



## LEVERAGING 3D-HST GRISM REDSHIFTS TO QUANTIFY PHOTOMETRIC REDSHIFT PERFORMANCE

RACHEL BEZANSON<sup>1</sup>, DAVID A. WAKE<sup>2,3</sup>, GABRIEL B. BRAMMER<sup>4</sup>, PIETER G. VAN DOKKUM<sup>5</sup>, MARIJN FRANX<sup>6</sup>, IVO LABBÉ<sup>6</sup>,  
 JOEL LEJA<sup>5</sup>, IVELINA G. MOMCHEVA<sup>5</sup>, ERICA J. NELSON<sup>5</sup>, RYAN F. QUADRI<sup>7</sup>, ROSALIND E. SKELTON<sup>8</sup>,  
 BENJAMIN J. WEINER<sup>1</sup>, AND KATHERINE E. WHITAKER<sup>9,10</sup>

<sup>1</sup> Steward Observatory, Department of Astronomy, University of Arizona, AZ 85721, USA

<sup>2</sup> Department of Physical Sciences, The Open University, Milton Keynes, MK7 6AA, UK

<sup>3</sup> Department of Astronomy, University of Wisconsin-Madison, Madison, WI 53706, USA

<sup>4</sup> Space Telescope Science Institute, Baltimore, MD 21218, USA

<sup>5</sup> Department of Astronomy, Yale University, 260 Whitney Avenue, New Haven, CT 06511, USA

<sup>6</sup> Leiden Observatory, Leiden University, Leiden, The Netherlands

<sup>7</sup> George P. and Cynthia W. Mitchell Institute for Fundamental Physics & Astronomy, Department of Physics & Astronomy,  
 Texas A&M University, College Station, TX 77843, USA

<sup>8</sup> South African Astronomical Observatory, P.O. Box 9, Observatory, Cape Town, 7935, South Africa

<sup>9</sup> Department of Astronomy, University of Massachusetts, Amherst, MA 01003, USA

Received 2015 October 22; accepted 2016 February 24; published 2016 May 2

## ABSTRACT

We present a study of photometric redshift accuracy in the 3D-HST photometric catalogs, using 3D-HST grism redshifts to quantify and dissect trends in redshift accuracy for galaxies brighter than  $JH_{\text{IR}} > 24$  with an unprecedented and representative high-redshift galaxy sample. We find an average scatter of  $0.0197 \pm 0.0003(1+z)$  in the Skelton et al. photometric redshifts. Photometric redshift accuracy decreases with magnitude and redshift, but does not vary monotonically with color or stellar mass. The  $1\sigma$  scatter lies between 0.01 and 0.03  $(1+z)$  for galaxies of all masses and colors below  $z < 2.5$  (for  $JH_{\text{IR}} < 24$ ), with the exception of a population of very red ( $U - V > 2$ ), dusty star-forming galaxies for which the scatter increases to  $\sim 0.1(1+z)$ . We find that photometric redshifts depend significantly on galaxy size; the largest galaxies at fixed magnitude have photo- $z$ s with up to  $\sim 30\%$  more scatter and  $\sim 5$  times the outlier rate. Although the overall photometric redshift accuracy for quiescent galaxies is better than that for star-forming galaxies, scatter depends more strongly on magnitude and redshift than on galaxy type. We verify these trends using the redshift distributions of close pairs and extend the analysis to fainter objects, where photometric redshift errors further increase to  $\sim 0.046(1+z)$  at  $H_{F160W} = 26$ . We demonstrate that photometric redshift accuracy is strongly filter dependent and quantify the contribution of multiple filter combinations. We evaluate the widths of redshift probability distribution functions and find that error estimates are underestimated by a factor of  $\sim 1.1$ – $1.6$ , but that uniformly broadening the distribution does not adequately account for fitting outliers. Finally, we suggest possible applications of these data in planning for current and future surveys and simulate photometric redshift performance in the Large Synoptic Survey Telescope, Dark Energy Survey (DES), and combined DES and Vista Hemisphere surveys.

