT, CMNN, SKYNET and TPZ all show a PIT histogram with a 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

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 $PIT = 0.0$ 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 $PIT < 0.0001$ or $PIT > 0.9999$. Table 2 lists these fractions for each of

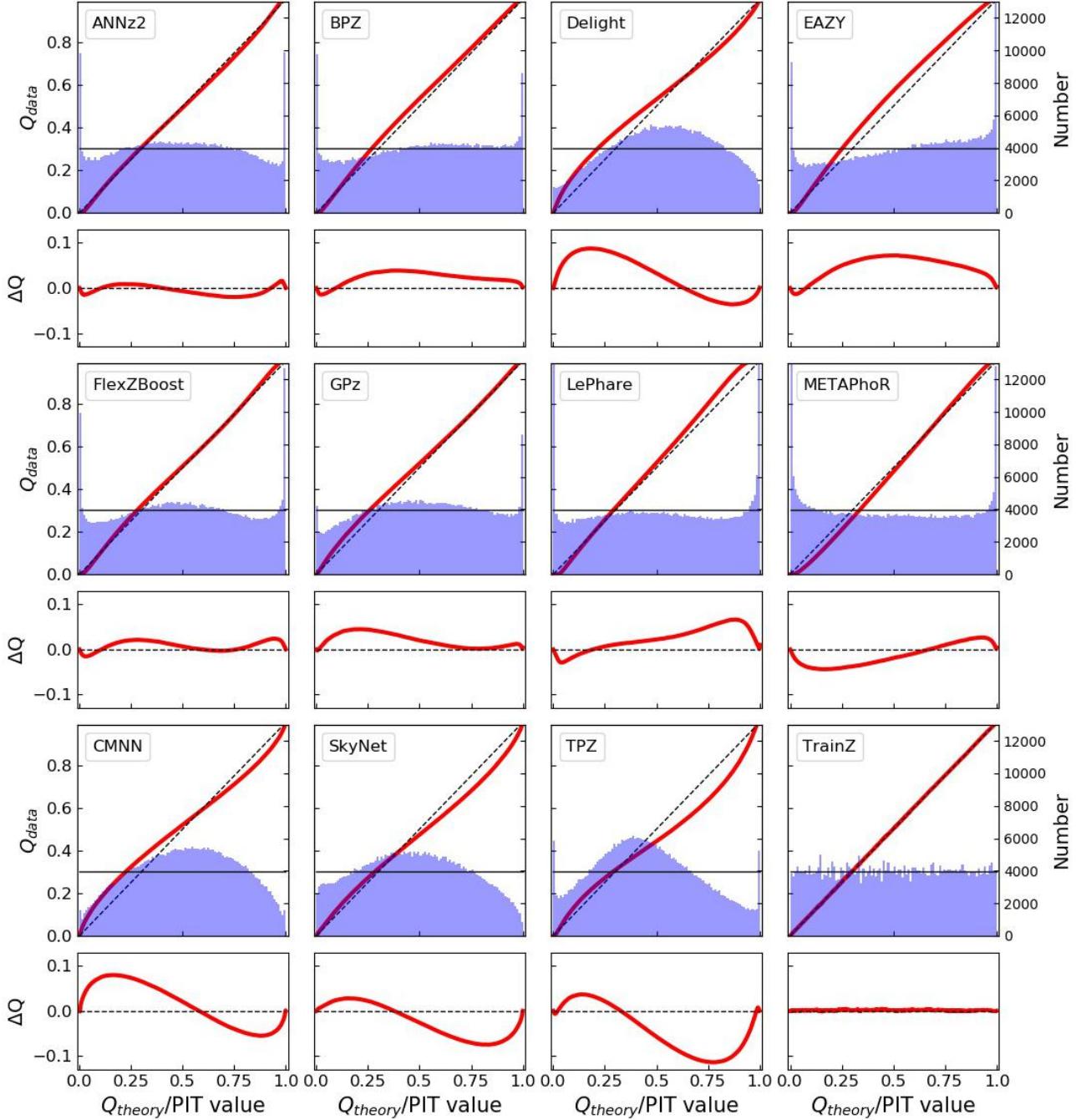


Figure 2. Summary plots for all twelve photo-z codes illustrating performance for the interim posterior statistics. The top panel of each pair shows both the Quantile-Quantile (QQ) plot (red) and the histogram of PIT values (blue). The desired behavior is a QQ plot that matches the diagonal dashed line, 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.

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	“catastrophic outlier” PIT fraction
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

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 $f_0 > 0.02$, 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_0 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 $[0,1]$. 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 extreme, 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 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.

