

# HOW TO OBTAIN THE REDSHIFT DISTRIBUTION FROM PROBABILISTIC REDSHIFT ESTIMATES

ALEX MALZ<sup>1</sup>, DAVID W. HOGG<sup>1,2,3,4</sup>, PHIL MARSHALL<sup>5</sup>

*Draft version January 18, 2018*

## ABSTRACT

A trustworthy estimate of the redshift distribution  $n(z)$  is crucial for weak lensing cosmology. Spectroscopically confirmed redshifts for the dim and numerous galaxies of weak lensing surveys are expected to be inaccessible, making photometric redshifts (photo- $z$ s) the next best alternative. The nontrivial systematics affecting photo- $z$  estimation have motivated the weak lensing community to favor photo- $z$  probability distribution functions (PDFs) as a more comprehensive alternative to photo- $z$  point estimates. However, analytic methods for utilizing these new data products in cosmological inference are still evolving. This paper presents the Cosmological Hierarchical Inference with Probabilistic Photometric Redshifts (CHIPPR) model, a novel approach to estimating the  $n(z)$  from a catalog of galaxy photo- $z$  PDFs based upon a probabilistic graphical model of hierarchical inference. We present the publicly available CHIPPR code implementing this technique, as well as its validation on mock data. The CHIPPR model yields a more accurate characterization of  $n(z)$  containing information beyond the best-fit estimator produced by traditional procedures.

**Keywords:** cosmology: cosmological parameters — galaxies: statistics — gravitational lensing: weak — methods: data analysis — methods: statistical

## 1. INTRODUCTION

After a brief literature review addressing how photo- $z$  PDFs are currently used in cosmology, this paper aims to answer the following questions:

- Why should we question existing methods?
- How can we improve the effectiveness of using photo- $z$  PDFs in inference?
- How does the result of CHIPPR compare to established estimators in terms of the accuracy of  $n(z)$ ?
- How significant is the effect of the discrepancy between  $n(z)$  estimators on cosmological constraints?

## 2. METHOD

This paper presents a mathematically consistent method for obtaining the posterior probability distribution over the redshift distribution  $n(z)$  using a catalog of photo- $z$  PDFs. We start by introducing some nomenclature, definitions, and symbols that will be used throughout Secs. 2.1, 2.2, and 2.3.

We shall say the parameters comprising  $\vec{\phi}$  define the redshift distribution under some functional form about which this method is agnostic, so long as the functional form evaluated with some true values of the parameters accurately describes the true  $n(z)$ . Since the redshift

distribution is itself a probability distribution, it may be written as  $p(z|\vec{\phi})$ .

Before jumping into the details of the model, it is crucial to settle on an interpretation of what a photo- $z$  PDF actually is. When a photo- $z$  PDF is reported, it is an object containing information about the redshift  $z_i$  of galaxy  $i$  based on its photometric data  $\vec{d}_i$ , making a photo- $z$  PDF a *posterior* probability distribution. The data  $\vec{d}_i$  may be fluxes, magnitudes, colors, or any combination thereof. We assume that the redshifts of the galaxies are not related to the redshifts or photometry of other galaxies in the survey, i.e. they are *independent*.

In addition to the data  $\vec{d}_i$ , the photo- $z$  PDF also contains information imparted by assumptions that went into the process by which the photo- $z$  was made. Those assumptions are only relevant inasmuch as they impose a preference on the redshift distribution. Thus, they can be reduced to a sort of prior on the distribution of redshifts, which can be parametrized under the chosen functional form by parameters  $\vec{\phi}^*$ . The redshift distribution  $p(z|\vec{\phi}^*)$  associated with the parameters in  $\vec{\phi}^*$  shall be called the *interim prior* because in many cases it was our best guess as to the true redshift distribution before observing any data. For example, in template-based photo- $z$  PDF methods, the interim prior is a linear combination of the redshift distributions of the classes of templates based on previous observations. (In training-based photo- $z$  PDF methods, the interim prior is related to the redshift distribution of the training set and thus may not be shared among all galaxies for which photo- $z$  PDFs are produced; however, in this work, we will assume that it is shared among all galaxies in a survey.) Because the choice of parameters  $\vec{\phi}^*$  is not causally connected to the data or the true redshift distribution but nonetheless contributes to the photo- $z$  PDF, we call it an *interim prior*.

A photo- $z$  PDF is then an *interim posterior* probability

aimalz@nyu.edu

<sup>1</sup> Center for Cosmology and Particle Physics, Department of Physics, New York University, 726 Broadway, 9th floor, New York, NY 10003, USA

<sup>2</sup> Simons Center for Computational Astrophysics, 162 Fifth Avenue, 7th floor, New York, NY 10010, USA

<sup>3</sup> Center for Data Science, New York University, 60 Fifth Avenue, 7th floor, New York, NY 10003, USA

<sup>4</sup> Max-Planck-Institut für Astronomie, Königstuhl 17, D-69117 Heidelberg, Germany

<sup>5</sup> [SLAC]

distribution of the redshift of a galaxy given its observed photometric data and the interim prior:  $p(z_i|\vec{d}_i, \vec{\phi}^*)$ . We present the two attributes of a probabilistic graphical model (PGM) for the redshift distribution  $n(z)$ : the directed acyclic graph (DAG) and its mathematical interpretation in terms of Bayesian hierarchical inference.

The fundamental assumption underlying the concept of photo- $z$  estimation is that each galaxy  $i$  has some observed photometric data  $\vec{d}_i$  (fluxes, magnitudes, or colors) that is drawn from a function of its redshift  $z_i$ , which is a parameter in this model. This function constitutes a forward model for the observations. The redshifts  $\{z_i\}$  for all galaxies in a survey are random draws from the redshift distribution, whose parameters  $\vec{\phi}$  under a chosen functional form are called *hyperparameters* because they are shared among all  $N$  galaxies in the survey. Fig. 1 illustrates these relationships. PGMs of this structure are *hierarchical* in that, though the data are only directly influenced by their redshift parameters, we can still use them to infer something about the global hyperparameters by way of a higher level in the graph.

The DAG of Fig. 1 translates directly into a mathematical expression for the posterior distribution  $p(\vec{\phi}|\{\vec{d}_i\})$  of the hyperparameters given the entire set of interim posteriors  $\{p(z_i|\vec{d}_i, \vec{\phi}^*)\}$  by way of the following derivation. From this point forward, we will use log probabilities.

According to Bayes' rule, the posterior distribution of interest is

$$\log[p(\vec{\phi}|\{\vec{d}_i\})] \propto \log[p(\{\vec{d}_i\}|\vec{\phi})] + \log[p(\vec{\phi})], \quad (1)$$

where  $p(\vec{\phi})$  is the *hyperprior* probability distribution over possible values of the hyperparameters. The hyperprior is a choice made by those performing the inference and is not to be confused with the interim prior parameters that influence the photo- $z$  PDF data product.