*Key words:* galaxies: evolution – galaxies: high-redshift – galaxies: photometry – techniques: photometric

## 1. INTRODUCTION

Studies of the high-redshift universe rely increasingly upon *photometric redshifts* to identify and map the distribution of distant galaxies. These photometric redshifts are estimated from the overall spectral shapes as traced by catalogs of photometric data, as opposed to fitting one or more spectroscopic features. Photometric redshift surveys dramatically extend the possibilities of cosmological and galaxy evolutionary studies by vastly increasing the numbers and variety of galaxies beyond more observationally expensive spectroscopic galaxy surveys.

Because galaxy redshift is such a fundamental property, understanding the errors in photometric redshift estimates is crucial for interpreting empirical findings. For example, redshift uncertainties have been demonstrated to severely impact the measured evolution of the mass function (e.g., Chen et al. 2003; Marchesini et al. 2009; Muzzin et al. 2013). Photometric surveys can allow for studies of large-scale structure and galaxy clustering that are inaccessible to spectroscopic surveys, but the modeling of results depends strongly on understanding the redshift uncertainties (e.g., Chen

et al. 2003; Quadri et al. 2008; Wake et al. 2011; McCracken et al. 2015; Sołtan & Chodorowski 2015). In order to fully model the effects of photometric redshifts we must quantify their accuracy, which itself can depend on redshift and galaxy properties.

Traditionally, photometric redshift accuracy is tested by direct comparison between measured redshifts and *true* redshifts for a subset of a catalog with follow-up spectroscopy (e.g., Dahlen et al. 2013; Skelton et al. 2014). Alternatively, several groups have identified novel methods of testing photometric redshift accuracy using the clustering properties of galaxies (e.g., Newman 2008; Benjamin et al. 2010; Quadri & Williams 2010; Aragon-Calvo et al. 2015). Finally, a number of studies of photometric redshift accuracy have been conducted based on simulated mock galaxy catalogs (e.g., Ascaso et al. 2015). The first method is the most direct, but is typically biased toward very specific samples and the brightest galaxies for which spectroscopic redshifts are feasible: primarily at  $z < 1$  and for star-forming galaxies with bright emission lines. The second class of methods have different possible implementations, but in general these require large data sets, can lack sensitivity to certain types of systematic

<sup>10</sup> Hubble Fellow.

redshift errors or to catastrophic failures, and the results may be difficult to interpret. Although mock catalogs are an attractive alternative and require no additional data, they are fundamentally limited by their ability to match the empirical diversity of an evolving galaxy population.

Several methods of fitting photometric redshifts and many software packages and libraries are currently available and used by the community. Given the same data, each method will produce subtly different results (e.g., Hogg et al. 1998; Hildebrandt et al. 2008, 2010; Abdalla et al. 2011). Recently, Dahlen et al. (2013) published an extensive study evaluating the accuracy of redshifts produced by various photometric codes, focusing on the direct comparison of objects with spectroscopic redshifts in the Cosmic Assembly Near-infrared Deep Extragalactic Legacy Survey (CANDELS) fields, including a sample with deeper *Hubble Space Telescope* (*HST*) grism spectroscopic redshifts to extend the analysis to high redshift. Although the study investigated some trends in photometric redshift accuracy with galaxy properties, it is fundamentally limited to the availability of spectroscopic redshifts.

The 3D-HST survey (Brammer et al. 2012; Skelton et al. 2014, PI: P. van Dokkum) provides a unique opportunity to directly test the photometric redshift accuracy in the five CANDELS (Grogin et al. 2011; Koekemoer et al. 2011)/3D-HST fields. The data from this *HST* Legacy program combined with those from the A Grism H-Alpha SpecTrosopic (AGHAST) survey (PI: B. Weiner) include low-resolution grism spectroscopy across  $\sim 70\%$  of the CANDELS/3D-HST imaging footprint. This uniform spectroscopic coverage allows for unprecedented grism spectroscopic estimates of the true redshifts for thousands of galaxies beyond  $z > 1$ . Using grism redshifts, we can quantify the redshift accuracy of photometric catalogs in these fields for a sufficiently large and unbiased sample of high-redshift ( $z < 3$ ) galaxies. In this paper, we evaluate the photometric redshift accuracy in the *HST*/Wide Field Camera 3 (WFC3)-selected photometric catalogs produced by the 3D-HST collaboration (Skelton et al. 2014). Although we focus our investigation on photometric redshifts derived by the EAZY code (Brammer et al. 2008), we expect the conclusions to be similar for different algorithms given that Dahlen et al. (2013) found no strong differences among different methodologies and codes for a similar data set. Additionally, although that study recommended median combining photometric redshifts using a multitude of fitting techniques, the EAZY code was run by three different groups and consistently produced relatively low scatter and outlier fractions among the suite of redshift tests. In this work, we aim to quantify trends in the scatter between photometric and true redshifts as a function of galaxy properties as well as the occurrence rates of catastrophic failures.

