

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

LSST-DESC Photometric Redshift Working Group

April 24, 2019

ABSTRACT

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

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

1 INTRODUCTION

- 2 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, [Abell et al. 2005](#)), the Kilo-Degree Survey (KiDS, [de Jong et al. 2013](#)), Large Synoptic Survey Telescope (LSST, [Abell et al. 2009](#)), Euclid ([Lau-reijs 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.
- 18 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 system-

atic errors are not to dominate the statistical errors. The LSST Science Requirements Document (SRD)¹ lists the photometric redshift goals for a magnitude limited sample with $i < 25$ as: root-mean-square error with a goal of $\sigma_z < 0.02(1+z)$; 3σ “catastrophic outlier” rate below 10%; bias below 0.003².

The tremendous size of LSST’s galaxy catalog will be enabled by its exceptional depth, pushing to fainter magnitudes and deeper imaging and including galaxies of lower luminosity and higher redshift than ever before. In addition to the contribution of low signal-to-noise photometry to the systematic error of photo- zs , these populations in-

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.

2 LSST-DESC Photometric Redshift Working Group

36 introduce major physical degeneracies, for example 100
 37 the Lyman break/Balmer break degeneracy, that 102
 38 were not present in the populations covered in 104
 39 previous large area surveys like the Sloan Digital 106
 40 Sky Survey (SDSS, York et al. 2000) and the Two 108
 41 Micron All Sky Survey (2MASS, Skrutskie et al.
 42 2006). In order to meet LSST’s demanding error 110
 43 budget, it will be necessary to fully characterize 112
 44 those degeneracies wherein multiple redshift solu- 114
 45 tions have comparable likelihood.

46 There is often a desire to have a single valued 110
 47 “point-estimate” redshift for an individual galaxy. 112
 48 However, the complex, non-linear (and often non- 114
 49 unique) nature of the mapping between broad 116
 50 band fluxes and redshift means that a single value 118
 51 is unable to capture the full redshift information 120
 52 encoded in a galaxy’s magnitudes. For example, 122
 53 a common point-estimate for a template-based 124
 54 method is taking the highest likelihood solution as 126
 55 the point photo- z . A single valued redshift ignores 128
 56 degenerate redshift solutions of lower probability, 130
 57 potentially biasing photometric redshift estimates 132
 58 both for individual galaxies and ensemble distribu- 134
 59 tions. Storing more information is necessary, most 136
 60 often photo- z codes output the redshift probability 138
 61 density function (PDF), also often referred to as 140
 62 $p(z)$, describing the relative likelihood as a func- 142
 63 tion of redshift. Early template methods such as 144
 64 Fernández-Soto et al. (1999) converted relative χ^2 146
 65 values of template spectra to likelihoods to esti- 148
 66 mate $p(z)$. Soon after, codes such as Benítez (2000) 150
 67 added a Bayesian prior and output a posterior 152
 68 probability distribution. While many early ma- 154
 69 chine learning based algorithms focused on a point- 156
 70 estimate, Firth et al. (2003) used a neural net with 158
 71 1000 realizations scattered within the photometric 160
 72 errors to estimate a $p(z)$. As more groups began to 162
 73 employ photometric redshifts in their cosmological 164
 74 analyses, realization that point-estimate photo- z ’s 166
 75 were inadequate for precision cosmology measure- 168
 76 ments (Mandelbaum et al. 2008). From around this 170
 77 point onward, most photo- z algorithms have at- 172
 78 tempted to implement some estimate of the over- 174
 79 all redshift probability in their outputs, and some 176
 80 surveys began supplying a full $p(z)$ rather than a 178
 81 simple redshift point-estimate and error (e. g. de 180
 82 Jong et al. 2017).

83 There are numerous techniques for deriv- 182
 84 ing photo- z PDFs from photometry, yet no one 184
 85 method has yet been established as clearly supe- 186
 86 rior. Quantitative comparisons of photo- z meth- 188
 87 ods have been made before. The Photo- z Accuracy 188
 88 And Testing (PHAT, Hildebrandt et al. 2010) ef- 190
 89 fort focused on point estimates derived from many 192
 90 photometric bands. DES compared several codes 194
 91 for point estimates (Sánchez et al. 2014) and a 196
 92 summary statistic of photo- z interim posteriors for 198
 93 tomographically binned galaxy subsamples (Bon- 198
 94 nett et al. 2016). This paper is distinguished by 198
 95 its inclusion of metrics of photo- z interim poste- 198
 96 riors themselves and consideration of both classic 198
 97 and state-of-the-art photo- z algorithms, compar- 198
 98 ing the performance of several of the most widely 198
 99 employed codes as well as some that have been 198

100 developed only recently on the basis of metrics 102
 101 appropriate for a probabilistic data product. The 103 results presented in this work are a major fo- 105
 104 cus of the Photometric Redshift working group 106
 105 of the LSST Dark Energy Science Collaboration 108
 106 (LSST-DESC). This work is laid out in the Science 110
 107 Roadmap (SRM)³ as one of the critical activities 112
 109 to be completed in preparation for dark energy sci- 114
 110 ence analysis on the first year LSST data. This is 116
 111 the first of multiple papers by the working group, 118
 112 which will grow in sophistication. In this initial 120
 113 paper we focus on evaluating the performance of 122
 114 photometric redshift codes and their ability to pro- 124
 116 duce accurate PDFs in the presence of representa- 126
 118 tive training sets. This can be thought of as an 128
 119 initial test under near perfect conditions, before 130
 120 further complexities are added in future papers. 132
 122 Comparing the relative performance of the codes 134
 124 enables us to evaluate whether each code is us- 136
 126 ing information in an optimal way, and may re- 138
 128 veal enhancements in some codes and deficiencies 140
 130 in others, either in the fundamental algorithm, or 142
 132 in specific implementation.

144 Certain science cases need redshift informa- 146
 145 tion on individual objects, e. g. identification of 147
 146 host galaxy redshift for supernova classification, 148
 147 or identifying potential cluster membership. Other 149
 148 science cases need only ensemble redshift informa- 150
 150 tion; for instance current weak lensing techniques 152
 151 require the overall redshift distribution $N(z)$ for 153
 152 tomographic redshift samples, but do not need sin- 154
 153 gle galaxy estimates. In the case of the multiple 155
 154 types of probes of cosmology enabled by the LSST 156
 155 cosmology sample of several billion galaxies, the 157
 156 number of redshift bins and their photo- z require- 158
 157 ments vary with the specific probe; 2-point angular 159
 158 correlations benefit from many bins, while weak 160
 159 lensing probes do not (due to the wide lensing 161
 160 kernel). Large photometric surveys such as LSST 162
 161 must develop algorithms that meet the needs of 163
 160 all such science cases. In order to meet these am- 164
 161 bitious goals for photo- z accuracy, every aspect of 165
 162 photo- z estimation will have to be optimized: the 166
 163 algorithms employed, both template and machine- 167
 164 learning based (both in design and implementa- 168
 165 tion); the spectroscopic data used as a training 169
 166 set for machine learning algorithms or to estimate 170
 167 template sets and train Bayesian priors; and prob- 171
 168 abilistic catalog compression schemes that balance 172
 169 information retention against limited storage 173
 170 resources.

171 Before moving forward, we must address how 172
 173 the best methods may be unique to the per- 174
 174 formance metrics and science cases considered 175
 175 and what distinguishes photo- z PDFs of different 176
 176 methods from one another. Though photo- z PDFs 177
 177 are often written simply as $p(z)$, the PDF itself 178
 178 must be an interim posterior distribution $p(z|d, I)$, 179
 179 the probability of redshift conditioned on photo- 180
 180 metric data d that has actually been observed and 181

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

160 the prior information I that guides how a redshift
 161 is extracted from the photometry. If we run mul-
 162 tiple photo- z codes on a single dataset, the photo-
 163 z interim posteriors will not be identical because
 164 each code is based on assumptions in the form of
 165 an interim prior — these assumptions form the
 166 premise for photo- z estimation as a whole and are
 167 the only way to introduce differences in estimates
 168 of what would otherwise be a shared photo- z pos-
 169 terior $p(z|d)$ regardless of the code used to obtain
 170 it. Though explicit knowledge of the interim prior
 171 is necessary to use photo- z interim posteriors self-
 172 consistently in physical inference, the interim prior
 173 of a particular methodology is often implicit and
 174 not necessarily shared among all galaxies in the
 catalog.

175 This paper therefore aims (1) to constrain the
 176 impact of the interim prior I by separating it into a
 177 component I_H representing the method itself and
 178 a component I_D representing physical information,
 179 such as a training set or SED template library and
 180 (2) to present a procedure for evaluating the per-
 181 formance of photo- z codes in generic tests that
 182 may include many more systematics in the interim
 183 prior I . In order to isolate the effects encapsulated
 184 by I_H of variation between codes from issues with
 185 the training set or template library encapsulated
 186 by I_D , we use an identical set of simulated galax-
 187 ies for every code and construct a template library
 188 and training sample that are *complete and rep-
 189 resentative* and shared among all codes; that is,
 190 our training sample for machine learning codes is
 191 drawn from the same underlying galaxy popula-
 192 tion as our test set, with no additional selections,
 193 and the SED library used for template-based codes
 194 is the same as the one used to generate the photo-
 195 metric data. We explore a number of performance
 196 metrics in this paper, not to make a conclusion re-
 197 garding the superiority or even relative favorability
 198 of each code but to establish a method for compar-
 199 ing photo- z PDFs derived by different methods.
 200 These test conditions set the stage for addressing
 201 in a future paper the crucial issue of incomplete
 and non-representative prior information.