Focusing on the likelihood term  $p(\{\vec{d}_i\}|\vec{\phi})$ , we employ the power of the hierarchical model through *marginalization* of the unobserved redshift parameters by integrating over them to obtain

$$\begin{aligned} \log[p(\{\vec{d}_i\}|\vec{\phi})] &= \log \left[ \int \exp \left[ \log[p(\{\vec{d}_i\}|\{z_i\})] \right. \right. \\ &\quad \left. \left. + \log[p(\{z_i\}|\vec{\phi})] \right] d\{z_i\} \right]. \end{aligned} \quad (2)$$

The redshifts  $\{z_i\}$  are statistically independent of each other, and each galaxy's photometry depends solely on its own redshift, so we may invoke statistical independence to get

$$\log[p(\{z_i\}|\vec{\phi})] = \sum_i^N \log[p(z_i|\vec{\phi})] \quad (3)$$

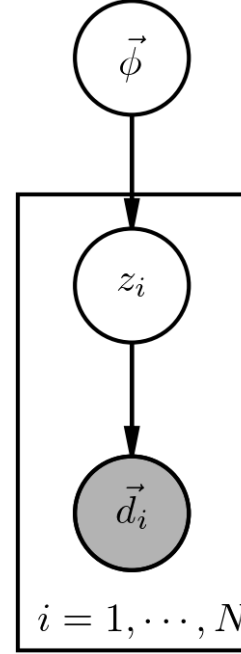
and

$$\log[p(\{\vec{d}_i\}|\{z_i\})] = \sum_i^N \log[p(\vec{d}_i|z_i)]. \quad (4)$$

Having chosen the functional form mapping  $\vec{\phi}$  to  $n(z)$ , the terms on the righthand side of Eq. 3 are known. However, the terms on the righthand side of Eq. 4 are

We would like to learn about  $p(z|\vec{\phi})$  from a catalog of photo- $z$  interim posteriors.

### 2.1. Model



**Figure 1.** This directed acyclic graph corresponds to a probabilistic graphical model for a Bayesian hierarchical inference of  $p(\vec{\phi}|\{\vec{d}_i\})$ . In this graph, all random variables are shown in circles, with observed random variables shown in shaded circles. Relationships between variables are indicated by arrows from parameters to the variables distributed according to functions of them. The box indicates that there are a number of copies of the relationships between boxed parameters, each independent of all others. The hyperparameters comprising  $\vec{\phi}$ , which define  $n(z)$ , are at the top. Independently drawn from a function of the hyperparameters  $\vec{\phi}$  are galaxy redshifts  $\{z_i\}$  below. The observed galaxy photometry  $\{\vec{d}_i\}$ , shown in shaded circles, is determined by the redshifts above.

likelihoods, and the photo- $z$  PDFs we have are interim posteriors.

To transform the interim posteriors we have into the likelihoods we need, we start with the vacuously true statement

$$\begin{aligned} \log[p(\vec{d}_i|z_i)] &= \log[p(\vec{d}_i|z_i)] + \log[p(z_i|\vec{d}_i, \vec{\phi}^*)] \\ &\quad - \log[p(z_i|\vec{d}_i, \vec{\phi}^*)]. \end{aligned} \quad (5)$$

We expand the last term using Bayes' Rule to get

$$\begin{aligned} \log[p(\vec{d}_i|z_i)] &= \log[p(\vec{d}_i|z_i)] + \log[p(z_i|\vec{d}_i, \vec{\phi}^*)] \\ &\quad - \log[p(\vec{d}_i|z_i, \vec{\phi}^*)] \\ &\quad + \log[p(\vec{d}_i|\vec{\phi}^*)] - \log[p(z_i|\vec{\phi}^*)]. \end{aligned} \quad (6)$$

Fig. 1 has no direct connection between  $\vec{d}_i$  and  $\vec{\phi}$ , so we can break up  $\log[p(\vec{d}_i|z_i, \vec{\phi}^*)]$  as follows:

$$\begin{aligned} \log[p(\vec{d}_i|z_i)] &= \log[p(\vec{d}_i|z_i)] + \log[p(z_i|\vec{d}_i, \vec{\phi}^*)] \\ &\quad - \log[p(\vec{d}_i|z_i)] - \log[p(\vec{d}_i|\vec{\phi}^*)] \\ &\quad + \log[p(\vec{d}_i|\vec{\phi}^*)] - \log[p(z_i|\vec{\phi}^*)]. \end{aligned} \quad (7)$$

Canceling terms, we find

$$\log[p(\vec{d}_i|z_i)] = \log[p(z_i|\vec{d}_i, \vec{\phi}^*)] - \log[p(z_i|\vec{\phi}^*)]. \quad (8)$$

Now we combine Eqs. 1, 2, 3, 4, and 8 to arrive at

$$\begin{aligned} \log[p(\vec{\phi}|\{\vec{d}_i\})] &\propto \log[p(\vec{\phi})] + \log \left[ \int \exp \left[ \sum_i^N ( \right. \right. \\ &\quad \left. \left. + \log[p(z_i|\vec{\phi})] - \log[p(z_i|\vec{\phi}^*)] \right. \right. \\ &\quad \left. \left. + \log[p(z_i|\vec{d}_i, \vec{\phi}^*)] \right) \right] d\{z_i\} \right]. \end{aligned} \quad (9)$$

Thus, if we have a catalog of photo- $z$  interim posteriors  $\{p(z_i|\vec{d}_i, \vec{\phi}^*)\}$ , the interim prior parameters  $\vec{\phi}^*$ , and a prior distribution  $p(\vec{\phi})$ , we may find the posterior  $\log[p(\vec{\phi}|\{\vec{d}_i\})]$  on the redshift distribution  $n(z)$ .

This framework entails a number of choices and assumptions that must be addressed explicitly:

1. The chosen functional form of  $n(z)$  with parameters  $\vec{\phi}$  must be capable of describing the true redshift distribution.
2. The photometric data and redshift for each galaxy must be independent of the photometric data and redshifts of all other galaxies.
3. There is one interim prior shared among all galaxies in the survey, and its parameters with known parameters  $\vec{\phi}^*$  are known.
4. We must choose a hyperprior probability distribution  $p(\vec{\phi})$  over the hyperparameters.
5. While we advocate for the approach of hierarchical inference, the probabilistic graphical model presented here is not the only one that could be proposed.

## 2.2. Alternative Approaches

It is useful to translate some popular existing methods for deriving  $n(z)$  from photo- $z$  PDFs into the mathematical framework of Sec. 2.1. We briefly discuss the conditions under which they are equivalent to Eq. 9, which are elaborated upon in the Appendix.

### 2.2.1. Stacking

The most common way to combine photo- $z$  interim posteriors into an estimator  $\hat{n}_{stack}(z)$  for  $n(z)$  is to "stack" them, which corresponds to

$$\hat{n}^{stack}(z) = \frac{1}{N} \sum_i^N \exp \left[ \log[p(z_i|\vec{d}_i, \vec{\phi}^*)] \right]. \quad (10)$$

### 2.2.2. Point Estimation

In some cases, estimators  $\hat{n}^{point}(z)$  are obtained by assuming the photo- $z$  PDFs are effectively delta functions by reducing them to point estimates according to

$$p(z_i|\vec{d}_i, \vec{\phi}^*) \approx \delta(z, \hat{z}_i^{point}) \quad (11)$$

before applying Eq. 10. The most popular redshift point estimators are the mean  $\hat{z}_i^{mean}$ , median  $\hat{z}_i^{median}$ , and mode  $\hat{z}_i^{mode}$  of the original photo- $z$  interim posterior.