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. All distributions are smoothed via kernel density estimation (KDE) with a common bandwidth chosen via Scott’s rule (Scott 1992) in order to minimize differences in small-scale features and make for a more uniform comparison between codes. 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. As expected, TRAINZ is in excellent agreement with the true redshift distribution: as the training sample is selected from the same underlying distribution as the test set, the redshift distributions are identical, up to Poisson fluctuations due to the finite number of sample galaxies. CMNN is also in excellent agreement for similar reasons: with a representative training sample of galaxies spanning the colour-space, the sum of the colour-matched neighbour redshifts should return the true redshift distribution. FLEXZBOOST and TPZ also show very good agreement, with only slight departures, with an over/under-prediction in the high redshift tail of $\hat{N}(z)$ evident around $z \sim 1.4$. In fact, several of the other 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 Å break passing through the gap between the z and y filters, resulting in a drastic change in $z - y$ colour for galaxies in this redshift range. With a relative dearth of strong features blue-ward of the 4000 Å break in most galaxy SEDs, the colour change in the two reddest filter bandpasses of a survey has a large influence on the redshift determination. The $z \sim 1.4$ feature is one of the most prominent sources of larger uncertainty in individual galaxy $p(z)$. In our sample individual galaxy $p(z)$ ’s tend to be broader around $z \sim 1.4$ and point estimates are more uncertain in this regime, as is readily seen in the point-estimate plots shown in Fig. A1 and described in the Appendix. This feature 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.

Two of the machine learning based codes, ANNz2 and METAPHOR, appear to be over-trained, adding excess galaxy probability to the redshift peaks with the largest number of training galaxies, and missing probability in the troughs where training galaxies are of fewer number. Given that our training data is drawn from the same galaxy pop-

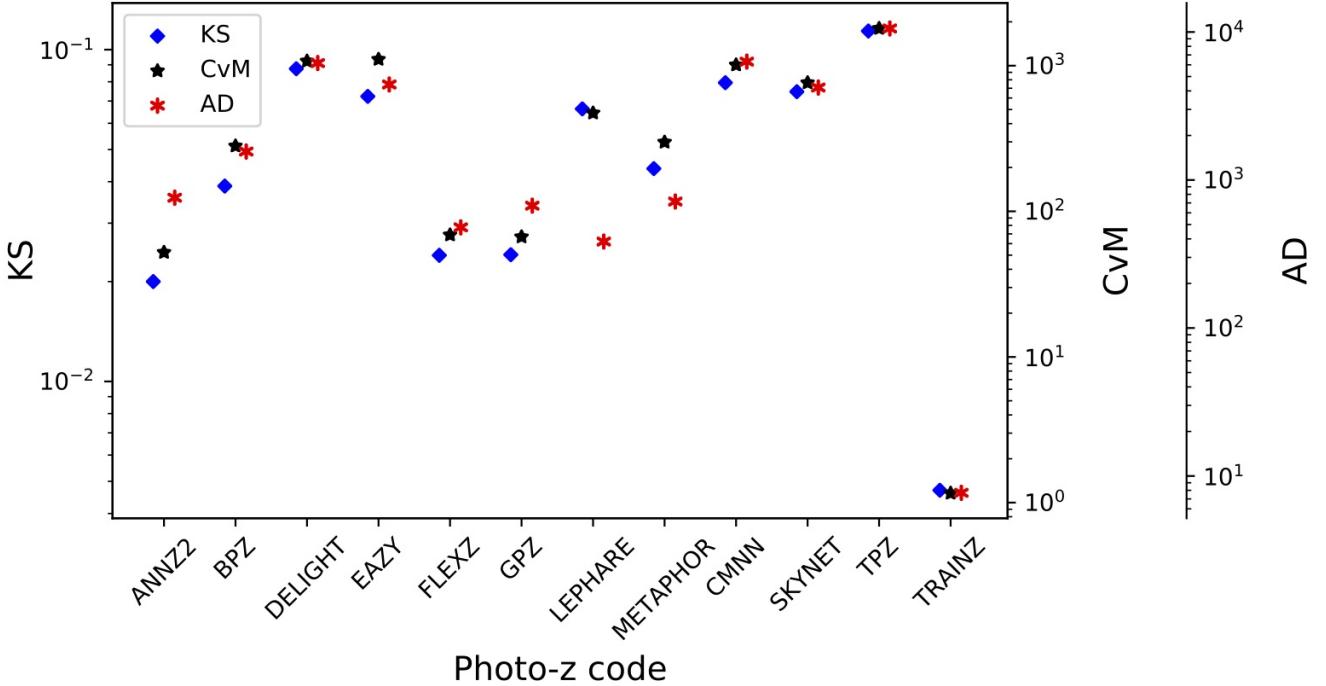


Figure 3. A visual representation of the Kolmogorov-Smirnov (KS, blue diamond), Cramer-von Mises (CvM, black star), and Anderson-Darling (AD, red asterisk) statistics for the PIT distributions. The statistics are often highly correlated, with similar relative values between metrics for each code; however, the AD statistic excludes some values at the extrema of the distributions, and can have disparate values compared to KS and CvM.

ulation as the test set, and our data has prominent peaks in $N'(z)$, perhaps it is not unexpected that such overtraining occurs in some codes, though the fact that it does not occur in all training-based codes indicates that it may be due to specifics of the implementations and bandwidth choices of ANN22 and METAPHOR. This, once again, emphasizes that care much be taken in choice of bandwidth parameters in individual codes in order to avoid such overtraining. 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. Unlike the previous implementation of SKYNET in 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-Smirnov (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.

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

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 $N'(z)$.