Given the ultimate goal of quantifying photometric redshift performance in the 3D-HST catalogs, this paper is organized as follows. Section 2 briefly describes the 3D-HST data set. Section 3 quantifies the accuracy of photometric redshifts of the full detected sample and as a function of galaxy properties by comparison with spectroscopic and grism redshifts in addition to an analysis of close pairs. Section 4 discusses the relationship between photometric redshift accuracy and photometric bandpasses included in the redshift fitting. Section 5 addresses the use of the full photometric probability distribution function (PDF) of redshift as opposed to single-valued photometric redshifts. Section 6 extends the analysis of filter-

dependence to simulate photometric redshift performance in the Dark Energy Survey (DES), DES plus Vista Hemisphere Survey (VHS), and Large Synoptic Survey Telescope (LSST) surveys. Finally, we summarize the major results of the study in Section 7.

Throughout this paper we assume a concordance cosmology ( $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$ ,  $\Omega_M = 0.3$ , and  $\Omega_\Lambda = 0.7$ ) and quote all magnitudes in the AB system.

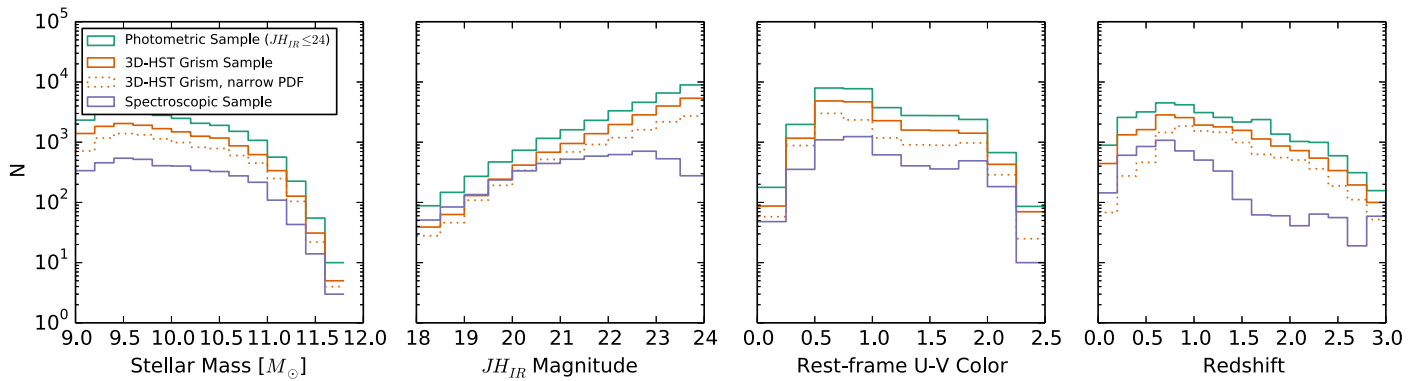
## 2. DATA

### 2.1. Sources of Data

The primary data in this paper are collected from the *HST*/WFC3-selected v4.1 photometric (Skelton et al. 2014) and grism catalogs (Momcheva et al. 2015) produced by the 3D-HST collaboration over  $\sim 900$  square arcminutes in five extragalactic fields: AEGIS, COSMOS, GOODS-North, GOODS-South, and UDS. The photometric catalogs include PSF-matched aperture photometry from a multitude of multi-wavelength ( $0.3\text{--}8.0 \mu\text{m}$ ) ground- and space-based images (Dickinson et al. 2003; Steidel et al. 2003; Capak et al. 2004; Giavalisco et al. 2004; Erben et al. 2005, 2009; Hildebrandt et al. 2006, 2009; Sanders et al. 2007; Taniguchi et al. 2007; Barmby et al. 2008; Furusawa et al. 2008; Wuyts et al. 2008; Nonino et al. 2009; Cardamone et al. 2010; Retzlaff et al. 2010; Grogin et al. 2011; Kajisawa et al. 2011; Koekemoer et al. 2011; Whitaker et al. 2011; Bielby et al. 2012; Brammer et al. 2012; Hsieh et al. 2012; McCracken et al. 2012; Ashby et al. 2013). Objects are detected from combined CANDELS/3D-HST *HST*/WFC3 images ( $J_{F125W}$ ,  $H_{F140W}$ , and  $H_{F160W}$ ). Photometric catalogs were produced using the MOPHONGO (Multiresolution Object PHotometry ON Galaxy Observations) code (I. Labbé et al., in preparation). The MOPHONGO code takes into account large differences in point-spread functions (PSFs) between *HST*/WFC3 detection images and ground-based and *Spitzer* data and corrects for source blending and confusion as described in Labbé et al. (2006), Wuyts et al. (2007), and Whitaker et al. (2011).