### 2.3. Implementation

The publicly available **CHIPPR** code implements the probabilistic graphical model presented in Sec. 2.1.

In addition to the choices and assumptions underlying the probabilistic graphical model, the implementation of **CHIPPR** makes choices and assumptions of its own.

1. **CHIPPR** is only applicable if the redshift interim posteriors and interim prior are accurate.
2. **CHIPPR** currently only accepts photo- $z$  PDFs and produces  $n(z)$  samples of a single format, that of the piecewise constant parametrization, also referred to as a binned histogram parametrization and a sum of top hat functions.

## 3. VALIDATION ON SIMPLE MOCK DATA

We demonstrate the superiority of **CHIPPR** over alternative approaches in a number of compelling test cases on mock data. Each experiment is characterized by a single change to a fiducial case in order to isolate the influence of systematic effects known to be relevant to photo- $z$  estimation and propagation in analysis.

### 3.1. Mock Data & Metrics

The mock data in these tests consists of photo- $z$  interim posteriors rather than photometric data because the various existing methods for deriving photo- $z$  interim posteriors do not in general yield results that are consistent with one another, indicating that their systematics are not well-understood. Because the mock data is independent of any choice of photo- $z$  PDF production method, we not only ensure that our photo- $z$  interim priors are perfectly understood but also deter readers from assuming that **CHIPPR** has any preference over the method by which photo- $z$  interim posteriors are derived from photometric data.

#### 3.1.1. Mock Data

The mock data used here are produced by the following steps.

1. Choose a true  $n(z)$  that is a continuous function with known parameters  $\vec{\phi}'$ .
2. Sample the true  $n(z)$  to create a catalog of  $N$  true redshifts  $z'_i$ .
3. Choose a true intrinsic scatter parameter  $\sigma$ , and sample Gaussians  $p(z''_i|z'_i, \sigma) = \mathcal{N}(z'_i, \sigma)$  to get "observed" redshifts  $z''_i$  defining Gaussian likelihoods  $p(\vec{d}_i|z''_i, \sigma) = \mathcal{N}(z''_i, \sigma)$ .

4. Choose a parametrization and the parameters  $\vec{\phi}^*$  of the interim prior  $p(z_i|\vec{\phi}^*)$ .
5. Create interim posteriors  $\log[p(z_i|\vec{d}_i, \vec{\phi}^*)] = \log[p(\vec{d}_i|z_i)] \log[p(z_i|\vec{\phi}^*)]$  in this parametrization.

In all of the following validation tests, we use a piecewise constant parametrization in log-space with 10 bins.  $N = 10,000$  galaxies. In Secs. 3.2 and 3.3, we use a flat interim prior. This method for deriving mock data is referred to as the fiducial case, and variations on it will refer directly to the steps that are altered.

### 3.1.2. $n(z)$ Accuracy Metric

The Kullback-Leibler divergence is our primary measure of the accuracy of estimators of  $n(z)$  in cases of mock data with known true redshifts.

**Review precision and bias from `kld.ipynb` and interpret in terms of a % difference.**

### 3.2. Underlying $n(z)$ Effects

Existing  $n(z)$  estimators are systematically smoother than the true  $n(z)$ . Here we show that the traditional estimators perform better when the true  $n(z)$  is weakly featured than when the true  $n(z)$  is strongly featured by experimenting with Step 1 of Sec. 3.1.1. The implication of this issue is quite serious; the consistently smooth, unimodal  $n(z)$  estimates appearing in the literature could result from much more featured true redshift distributions, and there would be no way to catch this error without using a fully probabilistic method.

#### 3.2.1. Featureless $n(z)$

In this test, we choose a weakly featured true  $n(z)$  of Fig. ?? with a smooth, unimodal shape, based on the interim prior used for the SDSS DR8 photo- $z$  PDFs. [check and include citation!]

#### 3.2.2. Featured $n(z)$

In this test, we choose the true  $n(z)$  of Fig. 3 with nontrivial structure.

This featured true  $n(z)$  will henceforth be referred to as the fiducial  $n(z)$ .

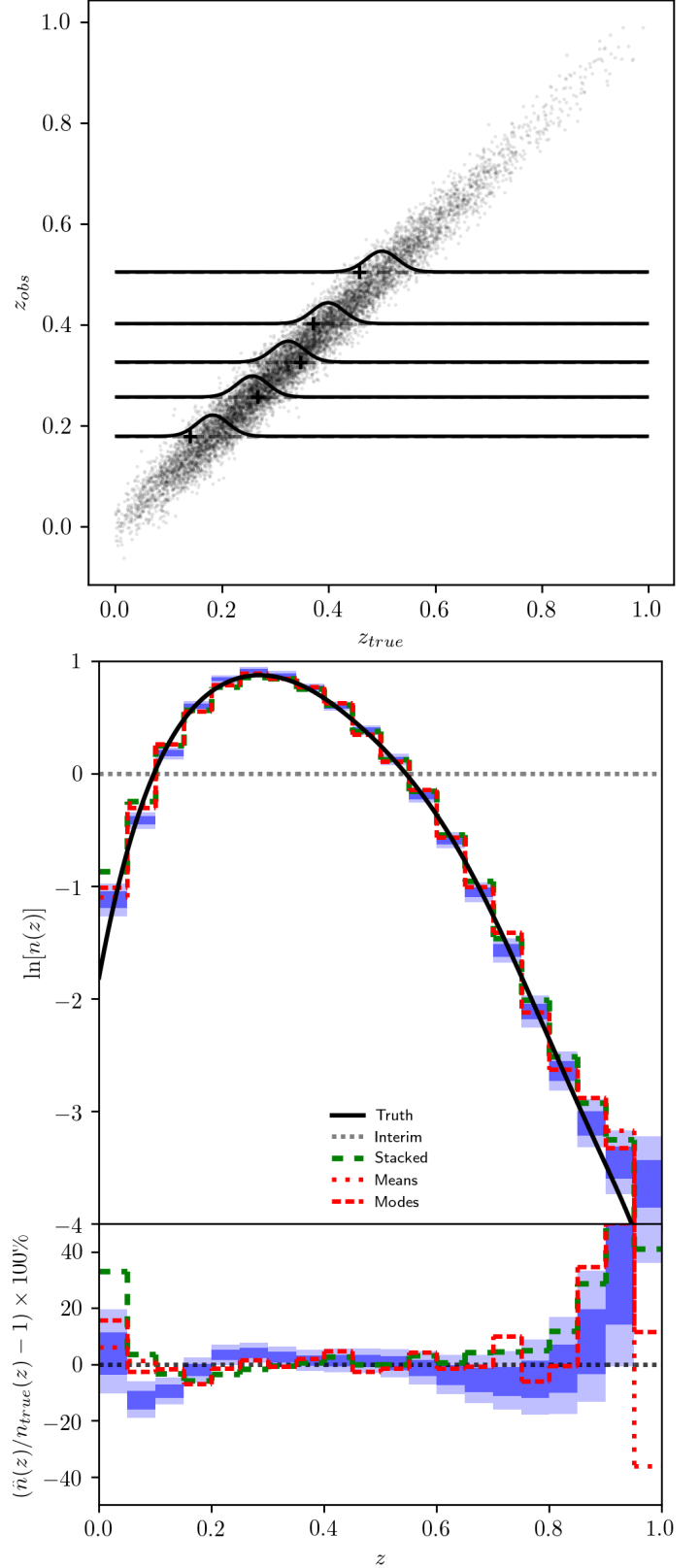
### 3.3. Emulated Data Quality Effects

In the following test cases, we vary the properties of the mock photo- $z$  likelihoods in an effort to emulate known systematics in photo- $z$  estimation. These tests vary Step 3 of Sec. 3.1.1.