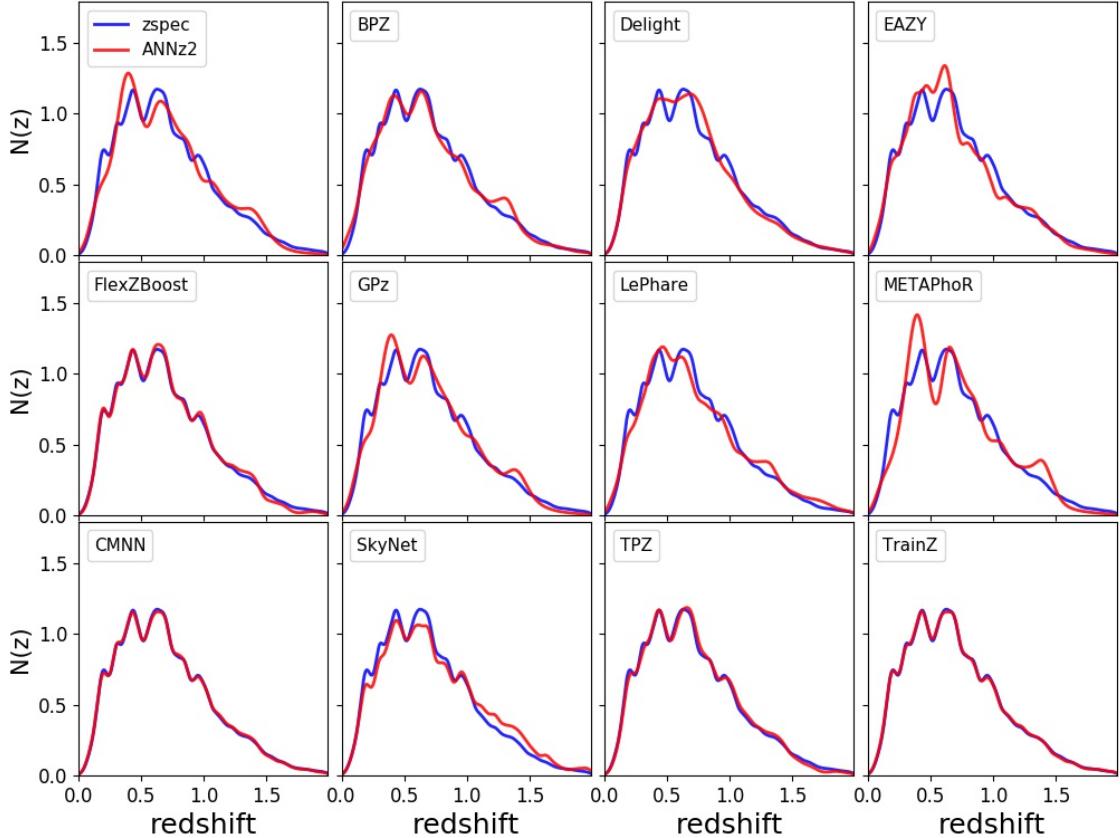


Figure 4. The stacked $p(z)$ produced by each photo-z code ($\hat{N}(z)$, red) compared to the spectroscopic redshift distribution ($N'(z)$, blue). Varying levels of agreement are seen in the codes. Both $\hat{N}(z)$ and $N'(z)$ in all codes are smoothed using a single bandwidth chosen via Scott's rule.

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. This implies that we may need multiple photo-z estimators, tuned to the particular metric, in order to maximize science returns in upcoming surveys. 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.

5.3 Interpretation of metrics

Samples from accurate photo-z posteriors should reproduce the space of $p(z, \text{data})$. However, it is difficult to test this reconstruction given our data set, as the galaxy distributions arise from mock objects pasted on to an underlying dark matter halo catalogue with properties designed to match empirical relations, rather than being drawn from statistical distributions in redshift. In previous sections we have mentioned that optimizing for a specific metric does not guarantee good performance on other metrics, nor is there any guarantee that good performance by our metrics corre-