202 The outline of the paper is as follows: in § 2
 203 we present the simulated data set; in § 3 we de-
 204 scribe the current generation codes employed in
 205 the paper; in § 4 we discuss the interpretation of
 206 photo- z PDFs in terms of metrics of accuracy; in
 207 § 5 we show our results and compare the per-
 208 formance of the codes; in § 6 we offer our conclusions
 209 and discuss future extensions of this work.

212 2 THE SIMULATION AND MOCK 213 GALAXY CATALOG

214 In order to test the current generation codes, we
 215 employ a simulated galaxy catalogue. The simula-
 216 tion is completely catalogue-based, with no image
 217 construction or mock measurements made. We de-
 218 scribe these in detail below.

219 2.1 Buzzard-v1.0 simulation

220 The BUZZARD-HIGHRES-V1.0 put in cites to
 221 in prep Buzzard papers catalogue construction
 222 started with a dark matter only simulation. This
 223 N-body simulation contained 2048^3 particles in a
 224 400 Mpc h^{-1} box. [N] snapshots (with smoothing
 225 and interpolation between snapshots) were saved
 226 in order to construct a lightcone. Dark matter ha-
 227 los were identified using the ROCKSTAR software
 228 package (Behroozi et al. 2013). These dark matter
 229 halos were populated with galaxies with a stellar
 230 mass and absolute r -band magnitude in the SDSS
 231 system determined using a sub-halo abundance
 232 matching model constrained to match both pro-
 233 jected two-point galaxy clustering statistics and an
 234 observed conditional stellar mass function (Red-
 235 dick et al. 2013).

236 To assign an SED to each galaxy, the *Adding
 237 Density Dependent Spectral Energy Distributions*
 (ADDSEDS, deRose in prep.)⁴ procedure was
 238 used. This consisted of training an empirical
 239 relation between absolute r -band magnitude, local
 240 galaxy density, and SED using a sample of $\sim 5e^5$
 241 galaxies from the magnitude-limited Sloan Digital
 242 Sky Survey Data Release 6 Value Added Galaxy
 243 Catalog (Blanton et al. 2005)[Note: is this the
 244 proper reference to SDSS-NYU VAGC? File is
 245 called combined_dr6_cooper.fits, but I don't see
 246 which Cooper et al 2006 this is supposed to refer
 247 to]. Each SDSS spectrum is fit with a sum of five
 248 SED components using the K-CORRECT v? soft-
 249 ware package⁵ (Blanton & Roweis 2007), thus each
 250 galaxy SED is parameterized as five weights for the
 251 basis SEDs. The distance to the spatial projected
 252 fifth-nearest neighbour was used as a proxy for lo-
 253 cal density in the SDSS training sample. For each
 254 simulated galaxy, a “random” [details] galaxy with
 255 “similar” [details] absolute r -band magnitude and
 256 local galaxy density was chosen from the training
 257 set, and that training galaxy’s SED was assigned
 258 to the simulated galaxy. Given the SED, abso-
 259 lute r -band magnitude and redshift, we computed
 260 apparent magnitudes in the six LSST filter pass-
 261 bands, $ugrizy$. We assigned magnitude errors in
 262 the six bands using the simple model described in
 263 in Ivezić et al. (2008), assuming full 10-year depth ob-
 264 servations had been completed. The number of to-
 265 tal 30-second visits assumed when generating the
 266 photometric errors differs slightly from the fiducial
 267 numbers assumed for LSST: we assume 60 visits in
 268 u-band, 80 visits in g-band, 180 visits in r-band,
 269 180 visits in i-band, 160 visits in z-band, and 160
 270 visits in y-band.

272 2.1.1 Selection of training and test sets

273 The total catalogue covered 400 square degrees
 274 and contained 238 million galaxies to an apparent

⁴ <https://github.com/vipasu/addseds>

⁵ <http://kcorrect.org>

4 LSST-DESC Photometric Redshift Working Group

magnitude limit of $r = 29$ and spanning the redshift range $0 < z \leq 8.7$. This catalogue contained two orders of magnitude more galaxies than were needed for this study, so only ~ 8 square degrees were used. Systematic problems with galaxy colors above $z > 2$ were observed, so the catalogue was trimmed to include only galaxies in the redshift range $0 < z \leq 2.0$. A random subset of the remaining galaxies was chosen, and placed at random into either a “training” set (10 per cent of the sample), for which the galaxies true redshifts will be supplied, or a “test” set (the remaining 90 per cent of the sample), for which each code will need to predict a redshift PDF for each galaxy. The resulting catalogues contain 111 171 training galaxies and 1 000 883 test galaxies. We restrict our analysis to a sample with $i < 25.3$, which give a signal-to-noise ~ 30 for most galaxies, a cut often referred to as the expected “LSST Gold Sample”. This magnitude cut results in a training set with 44 404 galaxies and a test set containing 399 356 galaxies. All subsequent results will evaluate this “gold sample” test set.

2.1.2 Templates

As mentioned in Section 2.1, the SEDs in the Buzzard simulation are drawn from an empirical set of SEDs taken from the SDSS DR6 NYU-VAGC, a sample of roughly $\sim 5e^5$ galaxies with spectra in SDSS. To determine a finite set of templates to use with template fitting codes we take the five SED weight coefficients for each of the $\sim 500\,000$ galaxies in the SDSS sample and run a simple K-means clustering algorithm on this five dimensional space. The K-means cluster centres span the space of coefficients and properly reflect the underlying density in the coefficient space, thus providing a reasonable approximation for a spanning SED set. An ad-hoc number of $K = 100$ was chosen and the 100 K-means centre positions are taken as the weights for the K-CORRECT SED components to construct one hundred template SEDs. These 100 templates were provided, however not every template code uses this set of one hundred templates: because EAZY was designed and written to use the same five basis templates employed by K-CORRECT when constructing our mock galaxies, EAZY was run using linear combinations of these five templates rather than using the 100 discrete templates.

2.1.3 Limitations

For our initial investigation of photometric redshift codes, we begin with a data set that is somewhat idealized, and does not contain all of the complicating factors present in real data. In several cases, the simplification is done with a purpose, with potentially confounding effects excluded in order to better isolate the differences between current-generation photo- z codes, and their causes. We list several of the simulations limitations in this section. As the simulation is catalogue-based, no im-

age level effects, such as photometric measurement effects, object blending, contamination from sky background (Zodiacal light, scattered light, etc...), lensing magnification, or Galactic reddening are included. No stars are included in the catalogue, nor are the effects of AGN. As all SEDs are constructed from only five basis templates, properties of the galaxy population will be restricted to follow linear combinations of the characteristics of the five basis templates, so certain non-linear features, for example the full range of emission line fluxes relative to the continuum, will not be included in the model galaxy population. No additional dust reddening intrinsic to the host galaxy is included, the only approximation of dust extinction comes in the form of dust encoded in the five basis SEDs via the training set used to create the basis templates. Simple linear combinations of these basis templates will, once again, not explore the full range of realistic dust extinction observed in galaxy populations.

3 METHODS

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

The questions that must be answered for each code are: what unique features are included in the specific implementation that influence the output $p(z)$. What form of validation was performed with the training data, how were photometric uncertainties employed in the analysis, how were negative fluxes treated, what specific prior form was employed (for template based codes), or what specific machine learning architecture was used (for ML codes)?

3.1 Template-based Approaches

3.1.1 BPZ

BPZ⁶ (Bayesian Photometric Redshift, Benítez 2000) is a template-based photo- z code that compares the expected colors (C) calculated for a set of spectral energy distribution (SED) types/templates (T) to the observed colors to calculate the likelihood of observing colors at each redshift for each type, $p(C|z, T)$. The code employs an empirically determined Bayesian

in apparent magnitude (m_0) and SED-type. Assuming that the SED-types are spanning and exclusive, we can determine the redshift posterior $p(z|C, m_0)$ by marginalizing over all SED-types with a simple sum (Eq. 3 from Benítez 2000):

⁶ <http://www.stsci.edu/~dcoe/BPZ/>

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

Code	Type	Paper	Website
BPZ	template	Benítez (2000)	http://www.stsci.edu/~dcoe/BPZ/
EAZY	template	Brammer et al. (2008)	https://github.com/gbrammer/eazy-photoz
LePHARE	template	Arnouts et al. (1999)	http://www.cfht.hawaii.edu/~arnouts/lephare.html
ANNz2	ML	Sadeh et al. (2016)	https://github.com/IftachSadeh/ANNz2
DELIGHT	ML/template	Leistedt & Hogg (2017)	https://github.com/ixkael/Delight
FLEXZBOOST	ML	Izbicki & Lee (2017)	https://github.com/tpospisi/flexcode; https://github.com/rizbicki/FlexCoDE
GPz	ML	Almosallam et al. (2016a)	https://github.com/OxfordML/GPz
METAPHoR	ML	Cavuoti et al. (2017a)	http://dame.dsfc.unina.it
CMNN	ML	Graham et al. (2018)	-
SKYNET	ML	Graff et al. (2014)	http://ccforge.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	

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

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

In this paper we use BPZ v 1.99.3. The template set employed here is the set of 100 discrete SEDs described in Section 2.1.2 To keep the number of free parameters to a manageable level the SEDs in the training set are sorted by the rest-frame $u-g$ colour and split into three “broad” SED classes, equivalent to the E, Sp and Im/SB types in Benítez (2000). We assume the same functional form for the Bayesian priors as used by Benítez (2000), and utilize the training-set galaxies with known SED-type, redshift, and apparent magnitude to determine the type fractions and the best fit for the eleven free parameters of the prior. For galaxies with negative flux in a measured band, the placeholder value is replaced with an estimate one σ detection limit in that particular band, i. e. a value close to the estimated sky noise threshold. The type-marginalized $p(z)$ is generated by setting the parameter PROBS_LITE=TRUE in the BPZ parameter file.