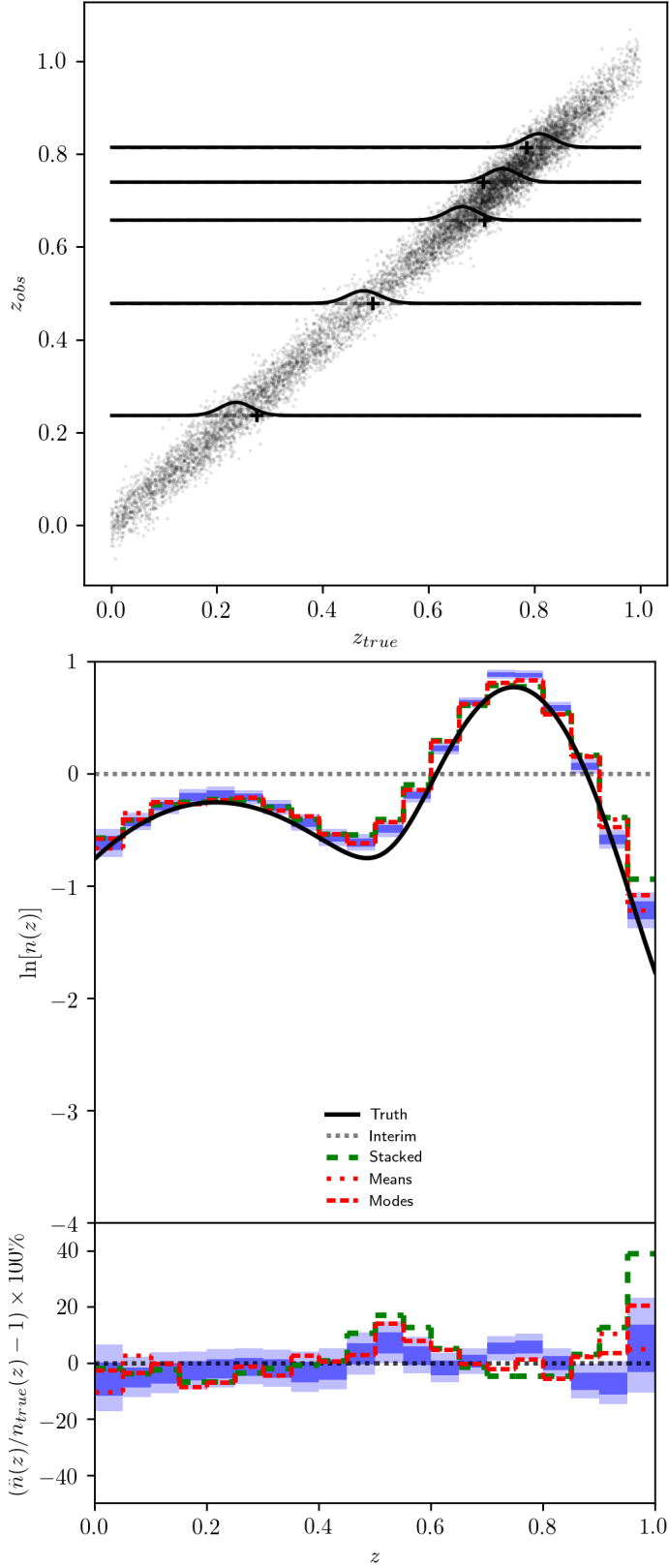
#### 3.3.1. Intrinsic Scatter

One major concern about photo- $z$ s is the intrinsic scatter of point estimators, including those derived from photo- $z$  PDFs, that is observed to varying extents with every existing photo- $z$  algorithm and illustrated in Fig. 4. To emulate intrinsic scatter, we modify the fiducial case to simply broaden the single Gaussian component of the likelihood. To enforce self-consistency, the mean is a random variable drawn from a Gaussian distribution with the newly increased variance.

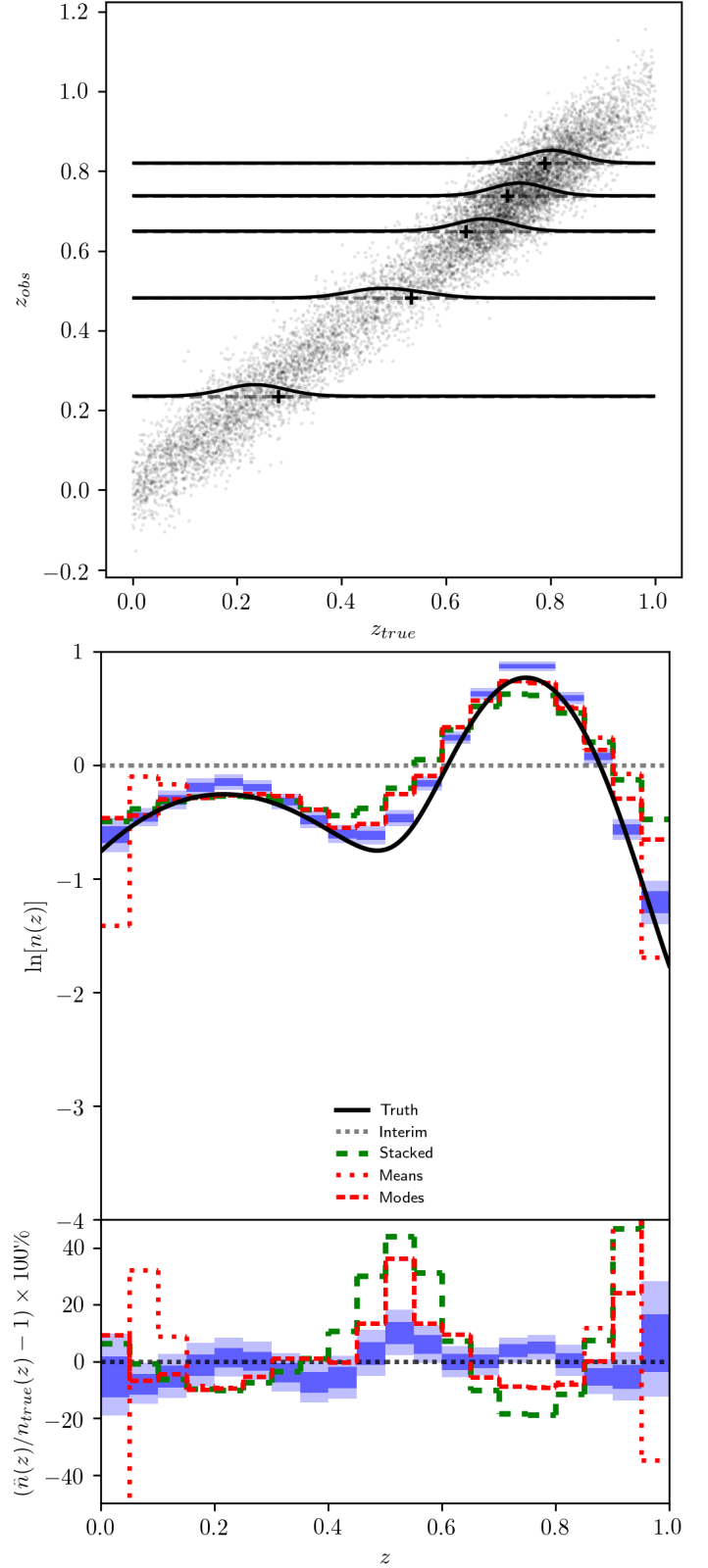
To emulate intrinsic scatter, we modify the fiducial case to simply broaden the single Gaussian component of the likelihood. To enforce self-consistency, the mean is a random variable drawn from a Gaussian distribution with the newly increased variance.