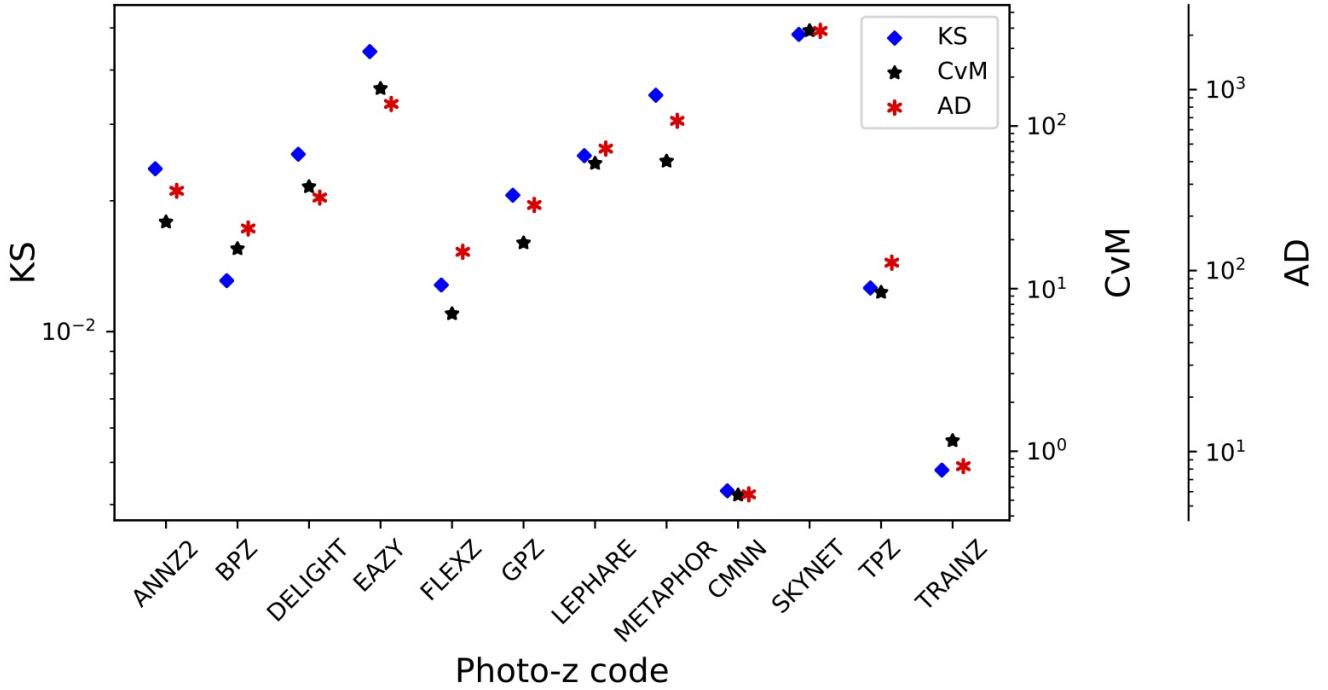


Figure 5. A visual representation of the Kolmogorov-Smirnov (KS, blue diamond), Cramer-von Mises (CvM, black star), and Anderson-Darling (AD, red asterisk) statistics for the $\hat{N}(z)$ distributions. The statistics are correlated, the codes with the lowest KS statistics tend to have the lowest CvM 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.

sponds to *accurate* photo- z posteriors, in the sense of predicting redshifts for individual galaxies. In other words, we can construct photo- z estimators that provide good coverage in many of our tests, but which have very little predictive power.

The TRAINZ estimator, which assigns every galaxy a $p(z)$ equal to $N(z)$ of the training set as described in Section 3.3, is introduced as a “null test” to demonstrate this point via *reductio ad absurdum*. TRAINZ outperforms all codes on the PIT-based metrics, and all but one code on the $N(z)$ based statistics. Because our training set is perfectly 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 presented in this paper. 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

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

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.

6 DISCUSSION AND FUTURE WORK

In this paper we presented results evaluating the computation of individual galaxy photometric redshift PDFs for twelve photo- z codes. As discussed in Section 4 each $p(z)$

Table 4. Root-Mean-Square-Error (RMSE) statistics for the twelve 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

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
METAPHOR	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

should accurately reflect the relative likelihood as a function of redshift for each galaxy in an informative way; that is, the estimates should provide useful information per individual galaxy, not just the ensemble. 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, or other complications 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 of a problem with either/both of the specific code implementation or the underlying algorithm. In order to meet the stringent LSST goals for photo-z performance, identifying and correcting such problems is an important first step before tackling more realistic data in

future challenges, and codes that do not perform well may not be worth pursuing in future challenges. Many of the codes tested performed well; however, several did not meet the stringent goals that have been laid out for LSST photometric redshift performance for individual galaxies, as laid out in the LSST SRD (See Section 1). This is a cause for concern, given the idealized conditions, and the individual code responses will be studied in detail moving forward. If methods can not reach the goals on idealized data, then they will almost surely not meet those same goals when the more complex problems that we expect to arise from real LSST data are included. The results presented in this paper enable an evaluation of which algorithms are the most promising moving forward, and potentially point to implementation choices or mistakes which could be improved or corrected in others.

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 $\hat{N}(z)$ 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, including propagation of redshift uncertainties in a set of fiducial tomographic redshift bins in order to estimate impact on cosmological parameter estimation.

In future papers a focus of the *LSST Dark Energy Science Collaboration Photo-z Working Group* will be to 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. The SEDs for the galaxy sample in this paper were constructed from linear combinations of five basis SED templates. Future simulations will also include more complex SED information, with a more realistic range of physical properties, and the inclusion of AGN effects, 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.

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. Current and upcoming surveys are putting 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. We plan 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.

The underlying dataset limited this work 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.

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, which will naturally include the complications of object blending, sensor effects, different observing conditions, amongst others. 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.

Finally, while this paper and future papers discussed above focus on photometric redshift codes and estimating accurate $p(z)$ from training data, we plan a separate, but complementary, project to examine calibration of the resultant redshifts via spatial cross-correlations (Newman 2008), which will be explored in a separate set of papers. The overarching plan describing everything laid out in this section is described in more detail in the LSST DESC Science Roadmap (see Footnote in Section 1). These plans will require significant effort, but they are necessary if we are to make optimal use of the LSST data for astrophysical and cosmological analyses.

Acknowledgments

Author contributions are listed below.