The 3D-HST Treasury Survey is primarily a 248 orbit grism spectroscopic survey, providing *HST*/WFC3 G141 near-infrared grism spectroscopy ( $\lambda = 1.1\text{--}1.7 \mu\text{m}$ ) in four of the five CANDELS/3D-HST (Grogin et al. 2011; Koekemoer et al. 2011) fields (AEGIS, COSMOS, GOODS-S, and UDS). Additional *HST*/WFC3 G141 grism spectroscopy in the GOODS-N field is included from the AGHAST survey (GO-11600, P.I.: B. Weiner). The combined data set covers a total of  $\sim 600$  square arcminutes with an average two-orbit depth. Objects selected from the 3D-HST photometric catalogs are matched in the grism data, extracted, and analyzed uniformly by the 3D-HST collaboration (Brammer et al. 2012; Momcheva et al. 2015). All extracted spectra are jointly fit along with the photometric data to provide grism redshifts for all objects brighter than  $JH_{\text{IR}} \leq 26$ , where  $JH_{\text{IR}}$  is based on flux in the combined  $F124W$ ,  $F140W$ , and  $F160W$  images. Grism spectra and redshift fits for all 23,564 galaxies brighter than  $JH_{\text{IR}} < 24$  are visually inspected to determine grism quality flags (use\_zgrism). Although redshift fits exist for fainter objects in the 3D-HST catalogs, we only include grism redshifts with these quality flags in this analysis. We adopt the term *grism redshift* ( $z_{\text{grism}}$ ) to describe these low-resolution spectroscopic redshifts to distinguish from traditional high-resolution *spectroscopic redshifts* ( $z_{\text{spec}}$ ). The uniform





**Figure 1.** Distribution of galaxies in the full photometric, grism, and spectroscopic redshift samples in stellar mass, apparent magnitude, rest-frame  $U - V$  color, and redshift. The grism redshift sample better reflects the distribution of properties of the photometric sample, particularly in faint and high-redshift galaxies. The dotted orange histogram indicates the sample of grism redshifts that provide estimates of  $z_{\text{true}}$  that are independent from photometric redshifts, as identified by decreased redshift uncertainty when the grism spectra are included in redshift fits.

spectroscopic coverage of the survey is crucial to the current investigation. For a complete description of the 3D-HST survey see Brammer et al. (2012), the photometric catalogs see Skelton et al. (2014), and the grism spectra see Momcheva et al. (2015).