**Figure 2.** All estimators perform well when the true  $n(z)$  is well behaved, exhibiting significant deviation only when  $n(z)$  is very small, as the sample size of true redshifts in that range will be small. Top panel: The traditional MAP reduction of a photo- $z$  PDF against the true redshifts with a few rescaled photo- $z$  interim posteriors are overplotted in solid lines, with a dotted line indicating zero probability. Bottom panel: Various estimators of  $\ln[n(z)]$ , the interim prior, and the true  $\ln[n(z)]$  as a continuous function and under a binned parametrization.



**Figure 3.** [This plot doesn't really show what I want because the intrinsic scatter is too low! I used  $\sigma = 0.03$  because that's what has been quoted as what LSST, etc. needs, but that's not what comes out of photo- $z$  estimation methods before they impose aggressive cuts. . .] Top panel: The traditional MAP reduction of a photo- $z$  PDF against the true redshifts with a few rescaled photo- $z$  interim posteriors are overplotted in solid lines, with a dotted line indicating zero probability. Bottom panel: Various estimators of  $\ln[n(z)]$ , the interim prior, and the true  $\ln[n(z)]$  as a continuous function and under a binned parametrization.



**Figure 4.** As the intrinsic scatter increases, the discrepancy between estimators increases. In particular, the stacked estimator and marginalized point estimators predict  $\ln[n(z)]$  to be smoother than the truth, while the [This would be a lot more compelling with more galaxies. Also, the weird edge effects in the top panel are real, because the point estimator is the MAP, not the center of the Gaussian likelihood, and there's no requirement that the mean of the likelihood be within the true redshift range.] Top panel: The traditional MAP reduction of a photo- $z$  PDF against the true redshifts with a few rescaled photo- $z$  interim posteriors are overplotted in solid lines, with a dotted line indicating zero probability. Bottom panel: Various estimators of  $\ln[n(z)]$ , the interim prior, and the true  $\ln[n(z)]$  as a continuous function and under a binned parametrization.

### 3.3.2. Template-like Catastrophic Outliers

In addition to intrinsic scatter, photo- $z$  methods employing template fitting tend to produce catastrophic outliers that are distributed to be broad in  $z_{\text{spec}}$  and narrow in  $z_{\text{phot}}$ , as in Fig. 5. The systematic behind these catastrophic outliers may be described as an attractor in the space of  $z_{\text{phot}}$ ; some galaxies at a range of  $z_{\text{spec}}$  map onto a single  $z_{\text{phot}}$  (with some scatter) if their true SED does not have sufficiently strong features (as is the case for blue galaxies), leading galaxies of that SED type at many  $z_{\text{spec}}$  to have similar colors.

### 3.3.3. Training-like Catastrophic Outliers

Data driven photo- $z$  methods tend to suffer from a different form of catastrophic outliers that are distributed to be narrow in  $z_{\text{spec}}$  and broad in  $z_{\text{phot}}$ , as in Fig. 6. The systematic behind these catastrophic outliers may be described as an attractor in the space of  $z_{\text{spec}}$ ; some galaxies near a particular  $z_{\text{spec}}$  map to a range of  $z_{\text{phot}}$  if the training set galaxies at that  $z_{\text{spec}}$  have inconsistent  $z_{\text{phot}}$ , as might occur if their SED's features fall between photometric filters.

### 3.4. Emulated Interim Prior Effects

The interim prior encapsulates the the relationship between observed photometry and redshift information upon which a photo- $z$  estimate is based. Interim priors are in general not identical to the true  $n(z)$  we wish to estimate; if they were, we would not need any data! For template fitting photo- $z$  methods, the interim prior is usually an input chosen by the researcher. However, for machine learning methods, the interim prior is some function of the training set data that in many cases may be influenced by random numbers and is rarely output with the redshift estimates. Interim priors for template fitting methods tend to have incomplete coverage in the space of true photometry, because they are limited by the choice of the library of SEDs. Interim priors for machine learning methods tend to have incomplete coverage in the space of redshifts, because there are fewer galaxies with spectroscopically confirmed redshifts at high redshifts than low redshifts.