S.J. Schmidt: Led the project. (conceptualization, data curation, formal analysis, investigation, methodology, project administration, resources, software, supervision, visualization, writing – original draft, writing – review & editing)

A.I. Malz: Contributed to choice of metrics, implementation in code, and writing. (conceptualization, methodology, project administration, resources, software, visualization, writing – original draft, writing – review & editing)

J.Y.H. Soo: Ran ANNz2 and Delight, updated abstract, edited sections 1 through 6, added tables in Methods and Results, updated references.bib and added references throughout the paper

M. Brescia: main ideator of METAPHOR and of MLPQNA; modification of METAPHOR pipeline to fit the LSST data structure and requirements

S. Cavaudi: Contributed to choice and test of metrics, ran

1633 METAPHOR, minor text editing
 1634 G. Longo: Scientific advise, test and validation of the modified METAPHOR pipeline, text of the METAPHOR section
 1635 I.A. Almosallam: vetted the early versions of the data set and ran many photo-z codes on it, applied GPz to the final version and wrote the GPz subsection
 1636 M.L. Graham: Ran the colour-matched nearest-neighbours photo-z code on the Buzzard catalog and wrote the relevant piece of Section 2; participated in discussions of the analysis.
 1637 A.J. Connolly: Developed the colour-matched nearest-neighbours photo-z code; participated in discussions of the analysis.
 1638 E. Nourbakhsh: Ran and optimized TPZ code on the Buzzard catalog and wrote a subsection of Section 2 for that
 1639 J. Cohen-Tanugi: contributed to running code, analysis discussion, and editing,reviewing the paper
 1640 H. Tratin: contributed to providing SkyNet results and writing the relevant section
 1641 P.E. Freeman: Contributed to choice of CDE metrics and to implementation of FlexZBoost
 1642 K. Iyer: assisted in writing metric functions used to evaluate codes
 1643 J.B. Kalmbach: Worked on preparing the figures for the paper.
 1644 E. Kovacs: Ran simulations, discussed data format and properties for SEDs, dust, and ELG corrections
 1645 A.B. Lee: Co-developed FlexZBoost and the CDE loss statistic, wrote text on the work, and supervised the development of FlexZBoost software packages
 1646 C. Morrison: Managerial support; Discussions with authors regarding metrics and style; Some coding contribution to metric computation.
 1647 J. Newman: Contributions to overall strategy, design of metrics, and supervision of work done by Rongpu Zhou
 1648 E. Nuss: contributed to running code, analysis discussion, and editing,reviewing the paper
 1649 T. Pospisil: Co-developed FlexZBoost software and CDE loss calculation code
 1650 M.J. Jarvis: Contributed text on AGN to Discussion section and portions of GPz work
 1651 R. Izbicki: Co-developed FlexZBoost and the CDE loss statistic, and wrote software for FlexZBoost

1652 The authors would like to thank their LSST-DESC publication review committee for comments that improved the paper draft.

1653 personal funding sources

1654 AIM is advised by David W. Hogg and was supported by National Science Foundation grant AST-1517237.

1655 In addition to packages cited in the text, analyses performed in this paper used the following software packages: 1656 NUMPY and SCIPY (Oliphant 2007), MATPLOTLIB (Hunter 2007), SEABORN (Waskom et al. 2017), MINFUNC (Schmidt 2005), PYSKYNET (Bonnett 2016), and PhotUtils from the 1657 LSST simulations package (Connolly et al. 2014).

1658 The DESC acknowledges ongoing support from the Institut National de Physique Nucléaire et de Physique des Particules in France; the Science & Technology Facilities Council in the United Kingdom; and the Department of En-

ergy, the National Science Foundation, and the LSST Corporation in the United States. DESC uses resources of the IN2P3 Computing Center (CC-IN2P3-Lyon/Villeurbanne - France) funded by the Centre National de la Recherche Scientifique; the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231; STFC DiRAC HPC Facilities, funded by UK BIS National E-infrastructure capital grants; and the UK particle physics grid, supported by the GridPP Collaboration. This work was performed in part under DOE Contract DE-AC02-76SF00515.

APPENDIX A: POINT ESTIMATE PHOTOMETRIC REDSHIFTS

While we assume that all science analysis will use full PDF information and do not recommend the use of single point estimates of redshift for most science applications, we include a brief evaluation as an Appendix. 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 historical comparisons from older 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 $0.05 \times z_{PEAK}$ 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.

1754 A1 Point Estimate Metrics

1755 We calculate the commonly used point estimate metrics of
 1756 the overall photo- z scatter (σ_z , the standard deviation of
 1757 the photo- z residuals), bias, and “catastrophic outlier rate”.
 1758 Specifically, we calculate the metrics as follows: we define e_z
 1759 as

$$1760 e_z = \frac{z_P - z_S}{1 + z_S} \quad (A1)$$

1761 where z_P is the point estimate and z_S is the true redshift. In
 1762 practice, because the standard deviation calculation is quite
 1763 sensitive to the outliers, we define the photo- z scatter, σ in
 1764 terms of the Interquartile Range (IQR), the difference be-
 1765 tween the 75th and 25th percentiles of the e_z distribution.
 1766 In order to match the usual meaning of a 1σ interval, we
 1767 scale the IQR and define $\sigma_{IQR} = IQR/1.349$, as there is a
 1768 factor of 1.349 difference between the IQR and the standard
 1769 deviation of a Normal distribution. While many other stud-
 1770 ies define the bias based on the *mean* offset between true
 1771 and estimated redshift, in this study we define the bias as
 1772 the median value of e_z for the sample. We use median as
 1773 it is, once again, less sensitive to outliers than the mean.
 1774 The catastrophic outlier fraction is defined as the fraction
 1775 of galaxies with e_z greater than the *larger* of $3\sigma_{IQR}$ or 0.06,
 1776 i.e. 3σ outliers with a floor of $\sigma_{IQR}=0.02$. For reference, the
 1777 goals stated in Section 3.8 of the LSST Science Book (Abell
 1778 et al. 2009) for photo- z performance in these metrics, as-
 1779 suming perfect training knowledge (as we are testing in this
 1780 paper) are:

- 1781 • RMS scatter < $0.02(1+z)$
- 1782 • bias < 0.003
- 1783 • catastrophic outlier rate < 10%

1784 These definitions are similar, but not exactly the same, as
 1785 the σ_{IQR} and median bias calculated here, but are similar
 1786 enough for qualitative comparisons to the LSST goals.