The 3D-HST catalogs also include a vast collection of spectroscopic redshifts from ground-based spectroscopic surveys of these well-studied fields. In the AEGIS field, spectroscopic redshifts are matched with the DEEP2 DR4 survey (Cooper et al. 2012; Newman et al. 2013). In COSMOS, redshifts are collected from the zCOSMOS survey (Lilly et al. 2007), and a collection of MMT/Hectospec redshifts (M. Kriek et al. 2016, in preparation). GOODS-N redshifts are included from Kajisawa et al. (2010), who include data from a number of other surveys (Cohen et al. 2000; Cohen 2001; Dawson et al. 2001; Cowie et al. 2004; Wirth et al. 2004; Treu et al. 2005; Reddy et al. 2006; Barger et al. 2008; Yoshikawa et al. 2010). In GOODS-S, redshifts are collected from the FIREWORKS catalog (Wuyts et al. 2008). Finally, redshifts in the UDS are collected from the UDS Nottingham webpage, including data from Yamada et al. (2005), Simpson et al. (2006, 2012), Geach et al. (2007), van Breukelen et al. (2007), Ouchi et al. (2008), Smail et al. (2008), Ono et al. (2010), Akiyama et al. (2015), IMACS/Magellan redshifts (Papovich et al. 2010), an VLT X-shooter redshift from van de Sande et al. (2013), and Keck/DEIMOS redshifts (Bezanson et al. 2013, 2015).

Photometric redshifts from Skelton et al. (2014) catalogs are determined using the EAZY code (Brammer et al. 2008), which fits the spectral energy distribution (SED) of each galaxy with a library of galaxy templates and outputs the full PDF with redshift; see Skelton et al. (2014) for a complete description of this fitting. These fits utilize the default EAZY template set, which includes five PÉGASE (Fioc & Rocca-Volmerange 1997) stellar population synthesis models, a young, dusty template, and an old, red galaxy template that is described in Whitaker et al. (2011). Finally, EAZY includes an apparent magnitude prior,  $p(z|m_o)$ , which is the redshift distribution of galaxies with  $K$  band apparent magnitude,  $m_o$ . We note that these fits do not include active galactic nucleus templates and X-ray sources have not been excluded from this analysis. Skelton et al. (2014) found that removing sources with strong X-ray emission only marginally improves photometric redshift accuracy with respect to the spectroscopic redshift sample; only a small fraction of these bright sources are outliers in the  $z_{\text{spec}}$

versus  $z_{\text{phot}}$  comparison. We adopt  $z_{\text{peak}}$ , or the peak redshift marginalized over the PDF as a galaxy’s *photometric redshift* ( $z_{\text{phot}}$ ) in Sections 3 and 4 of this paper. In Section 5 we return to investigate the accuracy of the full photometric PDFs, assessing their overall widths. Grism redshifts that are obtained from joint fits to the photometry and *HST*-WFC3 slitless grism spectra from the 3D-HST survey. A full discussion of the redshift fitting can be found in Momcheva et al. (2015). In short, each two-dimensional grism spectrum is fit with a combination of EAZY continuum templates and a Dobos et al. (2012) emission line template, with a prior imposed by the photometric redshift PDF.

Derived properties are included from the version 4.1.5 3D-HST catalogs (Momcheva et al. 2015). These fits assume either the spectroscopic or grism redshift of each galaxy when available (Momcheva et al. 2015) or photometric redshifts (Skelton et al. 2014) in the full CANDELS footprint, derived as follows. Rest-frame colors are estimated for all galaxies following Brammer et al. (2011), also using the EAZY code. Stellar population parameters are calculated using the FAST code (Kriek et al. 2009) using Single Stellar Population (SSP) models from Bruzual & Charlot (2003) and assuming exponentially declining star formation histories, solar metallicity, and a Chabrier (2003) initial mass function. Galaxies with good photometry are identified by a use flag (use\_phot = 1 flag in the 3D-HST catalogs), which indicates that an object is not a star, is not near a bright star, has at least two exposures in  $F125W$  and  $F160W$  images, is detected in  $F160W$ , and has non-catastrophic redshift and stellar population fits. Finally, galaxy structural parameters are measured by single-component Sérsic fits to both  $F125W$  and  $F160W$  imaging (van der Wel et al. 2014). We adopt the definition of galaxy size as the effective radius along the semimajor axis.