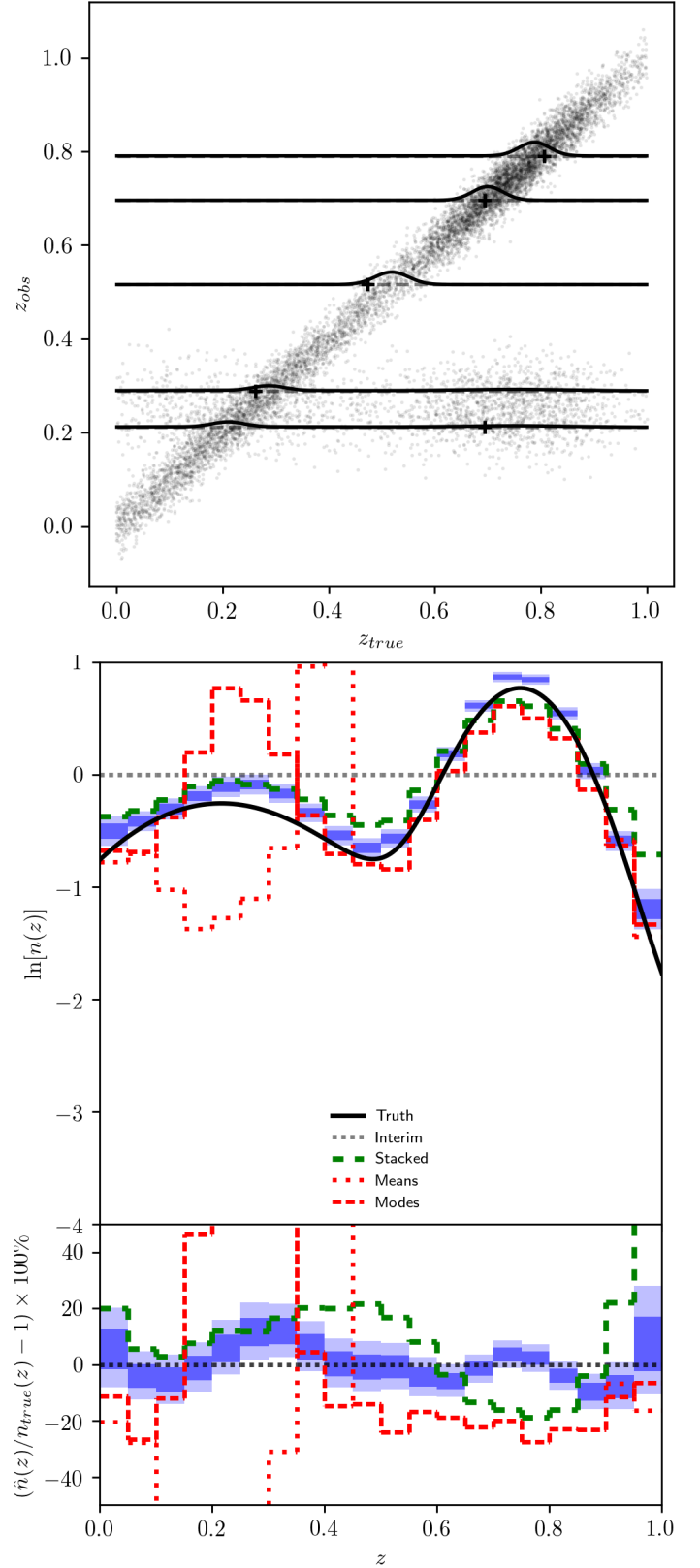
Existing  $n(z)$  estimation routines will always produce a biased estimator when the interim prior is not equal to the true  $n(z)$ . We demonstrate here that regardless of the appropriateness of the interim prior as an approximation to the true  $n(z)$ , CHIPPR is not affected by the choice of the interim prior so long as it has nontrivial coverage in the space of redshift. These tests modify Step 4 of Sec. 3.1.1.

#### 3.4.1. Template-like Interim Prior

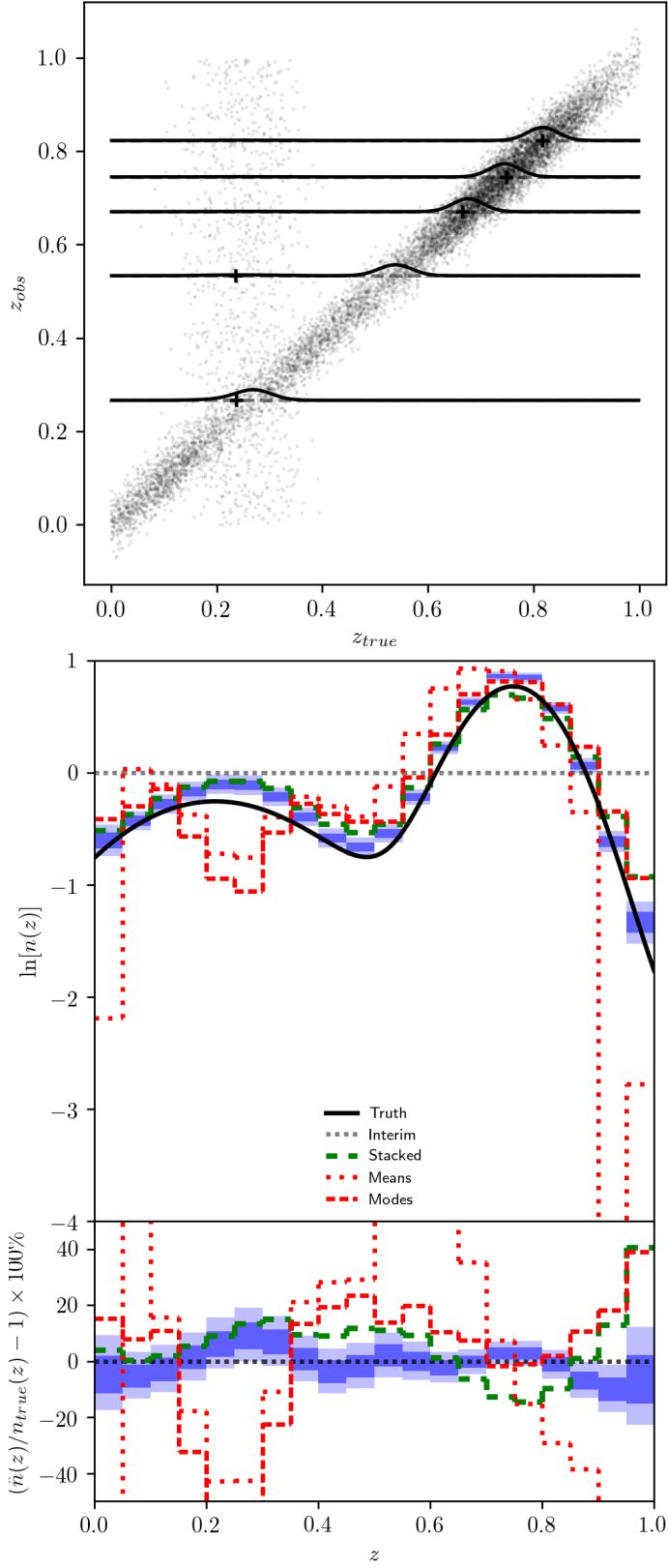
An interim prior based on a template library may be a sum of smooth functions representing  $n(z)$  for each SED type in the library. Template libraries do not include every possible galaxy SED, and the  $n(z)$  used for each SED type may not be accurate. The interim prior shown in Fig. 7 is an emulation of an interim prior corresponding to a template library of this type.