3.1.2 EAZY

EAZY⁷ (Easy and Accurate Photometric Redshifts from Yale, Brammer et al. 2008) is a template-based photo-z code that includes several features that improve on the basic χ^2 fit used in many template codes. It can fit the observed photometry with SEDs created from a linear combination of a set of templates at each redshift, and the best-fit SED is found by simultaneously fitting one, two or all of the templates by minimizing χ^2 .

The minimized $\chi^2(z)$ is then combined with an apparent magnitude prior to obtain the posterior redshift probability distribution, although some argue that this is not the mathematically correct way of calculating the posteriors. EAZY can also account for the uncertainties in the templates by adding an empirically derived template error in quadrature as a function of redshift to the flux errors.

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

3.1.3 LePhare

LEPHARE⁸ (Photometric Analysis for Redshift Estimate, Arnouts et al. 1999; Ilbert et al. 2006) is a photo-z reconstruction code based on a χ^2 template-fitting procedure. The observed colors are matched with the colours predicted from a set of spectral energy distribution (SED) which can be either synthetic or based on a semi-empirical approach. LEPHARE has been used to produce the COSMOS2015 photo-z catalogue (Laigle et al. 2016).

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

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

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

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

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

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

3.2 Training-based Codes

3.2.1 ANNz2

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

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

In this study, ANNz2 v. 2.0.4 was used. The full PDF for each galaxy is also produced with a linear stepsize of $z = 0.01$ for $0 \leq z \leq 2$. A set of 5 ANNs with architecture 6 : 12 : 12 : 1 (6 *ugrizy* inputs, 2 hidden layers with 12 nodes each, and 1 output) with different random seeds are used during each training. Half of the training set is used as a validation set to prevent overtraining. All training objects are set to have detected magnitudes, however the non-detections ($\text{mag} = -99$)

in the testing set are replaced with the mean of that particular band.

3.2.2 Color-Matched Nearest-Neighbours

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

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

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

¹⁰ <http://tmva.sourceforge.net/>

robust estimate of $p(z)$ is desired for this work we make several changes to our implementation of the CMNN photo- z estimator. We continue to use a percent point function of PPF = 0.95 to generate the subset of color-matched training galaxies, but weight them by the inverse of their Mahalanobis distance. This weighting maintains some of the accuracy that was previously achieved by simply using the nearest neighbour in color-space. We then use the weights to create the $p(z)$ instead of having the redshift of each color-matched training-set galaxy count equally. To obtain a robust estimate of the $p(z)$ for galaxies with a small number of color-matched training set galaxies, when this number is less than 20 the nearest 20 neighbours in color-space are used instead, and we convolve the $p(z)$ with a Gaussian with a standard deviation of:

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

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

3.2.3 Delight

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

This conceptually novel approach is done by using Gaussian processes operating in flux-redshift

space. The posterior distribution on the redshift of a target galaxy is obtained via a pairwise comparison with training galaxies,

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

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

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

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

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

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

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

3.2.4 FlexZBoost

FLEXZBOOST¹² (Izbicki & Lee 2017) is a particular realization of FlexCode, which is a general-purpose methodology for converting any conditional mean point estimator of z to a conditional density estimator $f(z|\mathbf{x})$, where \mathbf{x} here represents

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

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

8 LSST-DESC Photometric Redshift Working Group

our photometric covariates and errors.¹³ The key idea is to expand the unknown function $f(z|\mathbf{x})$ in an orthonormal basis $\{\phi_i(z)\}_i$:

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

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

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

These coefficients can easily be estimated from data by regression.

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

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

714 3.2.5 GPz

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

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

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>

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

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

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

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

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

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

788 data set that includes galaxies fainter than those
 789 in the test set, which was truncated at $i < 25.3$.
 790 This mismatch in training vs. test results in a
 791 slightly different distribution in galaxy photome-
 792 try and redshift distribution between the training
 793 and test set, which most likely degraded results
 794 compared to a run where the training set more
 795 closely matched the test set.

3.2.6 METAPhoR

796 METAPhoR (Machine-learning Estimation Tool
 797 for Accurate Photometric Redshifts, [Cavuoti et al.
 798 2017a](#)) is a pipeline designed to provide photo-z's
 800 point estimates and a reliable PDF for machine
 802 learning (ML) based techniques. It includes pre-
 804 and post-processing phases, hosting a photo-z pre-
 806 diction engine based on the Multi Layer Percep-
 808 tron with Quasi Newton Algorithm (MLPQNA),
 810 already validated on photo-z's in several cases ([de
 Jong et al. 2017; Cavuoti et al. 2017b, 2015; Bres-
 812 cia et al. 2014, 2013; Biviano et al. 2013](#)). Due to
 814 its plug-in based modular nature, METAPhoR
 816 can be easily replaced by any other photo-z pre-
 818 diction kernel, regardless its implementation, by
 820 taking the I/O interface compliance as unique con-
 822 straint.

824 At a higher level, the pipeline mainly consists
 826 of three modules: (i) *data pre-processing*, includ-
 828 ing a catalogue cross-matching sub-module (based
 830 on the tool C3, [Riccio et al. 2017](#)), a sub-module
 832 for photometric evaluation and error estimation of
 834 the multi-band catalogue used as Knowledge Base
 836 (KB), and a sub-module dedicated to the per-
 838 turbation of the photometric KB, propaedeutic to the
 840 PDF estimation; (ii) *photo-z prediction*, which is
 842 the training/validation/test phase, producing the
 844 photo-z's point estimates, based on a pre-selected
 846 ML method; (iii) *PDF estimation*, specifically de-
 848 signed to calculate the PDF of the photo-z estima-
 850 tion errors. The last module includes also a post-
 852 processing tool, providing some statistics on the
 854 produced point estimates and PDFs.

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

880 As introduced, the photo-z point estimate pre-
 882 diction engine of METAPhoR is based on the
 884 MLPQNA model, whose photo-z regression train-
 886 ing error, used by the quasi Newton learning rule,
 888 is based on the least square error and Tikhonov
 890 L_2 -norm regularization ([Hofmann & Mathé 2018](#)).

892 As main prerogative, METAPhoR is able to
 894 provide a PDF for ML methods by taking into ac-

896 count the photometric errors provided with data,
 898 by running N trainings on the same training set,
 900 or M trainings on M different random extractions
 902 from the KB. The different test sets, used to pro-
 904 duce the PDF, are thus obtained by introducing a
 906 proper perturbation, parametrized from the pho-
 908 tomatic error distribution in each band, on the
 910 photometric data populating the original test set
[\(Brescia et al. 2018\)](#).

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

3.2.7 SkyNet

926 SKYNET¹⁵ ([Graff et al. 2014](#)) is a publicly avail-
 928 able neural network software, based on a 2nd order
 930 conjugate gradient optimization scheme (see [Graff
 932 et al. 2014](#), for further details). It has been used ef-
 934 ficiently for redshift PDF estimates ([Sánchez et al.
 936 2014; Bonnett 2015; Bonnett et al. 2016](#)).

938 The neural network is configured as a stan-
 940 dard multilayer perceptron with three hidden lay-
 942 ers and one input layer with 12 nodes (the 6 mag-
 944 nitudes and their errors). The classifier is laid out
 946 such that the hidden layers have 20:40:40 nodes
 948 each, all rectified linear units, and the output layer
 950 has 200 nodes (corresponding to 200 bins for the
 952 PDF) activated with a “softmax” function so that
 954 they automatically sum to 1.

956 To avoid over-fitting, a 30 per cent frac-
 958 tion of the training set is used as validation,
 960 and the training is stopped as soon as the er-
 962 rror rate begins to increase in the validation set.
 964 The weights are randomly initialized based on nor-
 966 mal sampling. The error function is a standard
 968 chi-square function for the regressor, and a cross-
 970 entropy function for the classifier. Finally, the data
 972 are all whitened before processing, with magni-
 974 tudes pegged to (45,45,40,35,42,42) and their er-
 976 rrors pegged to (20,20,10,5,15,15) for *ugrizy* filters,
 978 respectively.

3.2.8 TPZ

980 TPZ¹⁶ ([Trees for Photo-z, Carrasco Kind & Brunner 2013; Carrasco Kind & Brunner 2014](#)) is a
 982 parallel machine learning algorithm that gener-
 984 ates photometric redshift PDFs using prediction
 986 trees and random forest techniques. The code re-
 988 cursively splits the input data (i. e. the training
 990 sample), into two branches, one after another, until
 992 a terminal leaf is created that meets a termina-
 994 tion criterion (e. g. a minimum leaf size or a variance
 996 threshold). Bootstrap samples from the training

¹⁵ [http://ccforge.cse.rl.ac.uk/gf/project/
 skynet/](http://ccforge.cse.rl.ac.uk/gf/project/skynet/)

¹⁶ <https://github.com/mgckind/MLZ>

904 data and associated errors are used to build a set
 905 of prediction trees. In order to minimize correlation
 906 between the trees, the data is divided in such
 907 a way that the highest information gain among the
 908 random subsample of features is obtained at every
 909 point. The regions in each terminal leaf node cor-
 910 responds to a specific subsample of the entire data
 911 that possesses similar properties.

912 The training data is examined before running
 913 TPZ. Since TPZ does not handle non-detections
 914 (magnitudes flagged as 99.0), we replace these val-
 915 ues with an approximation of the 1σ detection
 916 threshold, i. e. a signal to noise ratio of 1 in
 917 terms of magnitude uncertainty using the equation
 918 $dm = 2.5 \log(1 + N/S)$ where $dm \sim 0.7526$ mag
 919 for $N/S = 1$. That is, for each band, we replace
 920 the non-detection with the magnitude correspond-
 921 ing to the error of 0.7526 from the error model
 922 forecasted for 10-year LSST data. The Out-of-
 923 Bag (Breiman et al. 1984; Carrasco Kind & Brunner
 924 2013) cross-validation technique is used within
 925 TPZ to evaluate its predictive validity and deter-
 926 mine the relative importance of the different input
 927 attributes. We employed this information to cali-
 928 brate our algorithm.