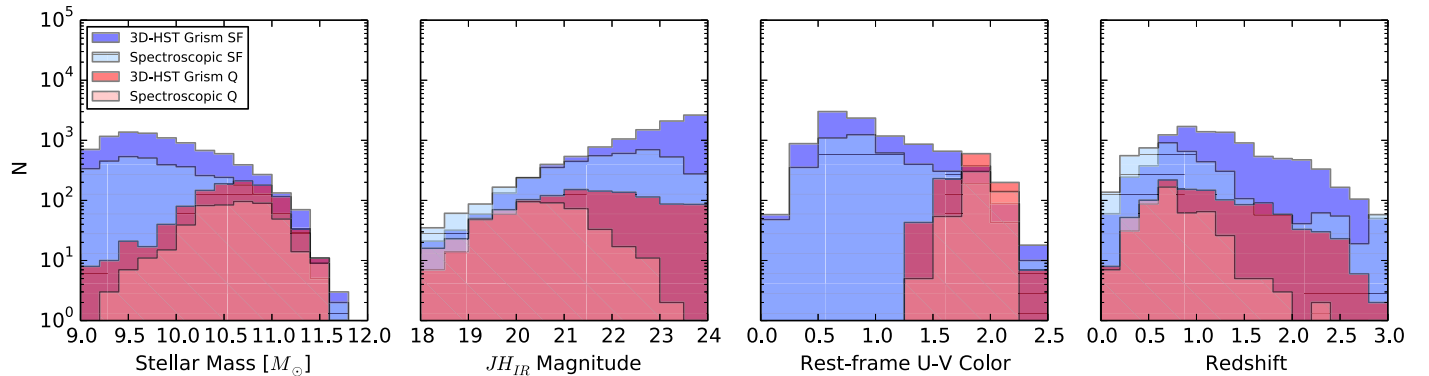
We adopt the maximum probability redshift,  $z_{\text{max\_gris}}$ , from the 3D-HST catalogs as the *grism redshift* ( $z_{\text{gris}}$ ) in this paper. A consequence of the inclusion of the photometric data in this fitting method is that the photometric and grism redshift estimates are not completely independent measurements. When investigating the scatter between the two correlated measurements, we only include galaxies for which the addition of the grism spectrum added significant information to the fit, as quantified by a tightened probability distribution, such that the 68% confidence interval for the  $z_{\text{gris}}$  is less than half of that for  $z_{\text{phot}}$  (discussed in more detail in Section 3.2).

**Table 1**  
Number of Galaxies in Each Sample

Total	AEGIS	COSMOS	GOODS-N	GOODS-S	UDS
Full Spectroscopic Sample					
4805	1094	420	1836	1280	175
Full Grism Sample					
18114	3244	3734	3565	4169	3402
Grism Sample with Narrowed PDFs <sup>a</sup>					
10550	2165	1720	2288	2250	2127
Quiescent Grism Sample with Narrowed PDFs <sup>a</sup>					
1080	192	180	215	274	219
Star-Forming Grism Sample with Narrowed PDFs <sup>a</sup>					
9470	1973	1540	2073	1976	1908

**Notes.** Total number of galaxies in each redshift sample. The narrowed grism redshift sample is over twice the size of the spectroscopic redshift sample and is much more representative of high redshifts ( $z \gtrsim 1$ ), faint magnitudes ( $JH_{\text{IR}} \gtrsim 21$ ), and for quiescent galaxies. The grism redshifts are evenly spread across the CANDELS/3D-HST fields.

<sup>a</sup> Narrowed PDFs refer to galaxies for which the 68% confidence interval for  $z_{\text{grism}}$  is less than or equal to half that of  $z_{\text{phot}}$  to minimize correlated measurement errors.



**Figure 2.** Distribution of star-forming (SF) and quiescent (Q) galaxies in the spectroscopic sample (hatched blue and red histograms) and the grism sample (solid blue and red histograms) with tightened redshift uncertainties in stellar mass, apparent magnitude, rest-frame  $U - V$  color, and redshift. In addition to its improved completeness at faint magnitudes ( $JH_{\text{IR}} \gtrsim 21$ ) and high redshifts ( $z \gtrsim 1$ ), it is clear that the sampling of the galaxy populations in those regimes is dramatically improved for the 3D-HST grism redshifts.