#### 3.4.2. Training-like Interim Prior

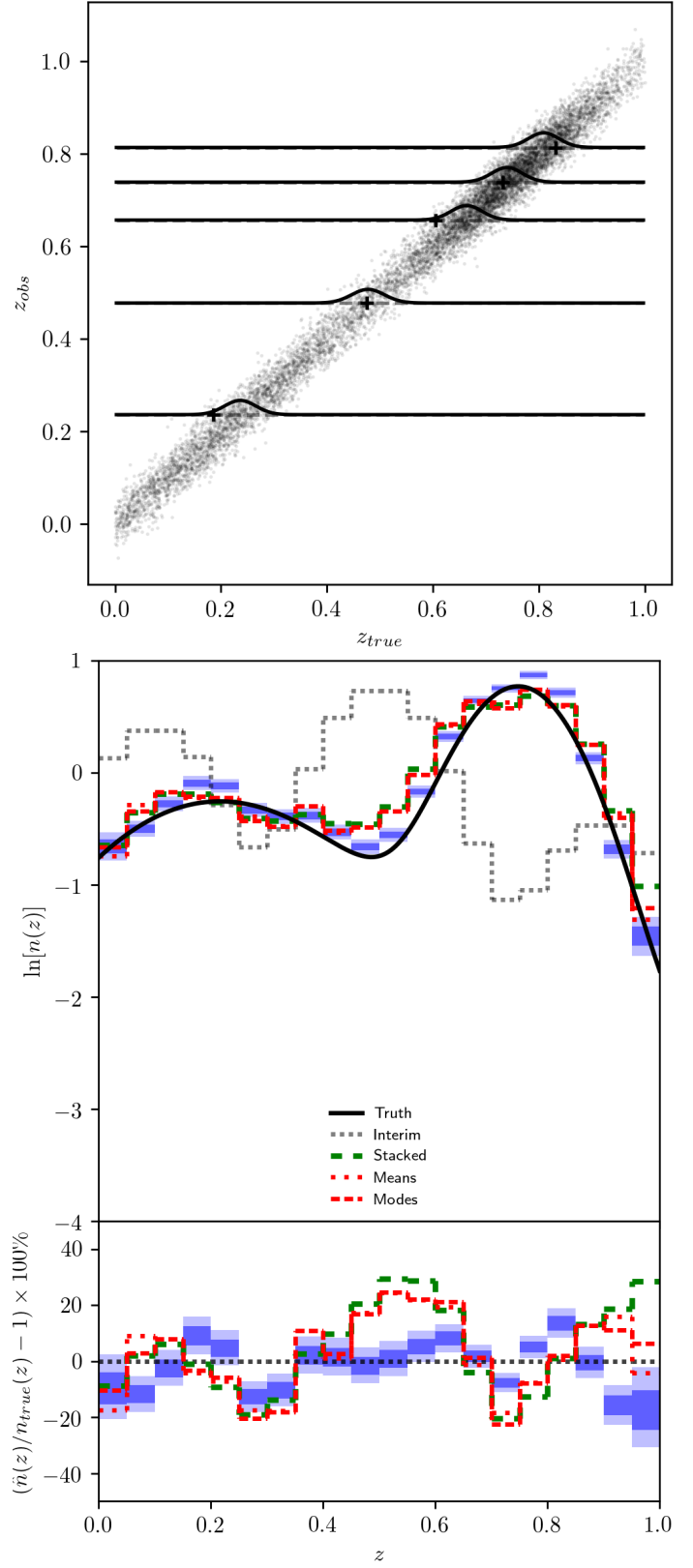
An interim prior based on a training set may be biased toward low redshifts due to the dearth of distant galaxies



**Figure 5.** Top panel: The traditional MAP reduction of a photo- $z$  PDF against the true redshifts with a few rescaled photo- $z$  interim posteriors are overplotted in solid lines, with a dotted line indicating zero probability. Bottom panel: Various estimators of  $\ln[n(z)]$ , the interim prior, and the true  $\ln[n(z)]$  as a continuous function and under a binned parametrization.



**Figure 6.** [This has proven to be the most challenging test case to implement, and there's clearly still a bug in the function that makes the likelihoods.]



**Figure 7.** [The case of a multimodal interim prior was a very compelling test in the previous version but somehow isn't anymore.]



with spectroscopic redshifts. The interim prior shown in Fig. 8 is an emulation of an interim prior corresponding to a training set biased in this way. We chose an interpolation of the interim prior used for the SDSS DR8 photo- $z$  PDFs.

#### 4. APPLICATION TO REALISTIC MOCK DATA

To show how the choice of  $n(z)$  estimator propagates to cosmological constraints, we apply CHIPPR to a data from a realistic cosmological simulation (probably **Buzzard**) with photo- $z$  PDFs produced by a popular method (probably BPZ).

##### 4.1. Mock Data & Metrics

As in Sec. 3.1, the mock data takes the form of photo- $z$  interim posteriors, but that is where the similarity ends. These photo- $z$  interim posteriors are derived from the photometry resulting from the **Buzzard** simulation by way of BPZ. Because the simulation begins with setting true values of the cosmological parameters, we can propagate the different estimators of  $n(z)$  through a forecasting code (which one?) to generate error ellipses on the cosmological parameters.

##### 4.1.1. Mock Data

We summarize the details of the **Buzzard** simulation here.

##### 4.1.2. Cosmological Constraint Metric

Because the realistic mock data of Sec. 4.1.1 is associated with true values of the cosmological parameters, we may compare the quality of cosmological constraints under different estimators of  $n(z)$ , creating a figure like Fig. 9. We perform this analysis using a forecasting code that takes  $n(z)$  and produces projected error ellipses in the space of cosmological parameters.

##### 4.2. Results

#### 5. DISCUSSION

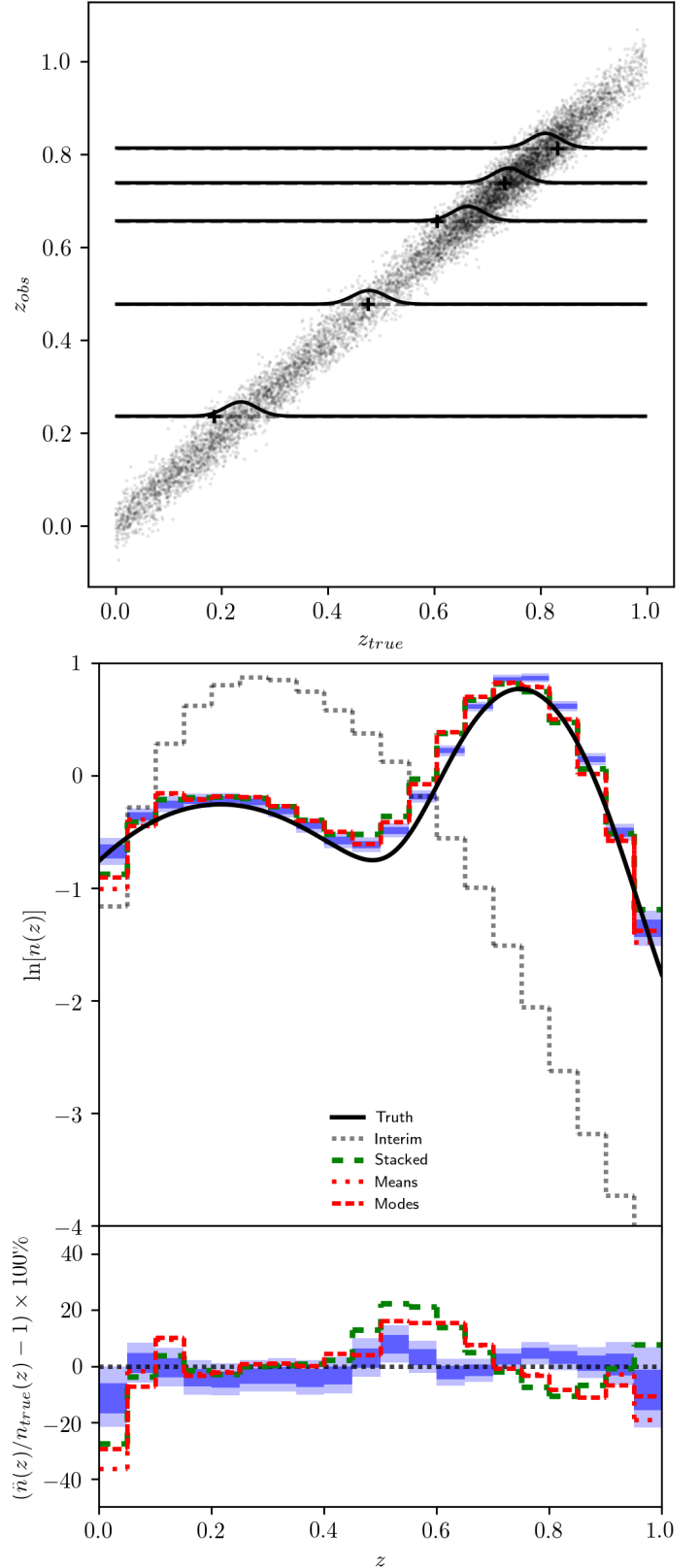
The results of Fig 10 have significant implications for the developing data analysis pipelines of next-generation telescope surveys. However, the method presented here has its own limitations, which are reiterated to discourage the community from applying this work inappropriately.