1787 Fig. A1 shows the point estimates for both z_{PEAK} and
 1788 z_{WEIGHT} . Point density is shown with mixed contours to
 1789 emphasize that most of the galaxies do fall close to the
 1790 $z_{phot} = z_{spec}$ line, while blue points show differing char-
 1791 acteristics of the outlier populations. The red dashed lines
 1792 indicated the cutoff for catastrophic outliers, defined as:
 1793 $\max(0.06, 3\sigma_{IQR})$. As with the full $p(z)$ results, a variety
 1794 of behaviours are evident in the different codes. Table A1
 1795 lists the scatter, bias, and catastrophic outlier fractions for
 1796 the codes. The performance of the codes for point met-
 1797 rics is highly correlated with performance on $p(z)$ based
 1798 tests, which is to be expected, given that the point-estimates
 1799 were derived from the $p(z)$. Some discretization is evident in
 1800 z_{PEAK} , particularly for SKYNET, due to the finite grid spac-
 1801 ing of the reported $p(z)$. These discreteness effects are miti-
 1802 gated by the weighting of z_{WEIGHT} , resulting in a smoother
 1803 distribution of redshift estimates. Several features perpen-
 1804 dicular to the main $z_{phot} = z_{spec}$ line are evident. These fea-
 1805 tures are due to the 4000 angstrom break passing through
 1806 the gaps between adjacent LSST filters. These features are
 1807 most prominent in template-based codes, but appear to
 1808 some degree in all codes tested.

1809 In even the best performing codes, there are visible oc-
 1810 cupied regions away from the $z_{phot} = z_{spec}$ line, correspond-
 1811 ing to degenerate redshift solutions for certain LSST magni-

1812 tudes and colors. While use of the full information available
 1813 via $p(z)$ mitigates their impact, a full understanding of the
 1814 outlier population is critical for LSST science, particularly
 1815 in tomographic applications

1816 Finally, we note that all twelve codes perform at or near
 1817 the goals for point-estimates as outlined in the LSST Science
 1818 Requirements Document¹⁹ and Graham et al. (2018). This
 1819 is to be expected, given that the requirements were designed
 1820 such that a point estimate photo- z would meet these require-
 1821 ments for perfect training data to a depth of $i < 25$. But, it
 1822 is still an encouraging sign, given an updated mock galaxy
 1823 simulation and the expanded set of photo- z codes tested.

1824 REFERENCES

- 1825 Abbott T., et al., 2005, preprint (arXiv:astro-ph/0510346)
 1826 Abell P. A., et al., 2009, preprint (arXiv:0912.0201),
 1827 Aihara H., et al., 2018a, *PASJ*, **70**, S4
 1828 Aihara H., et al., 2018b, *PASJ*, **70**, S8
 1829 Almosallam I. A., Lindsay S. N., Jarvis M. J., Roberts S. J.,
 1830 2016a, *MNRAS*, **455**, 2387
 1831 Almosallam I. A., Jarvis M. J., Roberts S. J., 2016b, *MNRAS*,
 1832 **462**, 726
 1833 Arnouts S., Cristiani S., Moscardini L., Matarrese S., Lucchin F.,
 1834 Fontana A., Giallongo E., 1999, *MNRAS*, **310**, 540
 1835 Behroozi P. S., Wechsler R. H., Wu H.-Y., 2013, *ApJ*, **762**, 109
 1836 Benítez N., 2000, *ApJ*, **536**, 571
 1837 Bernstein G., Huterer D., 2010, *MNRAS*, **401**, 1399
 1838 Blanton M. R., Roweis S., 2007, *AJ*, **133**, 734
 1839 Blanton M. R., et al., 2005, *AJ*, **129**, 2562
 1840 Bonnett C., 2015, *MNRAS*, **449**, 1043
 1841 Bonnett C., 2016, Python wrapper to SkyNet, <https://pyskynet.readthedocs.io/en/latest/>
 1842 Bonnett C., et al., 2016, *Phys. Rev. D*, **94**, 042005
 1843 Brammer G. B., van Dokkum P. G., Coppi P., 2008, *ApJ*, **686**,
 1844 1503
 1845 Breiman L., Friedman J. H., Olshen R. A., Stone C. J., 1984, Clas-
 1846 sification and Regression Trees, Statistics/Probability Series.
 1847 Wadsworth Publishing Company, Belmont, California, U.S.A
 1848 Brescia M., Cavuoti S., Amaro V., Riccio G., Angora G., Vellucci
 1849 C., Longo G., 2018, preprint, ([arXiv:1802.07683](https://arxiv.org/abs/1802.07683))
 1850 Carrasco Kind M., Brunner R. J., 2013, *MNRAS*, **432**, 1483
 1851 Carrasco Kind M., Brunner R. J., 2014, *MNRAS*, **442**, 3380
 1852 Cavuoti S., Amaro V., Brescia M., Vellucci C., Tortora C., Longo
 1853 G., 2017, *MNRAS*, **465**, 1959
 1854 Chen T., Guestrin C., 2016, in Proceedings of the 22Nd ACM
 1855 SIGKDD International Conference on Knowledge Discovery
 1856 and Data Mining. KDD '16. ACM, New York, NY, USA,
 1857 pp 785–794, doi:[10.1145/2939672.2939785](https://doi.org/10.1145/2939672.2939785), <http://doi.acm.org/10.1145/2939672.2939785>
 1858 Connolly A. J., et al., 2014, in Angeli G. Z., Dierickx P., eds,
 1859 Society of Photo-Optical Instrumentation Engineers (SPIE)
 1860 Conference Series Vol. 9150, Modeling, Systems Engineer-
 1861 ing, and Project Management for Astronomy VI. p. 14,
 1862 doi:[10.1117/12.2054953](https://doi.org/10.1117/12.2054953)
 1863 Dahlen T., et al., 2013, *ApJ*, **775**, 93
 1864 Dawson W. A., Schneider M. D., Tyson J. A., Jee M. J., 2016,
 1865 *ApJ*, **816**, 11
 1866 Fernández-Soto A., Lanzetta K. M., Yahil A., 1999, *ApJ*, **513**, 34
 1867 Firth A. E., Lahav O., Somerville R. S., 2003, *MNRAS*, **339**, 1195
 1868 Freeman P. E., Izbicki R., Lee A. B., 2017, *MNRAS*, **468**, 4556
 1869 Graff P., Feroz F., Hobson M. P., Lasenby A., 2014, *MNRAS*, **441**,
 1870 1741
 1871