929 In the present work, the LSST magnitudes
 930 u, g, r, i and colors $u-g, g-r, r-i, i-z, z-y$
 931 and their associated errors are used in the pro-
 932 cess of growing 100 trees with a minimum leaf size
 933 of 5 (the z and y magnitudes did not show sig-
 934 nificant correlation with the redshift in our cross-
 935 validation, so we did not use them when construct-
 936 ing our trees). We partitioned our redshift space
 937 into 100 bins from $z = 0.005$ to $z = 2.0$ and
 938 smoothed each individual PDF with a smoothing
 939 scale of twice the bin size.

3.3 Simple Ensemble Estimator

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

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

944 That is, we simply set the redshift PDF of every
 945 galaxy equal to the normalized $N(z)$ of the train-
 946 ing sample. As the training sample is drawn from
 947 the same underlying distribution as the test sam-
 948 ple, modulo small deviations due to sample size,
 949 the quantiles of the training and test distributions
 950 should be identical. This is a wildly unrealistic es-
 951 timator, as it assigns all galaxies, no matter their
 952 apparent magnitude, colour, or true redshift, the
 953 same redshift PDF, and is thus uninformative at
 954 the level of individual object redshifts, but is de-
 955 signed to perform very well for the ensemble of all
 956 objects. We will discuss this method and cautions
 957 relative to metrics in Section 5.3.

4 METRICS FOR QUANTIFYING PDF COMPARISONS

960 The overloaded “ $p(z)$ ” is a widespread abuse of no-
 961 tation; we would like the outputs of photo- z PDF
 962 codes to be interpretable as probabilities. Obvi-
 963 ously photo- z PDFs must not take negative values
 964 and must integrate to unity over the range of pos-
 965 sible redshifts. Additionally, an estimator derived
 966 by method H for the photo- z PDF of galaxy i must
 967 be understood as a posterior probability distribu-
 968 tion

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

969 conditioned not only on the photometric data d_i
 970 for that galaxy but also on parameters encompass-
 971 ing a number of things that will differ depending
 972 on the method H used to produce it, namely the
 973 assumptions I_H necessary for the method to be
 974 valid and any inputs I_D it takes as prior informa-
 975 tion, such as a template library or training set. Be-
 976 cause of this, direct comparison of photo- z PDFs
 977 produced by different methods is in some sense im-
 978 possible; even if they share the same prior informa-
 979 tion I_D , by definition they cannot be conditioned
 980 on the same assumptions I_H , otherwise they would
 981 not be distinct methods at all.

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

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

990 meaning that we assume comparisons of $\hat{p}_{i,H}(z)$
 991 isolate the effect of the method used to obtain the
 992 estimator, which should make examination of dif-
 993 ferences caused by specifics of the method imple-
 994 ments easier to isolate.

995 As mentioned previously, there are cosmology
 996 probes that require knowledge of individual galaxy
 997 $p(z)$ and others that require only knowledge of the
 998 ensemble redshift distribution, $N(z)$. Due to the
 999 paucity of principled techniques for using and val-
 1000 idating photo- z PDFs, there are few alternatives to
 1001 the common practice of reducing photo- z PDFs to
 1002 point estimates. Though this practice should not
 1003 be encouraged, we also calculate traditional met-
 1004 rics based on the most common point estimators
 1005 derived from photo- z PDFs. Those seeking to es-
 1006 tablish a connection to traditional ways of think-
 1007 ing about redshift estimation may consult the Ap-
 1008 pendix for these results.

1009 There are a number of metrics that can be
 1010 used to test the accuracy of a photo- z interim pos-
 1011 terior as an estimator of a true photo- z posterior
 1012 if it is known. Even for simulated data, the true
 1013 photo- z PDF is in general not accessible unless the

redshifts are in fact drawn from the true photo- z PDFs, 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 by far is the calculation 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 in prep.), stacking according to

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

is the most widely accepted method for estimating the redshift distribution from photo- z PDFs. If we assume that the response of estimators of $N(z)$ is uniform across all approaches H , then we may interpret metrics on the accuracy of $\hat{N}(z)$ obtained in this way. We must note, however, that this is a poor assumption in general. Under the setup of this paper, the true redshift distribution $N'(z) = p(z|I_D)$ (i.e. because our training data is representative, the interim prior is the truth). In this ideal case, 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 . This is the exact estimator, TRAINZ, that we have described in Section 3.3, and which will serve as a point of reference for the other codes.

The exact implementation of the stacked estimator $\hat{N}^H(z)$ will 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 such choices, we asked those running all eleven codes to output $p(z)$ parameterized with a generous ≈ 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 $p(z)$ for the final LSST Project tables.¹⁷ All $p(z)$ catalogues are processed using the QP software package (Malz et al. 2018)¹⁸ for manipulating and calculating metrics of 1-dimensional PDFs. We will discuss the choice of $p(z)$ parameterization further in Section 5.

4.1 Metrics of an ensemble of photo- z interim posteriors

4.1.1 Probability integral transform (PIT)

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

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

The distribution of PIT values quantifies the behavior of the *ensemble* of photo- z PDFs, enabling us to evaluate whether the $p(z)$ is, on average, accurate: The PIT value is the Cumulative Distribution Function (CDF) of the $p(z)$ evaluated at the 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 overrepresentation 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 is outside of the support of $p(z)$. Tanaka et al. (2017) 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

The quantile-quantile (QQ) plot is a graphical method for comparing two distributions, where the quantiles of one distribution are plotted against the quantiles of the other distribution (A quantile being defined by partitioning a distribution into consecutive intervals containing equal amounts of probability, or equal numbers of objects in each interval). In this paper we show the quantiles of the PIT values compared to the quantiles of the Uniform distribution that we expect the PIT values to match if $p(z)$ is an accurate probability distribution for all objects. 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. In this paper, QQ plots are used for two purposes: (1) for comparing $N(z)$ from photo- z PDFs (estimated using Eq. 14) with the true $N(z)$, i.e. comparing the estimated distribution of redshifts with the true redshift distribution, 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

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

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

1122 previous section) is compared to a uniform distribution
 1123 between 0 and 1. The QQ plot contains
 1124 very similar information to that shown in the PIT
 histogram plot, we include both forms, as visually
 1125 they each convey the information in a somewhat
 1126 distinct manner.

1128 4.1.3 Conditional density estimation loss

With the conditional density estimation loss (CDE)
 1129 loss) we can compare how well different methods
 1130 estimate individual PDFs for photometric covariates
 1131 \mathbf{x} rather than looking only at the ensemble
 1132 distribution. As in Section 3.2.4, we use the notation
 1133 $f(z|\mathbf{x})$ instead of $p(z)$ to explicitly show the
 1134 dependence on \mathbf{x} .

The CDE loss is defined as:

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

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

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

where 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 a 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 $p(z)$ as described in Eq. 14.

4.2.1 Root-mean-square error (RMSE)

We employ the familiar root-mean-square error:

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

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, and it can be used to quantify the otherwise qualitative metrics.

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:

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

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

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

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

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

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

where n is the sample size.

4.2.3 Moments

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

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

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

5 RESULTS

5.1 Ensembles of photo-z interim posteriors

Fig. 1 Shows the $p(z)$ produced by each of our eleven photo- z codes for four example galaxies

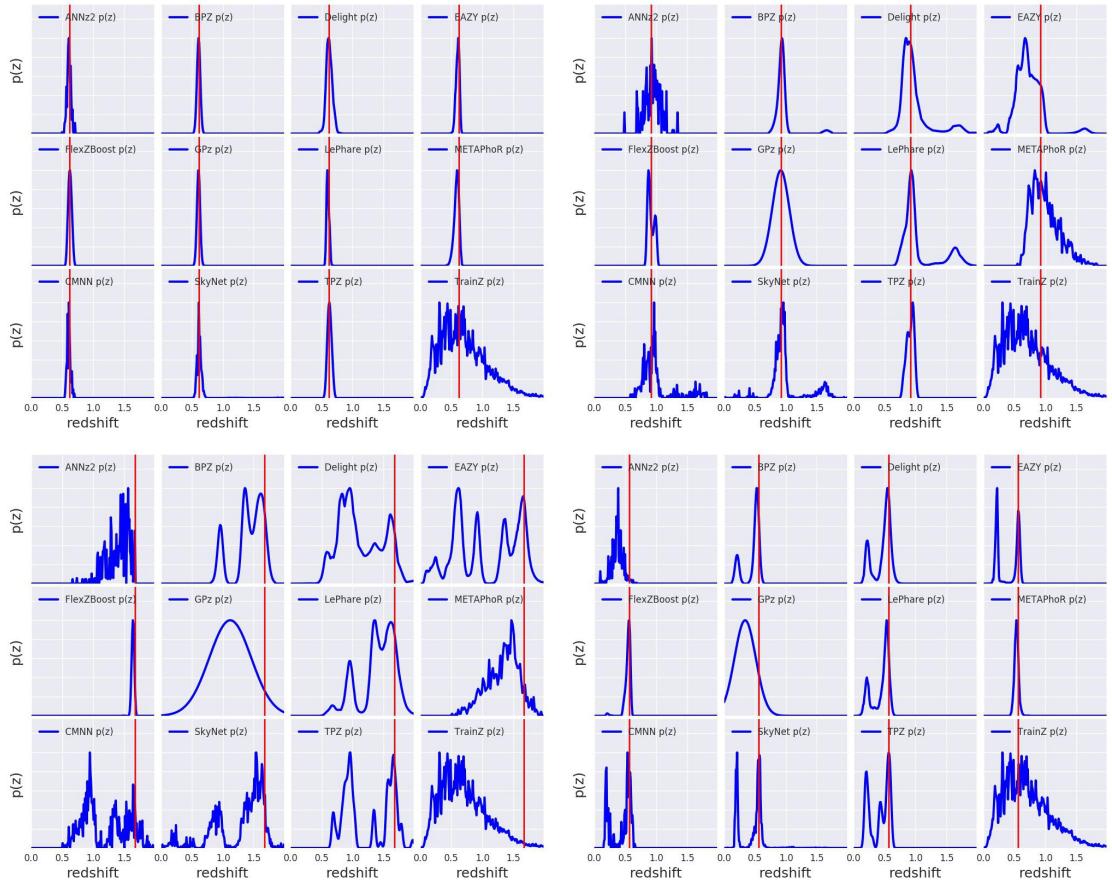


Figure 1. Four illustrative examples of individual $p(z)$ distributions produced by the codes. The red vertical line represents the true redshift. Examples are chosen with common features seen in PDFs: tight unimodal $p(z)$ (upper left), broad unimodal $p(z)$ (upper right), bimodal $p(z)$ (lower right),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