We intend to pursue a number of extensions of the work presented in this paper in future work.

#### 6. CONCLUSION

We now summarize answers to the questions posed in the introduction:

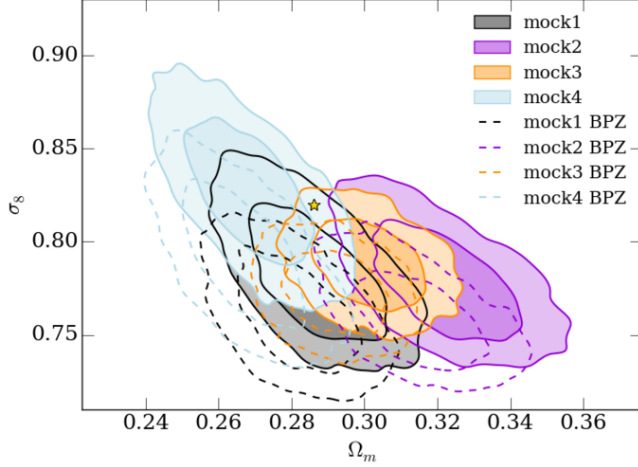
- Existing  $n(z)$  estimation methods produce biased estimators that propagate to inaccuracies in characterizing the cosmological parameters.
- Photo- $z$  PDFs are probabilistic data products so must be handled in a mathematically consistent manner such as the probabilistic graphical model outlined in this paper.
- In addition to coming with its own error distribution, the  $n(z)$  estimator produced by CHIPPR is quantifiably more accurate than established estimators.



**Figure 8.** [The case of a low- $z$  favoring interim prior was a very compelling test in the previous version but somehow isn't anymore.]



Not my figure! (N. MacCrann via J. DeRose)



**Figure 9.** [I'd like to be able to make a plot like this with the different  $n(z)$  estimators produced by `chippr` from the same set of photo- $z$  PDFs. Then I can calculate some metrics of the accuracy and precision of the error distributions relative to the true values of the cosmological parameters that produced the mock data.]

**Figure 10.** [moneyplot of error ellipses resulting from different  $n(z)$  estimators]

- Propagation of the CHIPPR result leads to a quantifiable improvement in the constraints on cosmological parameters.

In conclusion, we discourage the community from continuing to use the stacked estimator and reductions of photo- $z$  PDFs to redshift point estimates in obtaining estimators of  $n(z)$ . CHIPPR is freely available to the community for incorporation into evolving data analysis pipelines.

## APPENDIX

### Catalog Production

Sec. 3.1.1 outlines one way to make mock data, but it's not necessarily a logically consistent way to do it, a problem identified in trying to use it for the test cases in Sec. 3.3. If the goal is to emulate the traditional  $z_{\text{spec}}$  vs.  $z_{\text{phot}}$  scatterplots, we need to start from probability distributions in that space. This means there are two parameters  $z_{\text{spec}}$  and  $z_{\text{phot}}$  drawn from a two-dimensional distribution  $p(z_{\text{spec}}, z_{\text{phot}} | \underline{\sigma}, \vec{\phi})$ , where  $\underline{\sigma}$  contains parameters concerning the relationship between  $z_{\text{spec}}$  and  $z_{\text{phot}}$ , and  $\vec{\phi}$  contains hyperparameters concerning  $n(z)$ . We know from basic probability that

$$p(z_{\text{spec}}, z_{\text{phot}} | \underline{\sigma}, \vec{\phi}) = p(z_{\text{phot}} | z_{\text{spec}}, \underline{\sigma}, \vec{\phi}) p(z_{\text{spec}} | \underline{\sigma}, \vec{\phi}). \quad (1)$$

Because  $z_{\text{spec}}$  depends only on  $\vec{\phi}$  and  $z_{\text{phot}}$  is independent of  $\vec{\phi}$ , this becomes

$$p(z_{\text{spec}}, z_{\text{phot}} | \underline{\sigma}, \vec{\phi}) = p(z_{\text{phot}} | z_{\text{spec}}, \underline{\sigma}) p(z_{\text{spec}} | \vec{\phi}). \quad (2)$$

We know how to make reasonable choices for both terms on the righthand side of Eq. 2.

In the  $z_{\text{spec}}$  vs.  $z_{\text{phot}}$  scatterplot,  $p(z_{\text{phot}} | z_{\text{spec}}, \underline{\sigma}, \vec{\phi})$  would be represented by vertical slices, but we are interested in the horizontal slices that represent  $p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi})$ . Now we must distinguish the true hyperparameters  $\vec{\phi}_{\text{true}}$  from the interim hyperparameters  $\vec{\phi}_{\text{int}}$ . If the horizontal slices carry information about  $\vec{\phi}$ , how can we get a catalog of  $p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi}_{\text{int}})$  that's consistent with  $p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi}_{\text{true}})$ ? Somehow that information needs to be propagated through in order to recover it.

We achieve this goal by dividing out the true  $n(z)$  and multiplying by the interim  $n(z)$  as in

$$p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi}_{\text{int}}) = p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi}_{\text{true}}) \frac{p(z_{\text{spec}} | \vec{\phi}_{\text{int}})}{p(z_{\text{spec}} | \vec{\phi}_{\text{true}})}, \quad (3)$$

since  $p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi})$  is separable and the term  $p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma})$  changes independently of the  $p(z_{\text{spec}} | \vec{\phi})$  term. The information about  $\vec{\phi}_{\text{true}}$  is retained by the elements of the catalog: only  $(z_{\text{spec}}, z_{\text{phot}})$  pairs drawn from the true distribution will be represented in the catalog!

It appears as if this is the only fully self-consistent way to ensure that a draw from the posterior distribution  $p(z_{\text{spec}} | z_{\text{phot}}, \underline{\sigma}, \vec{\phi})$  does in fact follow that distribution. In the cases with fiducial posteriors  $p(z_{\text{spec}} | z_{\text{phot}}, \sigma, \vec{\phi})$ , the method of Sec. 3.1.1 was equivalent to this, but Sec. 3.3 requires a more sophisticated approach.

AIM thanks Elisabeth Krause for assistance with the `CosmoLike` code, Mohammadjavad Vakili for insightful input on statistics, Geoffrey Ryan for advice on debugging, and Boris Leistedt for helpful comments provided in the preparation of this paper.