1872 ¹⁹ available at: <http://ls.st/srd>

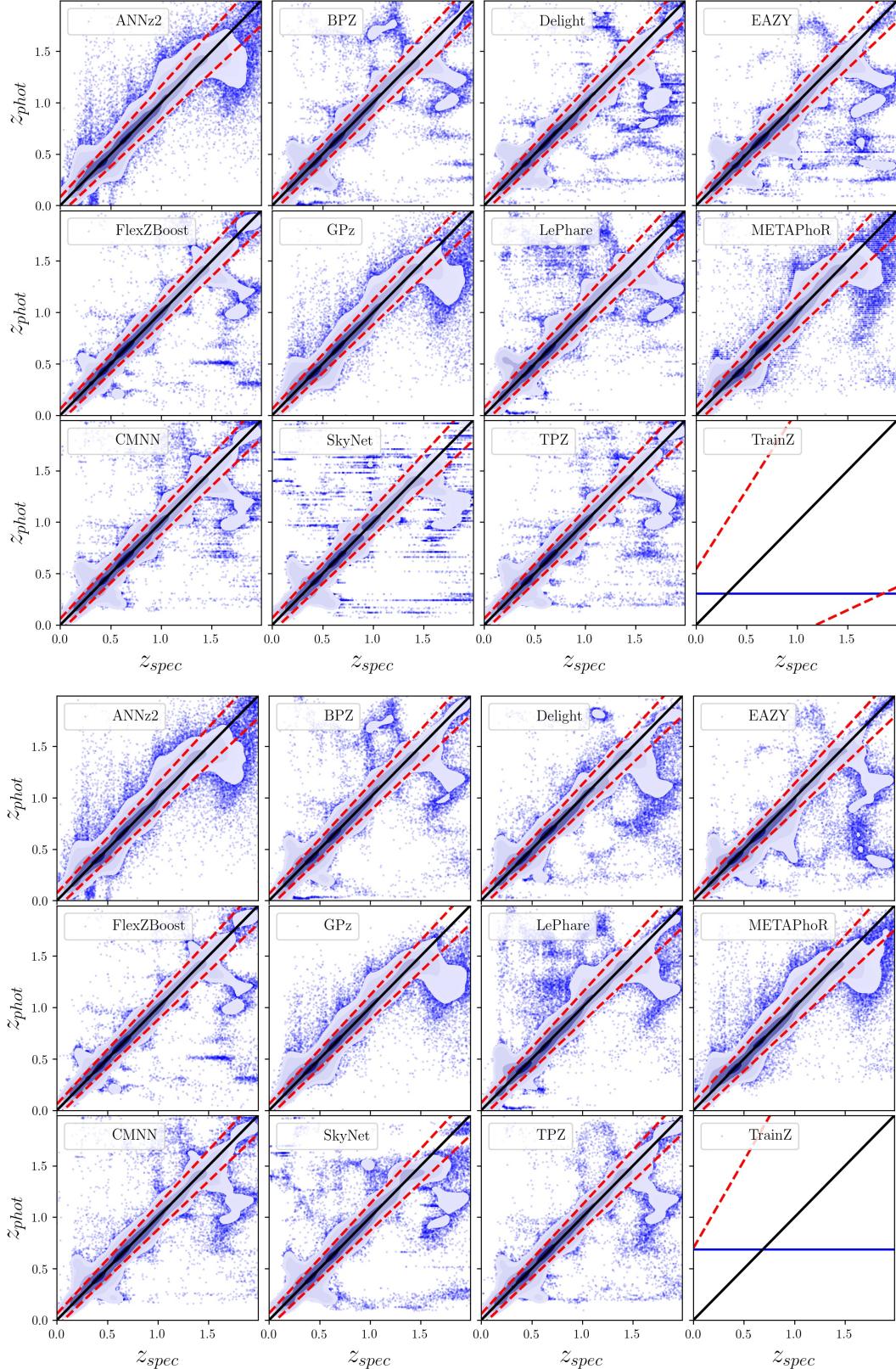


Figure A1. Point estimate photo-z's derived from the posteriors. Top panel shows z_{PEAK} , while bottom panel shows z_{WEIGHT} . Point estimate density is represented with fixed density contours, while outliers at lower density are represented by blue points. While use of point-estimate photo-z's is 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	<i>Z_{PEAK}</i>			<i>Z_{WEIGHT}</i>		
	$\frac{\sigma_{IQR}}{(1+z)}$	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

- 1873 Graham M. L., Connolly A. J., Ivezić Ž., Schmidt S. J., Jones 1919 preprint, ([arXiv:1809.01669](https://arxiv.org/abs/1809.01669))
 1874 R. L., Jurić M., Daniel S. F., Yoachim P., 2018, *AJ*, **155**, 1 1920 Waskom M., et al., 2017, [doi:10.5281/zenodo.824567](https://doi.org/10.5281/zenodo.824567)
 1875 Green J., et al., 2012, preprint (arXiv:1208.4012), 1921 York D. G., et al., 2000, *AJ*, **120**, 1579
 1876 Hildebrandt H., et al., 2010, *A&A*, **523**, A31 1922 de Jong J. T. A., Verdoes Kleijn G. A., Kuijken K. H., Valentijn
 1877 Hunter J. D., 2007, Matplotlib: A 2D Graphics Environment, 1923 E. A., 2013, *Exp. Astron.*, **35**, 25
 1878 [doi:10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55) 1924 de Jong J. T. A., et al., 2017, *A&A*, **604**, A134
 1879 Ilbert O., et al., 2006, *A&A*, **457**, 841
 1880 Ivezić Ž., et al., 2008, preprint (arXiv:0805.2366),
 1881 Izbicki R., Lee A. B., 2017, *Electron. J. Statist.*, **11**, 2800
 1882 Izbicki R., Lee A. B., Freeman P. E., 2017, *Ann. Appl. Stat.*, **11**, 1883 698
 1884 Laureijs R., et al., 2011, preprint (1110.3193),
 1885 Leistedt B., Hogg D. W., 2017, *ApJ*, **838**, 5
 1886 Malz A., Hogg D., in prep., CHIPPR, chippr