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

As stated in Section 4, $p(z)$ is parameterized as ≈ 200 piecewise constant bins covering $0 < z < 2$ for all eleven codes, giving a grid size of roughly $\delta z = 0.01$ for each code. A piecewise constant grid was a natural choice for some photo- z codes, for instance most template-based codes compute likelihoods on a fixed grid. In contrast, FlexZBoost, for example, can return estimates on any grid without compression errors as it's a ba-

sis expansion method where only the expansion coefficients need to be stored. Codes with a native output format other than the shared piecewise constant binning scheme (or one that can be losslessly converted to it) may suffer from loss of information when converting to it, which could artificially favor some codes over others.

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

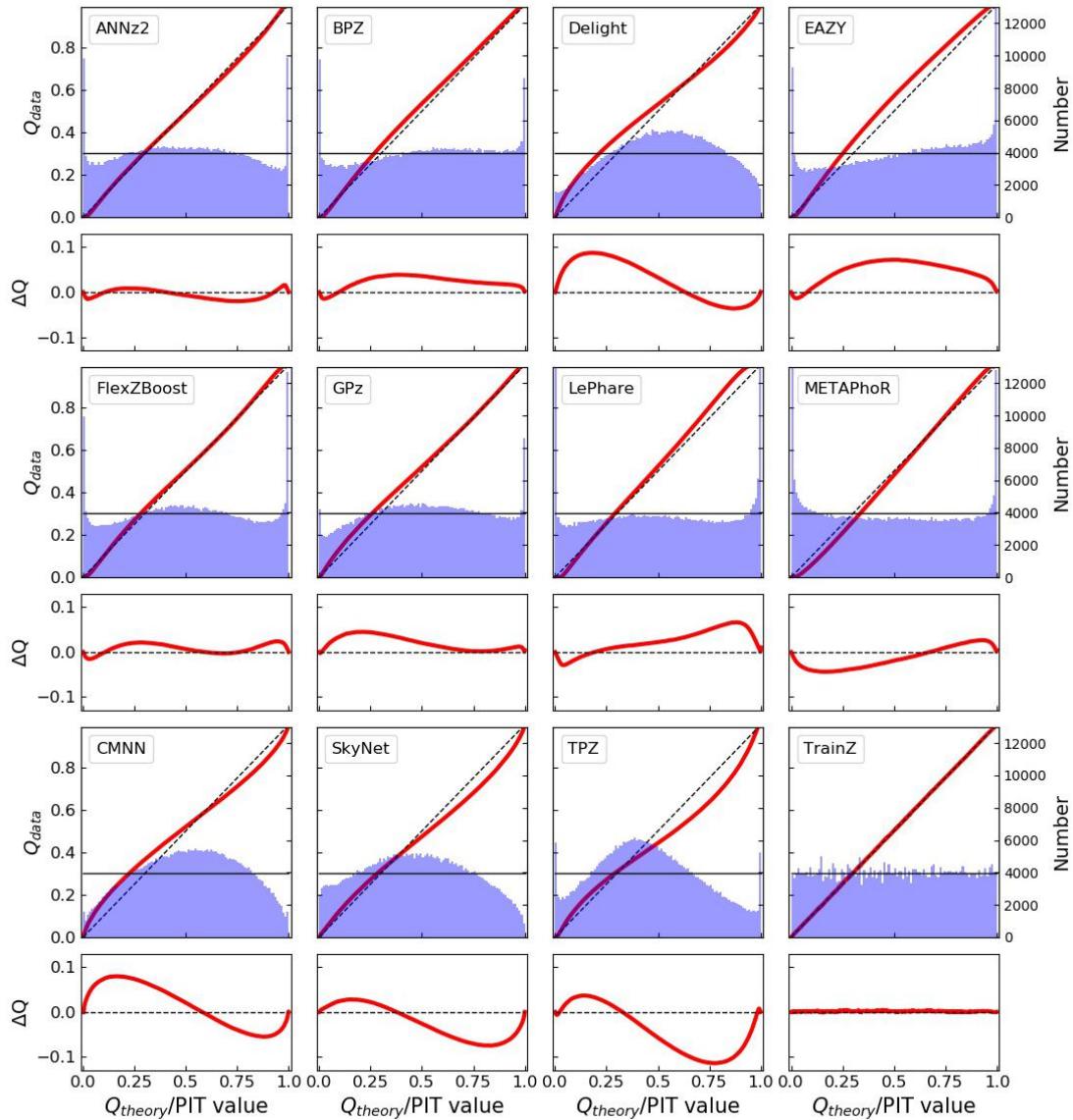


Figure 2. Summary plots for all eleven photo- z codes illustrating performance for the interim posterior statistics. The top panel of each pair shows both the Quantile-Quantile (QQ) plot (red) and the histogram of PIT values (blue). The desired behavior is a QQ plot that matches the diagonal dashed line,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

1256 tions and limits on storage resources may be con- 1266 sidered in future work. We will discuss this further
1258 in Section 6.

Fig. 2 shows both the quantile-quantile plots 1260 (red) and the histogram of PIT values (blue) sum- 1262 marizing the results from each photo- z code. The 1264 red line shows the measured quantiles, while the 1270 black diagonal represents the ideal QQ values if 1274 the distribution were perfectly reproduced. A sec-
ond panel below the main panel for each code

shows the difference between Q_{data} and Q_{theory} , i. e. the departure from the diagonal, for clarity. Biases and trends in whether the average width of the $p(z)$ values being over/under-predicted are evident. An overall bias where the predicted redshift is systematically low manifests as the measured QQ value falling above the diagonal, as is the case for BPZ and EAZY, while a systematic overprediction shows up as the measured QQ value falling below the diagonal, as seen in TPZ. In

terms of PIT histograms, a systematic underprediction of redshift corresponds to fewer PIT values at $PIT < 0.5$ and more at $PIT > 0.5$, while a systematic overprediction will show the opposite.

Examination of the PIT histograms and QQ plots shows that there are fairly generic issues with the width of $p(z)$ uncertainties: DELIGHT, NN, SKYNET and TPZ all show a PIT histogram with an dearth of low values and an excess of high values, signs that, on average, their $p(z)$ are more broad than the true distribution of redshifts. METAPHoR shows the opposite trend, indicating the $p(z)$ are more narrow than the distributions given by the true redshifts. In all of these code cases there is a free parameter or bandwidth that can be used to tune uncertainties. The sensitivity of multiple codes to this bandwidth choice emphasizes the fact that great care must be taken in setting user-defined parameters in photo- z codes, even in the presence of representative training/validation data. for FLEXZBoost the “sharpening” parameter (described in Section 3.2.4) plays a key role in improving the results, resulting in a QQ plot that is very nearly diagonal. A similar sharpening procedure could be beneficial for several codes. Interestingly, the three 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_0 as the fraction of objects with $PIT < 0.0001$ or $PIT > 0.9999$. Table 2 lists these fractions for each of the codes. For a proper Uniform distribution we expect a value of 0.0002. Several codes show a marked excess, with ANNz2, FLEXZBOOST, LEPHARE, AND METAPHoR with $f_0 > 0.02$, 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

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	0.0068
LEPHARE	0.0486
METAPHoR	0.0229
CMNN	0.0034
SKYNET	0.0001
TPZ	0.0130
TRAINZ	0.0002

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 extrema, and thus is calculated by excluding the edges of the distribution. We calculate the AD statistic over the range of PIT values $v = [0.01, 0.99]$. ANNz2 and FLEXZBoost score very well for the PIT metrics. METAPHoR and LEPHARE score very well in the PIT AD statistic, but both have 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.

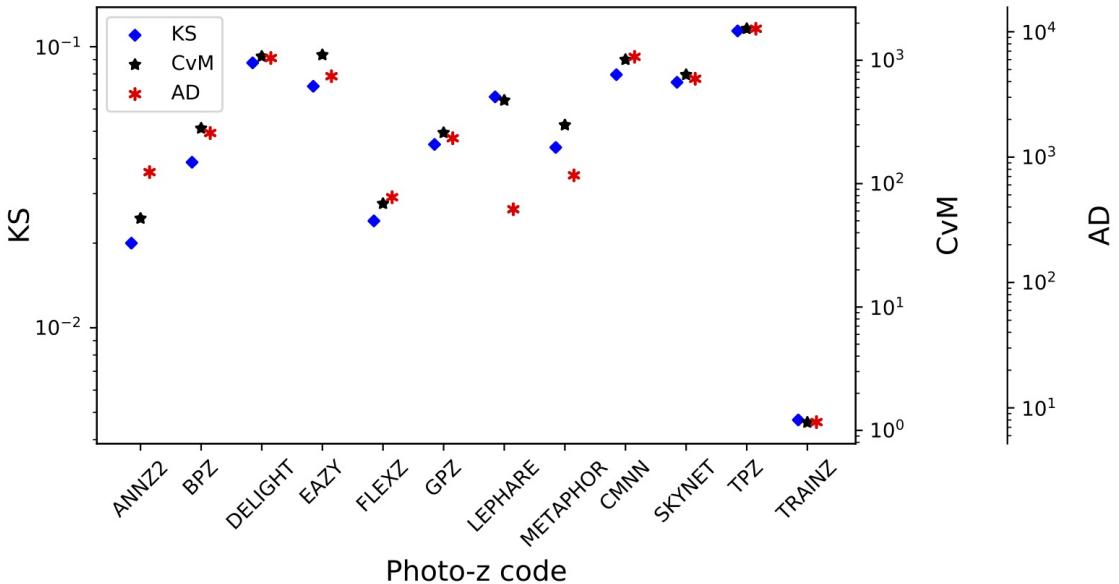


Figure 3. A visual representation of the Kolmogorov-Smirnov (KS, blue diamond), Cramer-von Mises (CvM, black star), and Anderson-Darling (AD, red asterisk) statistics for the PIT distributions. The statistics are often highly correlated, though the AD statistic truncates the extrema of the distribution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

5.2 Metrics of the stacked estimator of the redshift distribution

Fig. 4 shows the stacked $\hat{N}(z)$ distribution compared to the true redshift distribution $N'(z)$ for all tested codes. The red line indicates the summed $p(z)$ for each code, while the blue line shows the true redshift distribution smoothed via kernel density estimation (KDE), with a bandwidth chosen via Scott's rule (Scott 1992). While Scott's rule is used to display $N'(z)$ in the figure, all quantitative statistics are computed via the empirical CDF, and are thus unaffected by bandwidth/smoothing choice. Several of the codes show an excess at $z \sim 1.4$, particularly the template-based codes BPZ, EAZY, and LEPHARE. This is likely due to the 4000 angstrom break passing through the gap between the z and y filters. This feature is one of the most prominent in individual galaxy $p(z)$, and is readily seen in the point-estimate plots shown in Fig. A1 and described in the Appendix. Several of the machine learning based codes appear to be over-trained, adding excess galaxy probability to the redshift peaks and missing probability in the troughs. Given that our training data is drawn from the same galaxy population as the test set, and our data has prominent peaks in $N'(z)$, perhaps it is not unexpected that such overtraining occurs. A more extensive training/validation set might allow for a better choice of smoothing parameters in individual codes that would avoid such overtraining.

As with the $p(z)$ values in Figure 2, different levels of substructure are obvious for the different codes. While Scott's rule provides a relatively good general smoothing scale to represent the true $N'(z)$, there are smaller scale fluctuations:

while FLEXZBoost and CMNN appear somewhat discrepant in Fig. 4, they are actually the two most accurate in terms of their quantitative measurements. Interestingly, while ANNz2 shows an abundance of small scale structure in individual $p(z)$ measurements (see Fig. 1), the summed $\hat{N}(z)$ is rather smooth, where the small scale features average out. This is not the case for the two other codes that show an abundance of substructure in their individual $p(z)$: both CMNN and SKYNET show small scale features both in $p(z)$ and $\hat{N}(z)$. For CMNN the $p(z)$ are a simply a weighted histogram of all spectroscopic training galaxies in nearby colour space with no smoothing applied, so the substructure is not unexpected. The PIT histogram and shape of the QQ plot in Figure 2 show that CMNN is producing $p(z)$ that are overly broad, additional smoothing of the $p(z)$ would exacerbate this problem. While the $\hat{N}(z)$ plot shows more small scale features than other codes, these features are actually representative of real structure in the true $N'(z)$, as evidenced by the very good metric scores for CMNN. SKYNET $p(z)$ were also not smoothed: while previous implementations of the code such as Sánchez et al. (2014) and Bonnett (2015) (see Appendix C.3) implement a “sliding bin” smoothing, no such procedure was used in this study. In addition to excess substructure, SKYNET shows an obvious redshift bias, evident both visually in Figure 4 and in the first moment of $N(z)$ listed in Table 5, where it is clearly an outlier. SKYNET employed a method where a random sample of training galaxies was chosen, but there was no test that the subset was completely representative of the overall redshift distribution. Also unlike Bonnett (2015), no effort

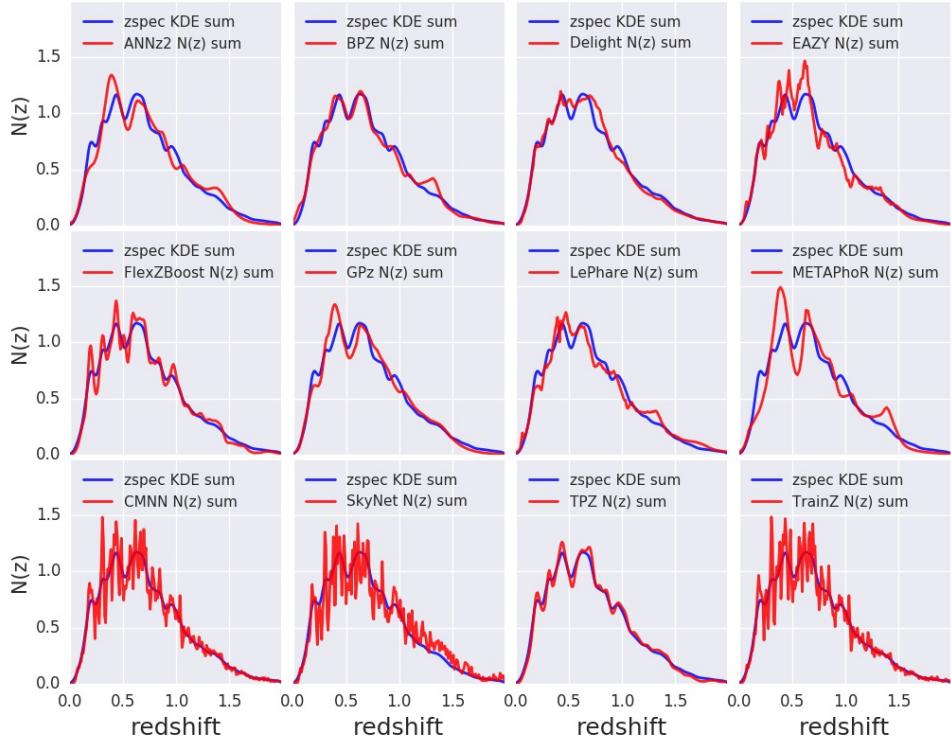


Figure 4. The stacked $p(z)$ produced by each photo- z code ($\hat{N}(z)$, red) compared to the spectroscopic redshift distribution ($N'(z)$, blue). Varying levels of small-scale structure are seen in the codes. $N'(z)$ is smoothed using a single bandwidth chosen via Scott’s rule for all codes.

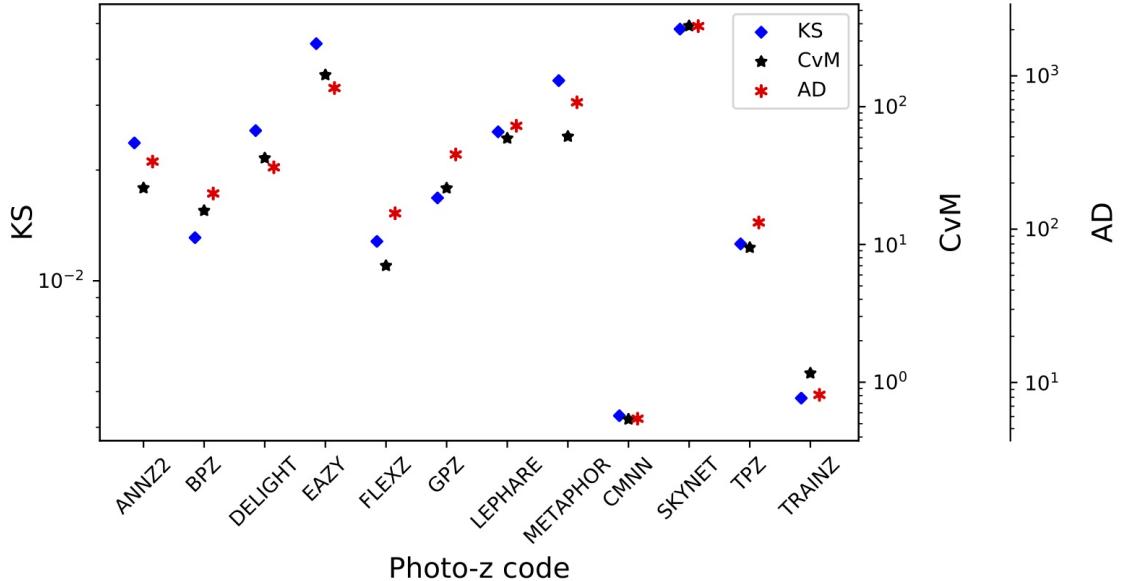


Figure 5. A visual representation of the Kolmogorov-Smirnov (KS, blue diamond), Cramer-von Mises (CvM, black star), and Anderson-Darling (AD, red asterisk) statistics for the $\hat{N}(z)$ distributions. The statistics are correlated, the codes with the lowest KS statistics tend to have the lowest CvM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

was made to add extra weight to more rare low and high redshift galaxies. Either of these decisions could be the cause of the bias seen in our results. Future runs of SKYNET will explore these implementation choices and their effects.