 1887 Malz A., Marshall P., DeRose J., Graham M., Schmidt S., Wechsler R., 2018, *AJ*, Accepted,
 1888 1889 Mandelbaum R., et al., 2008, *MNRAS*, **386**, 781
 1890 Massarotti M., Iovino A., Buzzoni A., 2001, *A&A*, **368**, 74
 1891 Masters D. C., Stern D. K., Cohen J. G., Capak P. L., Rhodes
 1892 J. D., Castander F. J., Paltani S., 2017, *ApJ*, **841**, 111
 1893 Newman J. A., 2008, *ApJ*, **684**, 88
 1894 Newman J. A., et al., 2015, *Astroparticle Physics*, **63**, 81
 1895 Oliphant T., 2007, Python for Scientific Computing,
 1896 [doi:10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55)
 1897 Oyaizu H., Lima M., Cunha C. E., Lin H., Frieman J., 2008, *ApJ*,
 1898 **689**, 709
 1899 Polsterer K. L., D'Isanto A., Gieseke F., 2016, preprint
 1900 (arXiv:1608.08016),
 1901 Rasmussen C., Williams C., 2006, Gaussian Processes for Machine
 1902 Learning. Adaptative computation and machine learning se-
 1903 ries, MIT Press, Cambridge, MA
 1904 Rau M. M., Seitz S., Brimioule F., Frank E., Friedrich O., Gruen
 1905 D., Hoyle B., 2015, *MNRAS*, **452**, 3710
 1906 Reddick R. M., Wechsler R. H., Tinker J. L., Behroozi P. S., 2013,
 1907 *ApJ*, **771**, 30
 1908 Sadeh I., Abdalla F. B., Lahav O., 2016, *PASP*, **128**, 104502
 1909 Sánchez C., Bernstein G. M., 2018, preprint, ([arXiv:1807.11873](https://arxiv.org/abs/1807.11873))
 1910 Sánchez C., et al., 2014, *MNRAS*, **445**, 1482
 1911 Schmidt M., 2005, minFunc: Unconstrained Differentiable Mul-
 1912 tivariate Optimization in Matlab, <http://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>
 1913
 1914 Scott D. W., 1992, Multivariate Density Estimation. Theory,
 1915 Practice, and Visualization. Wiley
 1916 Skrutskie M. F., et al., 2006, *AJ*, **131**, 1163
 1917 Tanaka M., et al., 2018, *PASJ*, **70**, S9
 1918 The LSST Dark Energy Science Collaboration et al., 2018,