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

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

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

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

5.3 Interpretation of metrics

Samples from accurate photo- z posteriors should reproduce the space of $p(z, data)$. 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 corresponds to *accurate* photo- z posteriors. 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

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	-9.56
LEPHARE	-1.66
METAPHoR	-6.28
CMNN	-10.43
SKYNET	-7.89
TPZ	-9.55
TRAINZ	-0.83

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. We note that many of the ensemble-based metrics are prominent in the photo-z literature despite their inability to identify problems such as those exemplified by TRAINZ.*

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

6 SUMMARY AND DISCUSSION

In this paper we presented results evaluating the photometric redshift PDF computation for eleven photo-z codes. As discussed in Section 4 the $p(z)$ should accurately reflect the relative likelihood as a function of redshift for each galaxy. All codes were provided a set of representative training data and tested on an idealized set of model galaxies with high signal-to-noise and photometry with no confounding effects due to blending, instrumental ef-

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

Root-Mean-Square-Error (RMSE) statistics			
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	0.048	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	0.692	0.605	0.619
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

fects, the night sky, etc... included. The goal was not to determine a “best” photo-z code: in many ways, this was a baseline test of a “best case scenario” to predict the expected photo-z performance if a stage IV dark energy survey was to obtain complete training samples and perfectly calibrated their multi-band photometry. Given these idealized conditions, any deficiencies observed in a photo-z code’s performance should be a cause for concern, and may be evidence in a problem with either/both of the specific code implementation or the underlying algorithm. In order to meet the stringent LSST requirements on photo-z performance, identifying and correcting such problems is an important first step before tackling more realistic data in future challenges. Most of the codes tested performed well, however, several did not meet the stringent goals that have been laid out for LSST photometric redshift performance. This is a cause for concern, given the idealized conditions, and the individual code responses will be studied in detail moving for-

ward. 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. [are there any references for this? I remember Gary Bernstein talking about this at a photo-z workshop in Japan, but I don't know that it was published. I believe Michael Troxel has discussed this as well.] We also plan to propagate the uncertainties measured in a set of fiducial tomographic redshift bins in order to estimate impact on cosmological parameter estimation.

In future papers we will add more and more complexity to our simulated data in order to test photo-z algorithms in increasingly realistic conditions. The most pressing concern is the impact of incomplete spectroscopic training samples. As discussed extensively in Newman et al. (2015) a representative set of spectroscopically confirmed galaxies spanning the full range of both redshift and apparent magnitude is necessary as a training set to characterize the mapping from broad-band fluxes to photometric redshifts. However, due to a combination of factors due to both the galaxy SEDs and limitations of spectrographic instruments, redshift samples are known to be systematically incomplete, where certain galaxy types and redshift intervals fail to yield a redshift even at the longest integration times on current and near-future instruments. The more representative the training data, the better the performance of photo-z algorithms will be. Current and upcoming surveys are putting in significant effort into obtaining these training samples (e.g. Masters et al. 2017), however we still expect significant incompleteness for LSST-like samples, particularly at faint magnitudes. One major focus of an upcoming LSST Dark Energy Science Collaboration Photo-z Working Group data challenge is to produce a realistically incomplete training set of spectroscopic galaxies, modeling the performance of spectrographs, emission-line properties, and expected signal-to-noise to determine which galaxies will fail to yield a secure redshift. In addition to outright redshift failures we will model the inclusion of a small number of falsely identified secure redshifts where misidentified emission lines or noise spikes cause an incorrect redshift solution to be marked as a high quality identification. Even sub-per cent level contamination by false redshifts can impact photo-z solutions at levels comparable to the stringent requirements of some LSST science cases. We expect different systematics to occur in different photo-z codes in response to training on incomplete data, particularly some of the machine learning methods. The response of the codes will inform future directions of code development.

This initial paper explored a data set that was constructed at the catalog level, with no inclusion of the complications that come from measuring photometry from images. Future data challenges will move to catalogs constructed from mock images, including effects that will have great impact

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

The underlying simulation used in this work was based on a light-cone constructed to a maximum redshift of $z = 2$. LSST imaging after 10 years of observations will include a significant number of $z > 2$ galaxies in expected cosmology samples, and their inclusion does have potential significant implications for photo-z measures: the high redshift galaxies lie at fainter apparent magnitudes and can have anomalous colours due to evolution of stellar populations and the shift to rest-frame magnitudes probing UV features of the underlying SED. More importantly, one of the most common “catastrophic outlier” degeneracies observed in deep photometric samples occurs when the Lyman break is mistaken for the Balmer break, leading to multiple redshift solutions at $z \sim 0.2 - 0.3$ and $z \sim 2 - 3$ (Massarotti et al. 2001). This degeneracy, along with other potential degeneracies, are currently not covered by the limited redshift range of this initial paper, which could mean that we are not probing the full range of potential extreme outlier populations and how our photo-z estimators respond to them. Extending simulations to include the high-redshift galaxy population will be a priority in future data challenges.

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

In addition to AGN dominated galaxies, those with lower levels of nuclear activity present a more

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

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

In addition to these improvements, the DESC Photo-z group plans to look at all potential methods to combine the results from multiple photo-z codes to improve $p(z)$ accuracy, similar to the work presented in Dahlen et al. (2013); Carrasco Kind & Brunner (2014); Duncan et al. (2018). Taking advantage of multiple algorithms that use observables in slightly different ways has shown promise, however we must be very conscious of whether a potential combination properly treats the covariance between the methods, given that they are estimating quantities based on the same underlying observables. Several science cases wish to estimate physical quantities along with redshift, for example galaxy stellar mass and star formation rate. Proper joint estimation of redshift and physical quantities requires an in depth understanding of galaxy evolution, and progress on accurate bivariate redshift probability distributions will go hand in hand with progress on understanding galaxies themselves. Parameterization and storage of a complex 2-dimensional probability surface for potentially billions of galaxies (or even subsets of hundreds of thousand of particular interest) pose a potential challenge. These issues will be examined in another future paper.

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 series of future 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.

ACKNOWLEDGEMENTS

[summary of effort breakdown], will need to be in authors.csv for auto-generation by start_paper [personal funding sources] /NERSC computation acknowledgement

The LSST Dark Energy Science Collaboration acknowledges generous ongoing support from the agencies and institutes supporting the collaboration. These include 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 Energy, the National Science Foundation, and the LSST Corporation in the United States. This research used resources of 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. Part of this work was undertaken on STFC DiRAC HPC Facilities, funded by UK BIS National E-infrastructure capital grants, and on the UK particle physics grid, supported by the GridPP Collaboration.

APPENDIX A: POINT ESTIMATE PHOTOMETRIC REDSHIFTS

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

There are several common point estimators of photo-z posteriors employed by different codes, e. g. the mode, mean, median of the $p(z)$ distribution. In addition, many of the machine learning based estimators can be set up to return a single redshift solution. For example, SkyNet can be configured to run as a regressor that returns a single float rather than a classifier that returns a 200-bin $p(z)$ estimate. The single value returned by a machine learning based code may not correspond to a particular measure such as the mode or mean, and so to avoid interpretation of results that might be introduced by variations in choice of specific point-estimate implementation per code, we discard the code-specific point estimates. We

instead calculate point estimates more uniformly across the codes directly from the $p(z)$ using two measures, z_{PEAK} and z_{WEIGHT} . z_{PEAK} is simply the maximum value attained for each galaxy $p(z)$, the mode of the probability distribution. z_{WEIGHT} is defined similarly to how it is defined in Dahlen et al. (2013), as the weighted mean of the redshift over the main peak of $p(z)$ containing the z_{PEAK} value. The main peak is defined by subtracting $0.05 \times z_{\text{PEAK}}$ from $p(z)$ and identifying the roots to isolate the peak containing z_{PEAK} , z_{WEIGHT} is defined as the weighted mean redshift within this peak. We restrict to a single peak in order to avoid confusion from bimodal and multimodal $p(z)$ such as those shown in bottom panels of Figure 1. For example, for a bimodal probability distribution a weighted mean calculated over both peaks would fall between the peaks, at a redshift where the probability is minimal. Restricting the weighting to a single peak ensures that the point estimate will fall in the region of maximum redshift probability.

A1 Point Estimate Metrics

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

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

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

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

These definitions are similar, but not exactly the same, as the σ_{IQR} and median bias calculated here,

but are similar enough for qualitative comparisons to the LSST goals.

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

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

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

References

- Abbott T., et al., 2005, preprint ([arXiv:astro-ph/0510346](#))
- Abell P. A. et al., 2009, preprint ([arXiv:0912.0201](#))
- Almosallam I. A., Jarvis M. J., Roberts S. J., 2016a, *MNRAS*, 462, 726
- Almosallam I. A., Lindsay S. N., Jarvis M. J., Roberts S. J., 2016b, *MNRAS*, 455, 2387
- Arnouts S., Cristiani S., Moscardini L., Matarrese S., Lucchin F., Fontana A., Giallongo E., 1999, *MNRAS*, 310, 540
- Behroozi P. S., Wechsler R. H., Wu H.-Y., 2013, *ApJ*, 762, 109
- Benítez N., 2000, *ApJ*, 536, 571
- Biviano A. et al., 2013, *A&A*, 558, A1
- Blanton M. R., Roweis S., 2007, *AJ*, 133, 734
- Blanton M. R. et al., 2005, *AJ*, 129, 2562
- Bonnett C., 2015, *MNRAS*, 449, 1043, *arXiv: 1312.1287*
- Bonnett C. et al., 2016, *Phys. Rev. D*, 94, 042005
- Brammer G. B., van Dokkum P. G., Coppi P., 2008, *ApJ*, 686, 1503
- Breiman L., Friedman J. H., Olshen R. A., Stone C. J., 1984, *Classification and Regression Trees, Statistics/Probability Series*. Wadsworth Publishing Company, Belmont, California, U.S.A
- Brescia M., Cavuoti S., Amaro V., Riccio G., Angora G., Vellucci C., Longo G., 2018, *ArXiv e-prints*
- Brescia M., Cavuoti S., D'Abrusco R., Longo G., Mercurio A., 2013, *ApJ*, 772
- Brescia M., Cavuoti S., Longo G., De Stefano V., 2014, *A&A*, 568
- Carrasco Kind M., Brunner R. J., 2013, *MNRAS*, 432, 1483
- Carrasco Kind M., Brunner R. J., 2014, *MNRAS*, 442, 3380
- Cavuoti S., Amaro V., Brescia M., Vellucci C., Tortora C., Longo G., 2017a, *MNRAS*, 465, 1959
- Cavuoti S., Brescia M., De Stefano V., Longo G., 2015, *Exp. Astron.*, 39, 45
- Cavuoti S. et al., 2017b, *MNRAS*, 466, 2039
- Chen T., Guestrin C., 2016, in *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, ACM, New York, NY, USA, pp. 785–794
- Dahlen T. et al., 2013, *ApJ*, 775, 93
- Dawson W. A., Schneider M. D., Tyson J. A., Jee M. J., 2016, *ApJ*, 816, 11
- de Jong J. T. A. et al., 2017, *A&A*, 604, A134
- de Jong J. T. A., Verdoes Kleijn G. A., Kuijken K. H., Valentijn E. A., 2013, *Exp. Astron.*, 35, 25
- Duncan K. J., Jarvis M. J., Brown M. J. I., Röttgering H. J. A., 2018, *Monthly Notices of the Royal Astronomical Society*, 940
- Fernández-Soto A., Lanzetta K. M., Yahil A., 1999, *ApJ*, 513, 34
- Firth A. E., Lahav O., Somerville R. S., 2003, *MNRAS*, 339, 1195
- Freeman P. E., Izbicki R., Lee A. B., 2017, *MNRAS*, 468, 4556
- Graff P., Feroz F., Hobson M. P., Lasenby A., 2014, *MNRAS*, 441, 1741
- Graham M. L., Connolly A. J., Ivezić Ž., Schmidt S. J., Jones R. L., Jurić M., Daniel S. F., Yoachim P., 2018, *AJ*, 155, 1
- Green J. et al., 2012, preprint ([arXiv:1208.4012](#))
- Hildebrandt H. et al., 2010, *A&A*, 523, A31
- Hofmann B., Mathé P., 2018, *Inverse Problems*, 34, 015007
- Ibert O. et al., 2006, *A&A*, 457, 841
- Ivezić Ž. et al., 2008, preprint ([arXiv:0805.2366](#))

¹⁹ available at: <http://ls.st/srd>

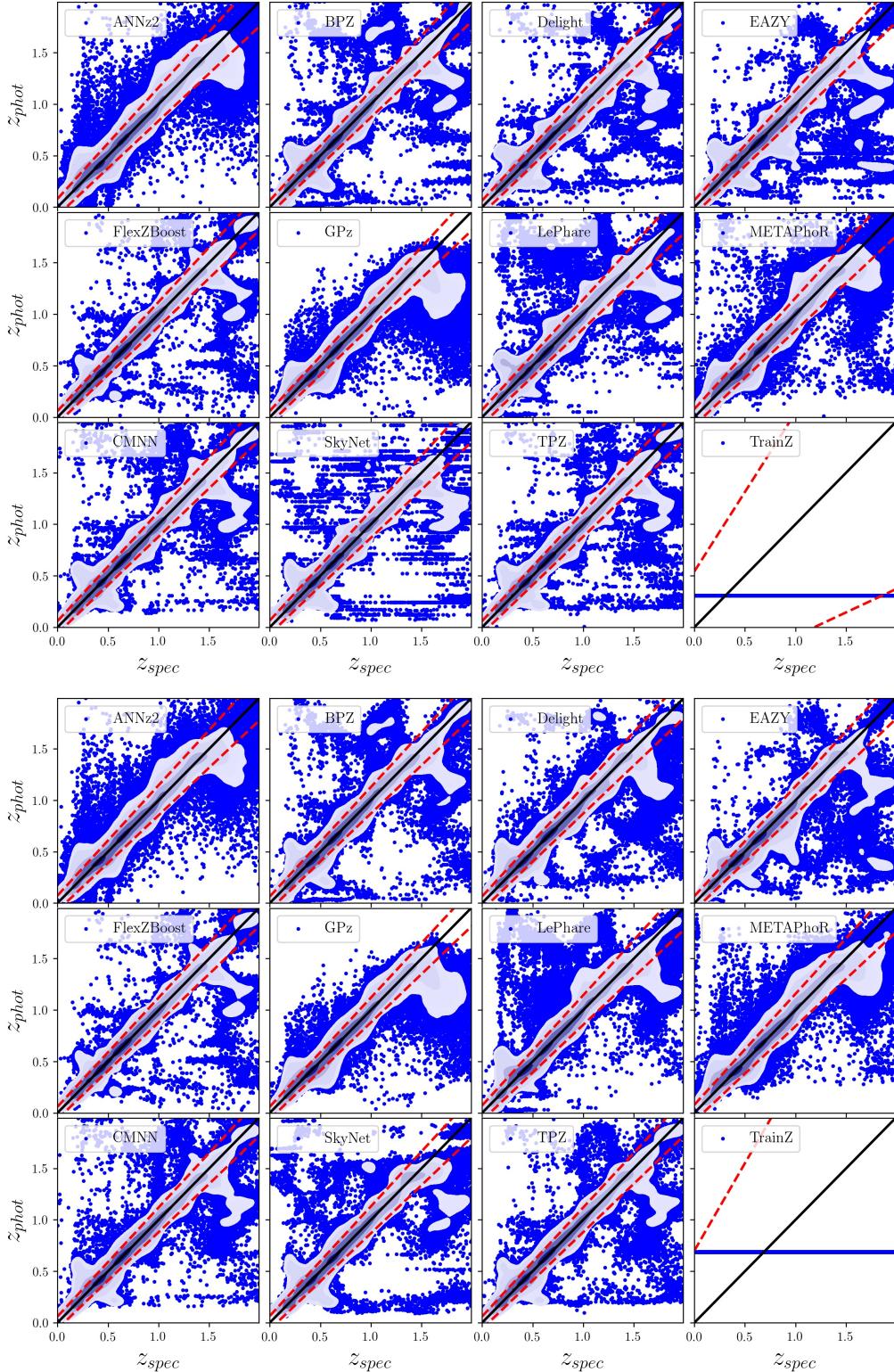


Figure A1. Point estimate photo-z's derived from the posteriors. Top panel shows z_{PEAK} , while bottom panel shows z_{WEIGHT} . Point estimate density is represented with fixed density contours, while outliers at lower density are represented by blue points. While use of point-estimate photo-z's is 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	Z_{PEAK}			Z_{WEIGHT}		
	$\frac{\sigma_{IQR}}{(1+z)}$	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	0.0202	-0.00091	0.036	0.0201	-0.000950	0.037
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

Izbicki R., Lee A. B., 2017, *Electron. J. Statist.*, 11, 2800

Izbicki R., Lee A. B., Freeman P. E., 2017, *Ann. Appl. Stat.*, 11, 698

Laigle C. et al., 2016, *ApJS*, 224, 24

Laureijs R. et al., 2011, preprint (1110.3193)

Leistedt B., Hogg D. W., 2017, *ApJ*, 838, 5

Maddox N., Hewett P. C., Warren S. J., Croom S. M., 2008, *MNRAS*, 386, 1605

Malz A., Hogg D., in prep., CHIPPR, chippr

Malz A., Marshall P., DeRose J., Graham M., Schmidt S., Wechsler R., 2018, *AJ*, Accepted

Mandelbaum R. et al., 2008, *MNRAS*, 386, 781

Massarotti M., Iovino A., Buzzoni A., 2001, *A&A*, 368, 74

Masters D. C., Stern D. K., Cohen J. G., Capak P. L., Rhodes J. D., Castander F. J., Paltani S., 2017, *ApJ*, 841, 111

Newman J. A., 2008, *ApJ*, 684, 88

Newman J. A. et al., 2015, *Astroparticle Physics*, 63, 81

Polsterer K. L., D'Isanto A., Gieseke F., 2016, preprint (arXiv:1608.08016)

Rasmussen C., Williams C., 2006, *Gaussian Processes for Machine Learning, Adaptative computation and machine learning series*. MIT Press, Cambridge, MA

Reddick R. M., Wechsler R. H., Tinker J. L., Behroozi P. S., 2013, *ApJ*, 771, 30

Riccio G., Brescia M., Cavuoti S., Mercurio A., di Giorgio A., Molinari S., 2017, *PASP*, 129

Richards G. T. et al., 2006, *ApJS*, 166, 470

Richards G. T. et al., 2015, *ApJS*, 219, 39

Sadeh I., Abdalla F. B., Lahav O., 2016, *PASP*, 128, 104502

Salvato M. et al., 2011, *ApJ*, 742, 61

Sánchez C. et al., 2014, *MNRAS*, 445, 1482

Scott D. W., 1992, *Multivariate Density Estimation. Theory, Practice, and Visualization*. Wiley

Skrutskie M. F. et al., 2006, *AJ*, 131, 1163

Tanaka M. et al., 2017, preprint (arXiv:1704.05988)

York D. G. et al., 2000, *AJ*, 120, 1579