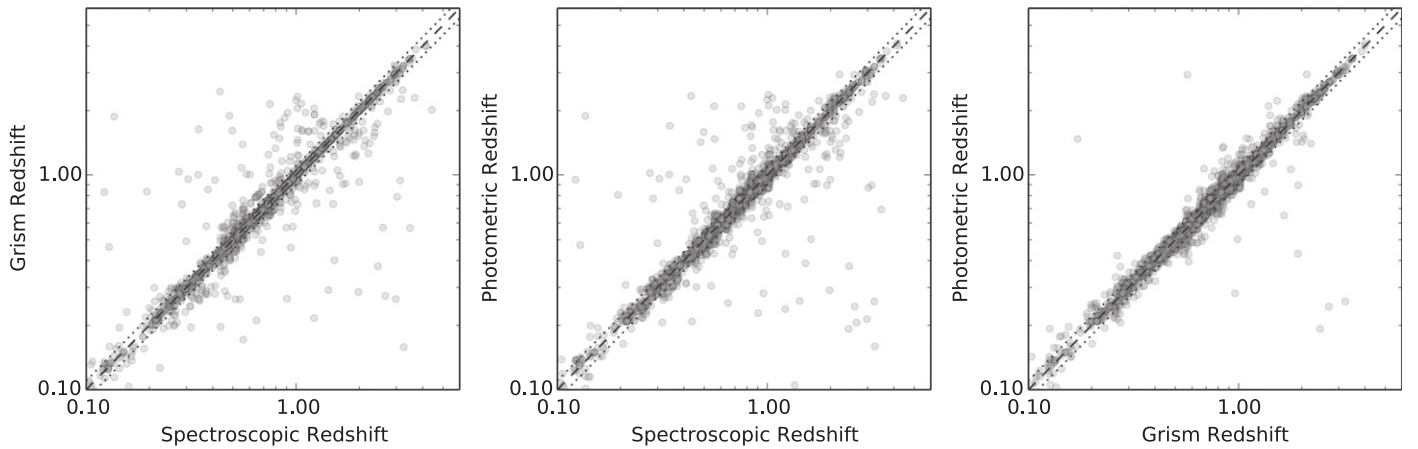
## 2.2. Properties of the Sample

Figure 1 indicates the distribution of  $JH_{\text{IR}} \leq 24$  galaxies in the 3D-HST catalogs with photometric redshifts (green), grism redshifts (orange), and spectroscopic redshifts (purple) as a function of stellar mass, apparent  $JH_{\text{IR}}$  magnitude, rest-frame  $U - V$  color, and redshift. The full grism sample is included as the orange histogram in Figure 1 and the effect of excluding possibly correlated redshift fits is indicated by the dotted orange histogram. The number of galaxies in each sample, both overall and in each field, is included in Table 1. Although this cut is roughly uniform across galaxy properties, this has the effect of preferentially excluding low redshift ( $z \lesssim 0.7$ ) galaxies, where the wavelength coverage of the G141 grism provides little spectral information. While this diminishes the utility of grism redshifts at low redshift, we emphasize that at these redshifts spectroscopic samples are much more representative of the overall population of galaxies. We further investigate the extent and consequences of possible correlations between photometric and grism redshifts in Section 3.2.

We highlight the bias of spectroscopic redshift surveys toward star-forming galaxies at the faint and high-redshift ends of the distributions. To demonstrate this, we use rest-frame  $U - V$  and  $V - J$  color criteria to distinguish between star-forming and quiescent galaxies in the 3D-HST catalogs, using

the thresholds defined by Whitaker et al. (2012). Solid histograms in Figure 2 show the number of star-forming galaxies (blue) and quiescent galaxies (red) with 3D-HST grism redshifts and narrowed redshift PDFs. The distribution of galaxies with spectroscopic redshifts is indicated by dotted lines and lighter histograms. Although the distributions of spectroscopic and photometric redshifts are similar in stellar mass and  $U - V$  color, the number of quiescent galaxies with spectroscopic redshifts dwindles dramatically fainter than  $JH_{\text{IR}} \gtrsim 21$  and at high redshift ( $z \gtrsim 1$ ). This is specifically the regime in which the grism redshifts are especially important.

The redshift distributions for spectroscopic, grism, and photometric redshift samples all differ and therefore the dependence of redshift accuracy on galaxy properties will also differ subtly; however, the grism sample will be more representative than a ground-based spectroscopic sample. It is clear from Figure 1 that the number of galaxies in the grism sample is nearly an order of magnitude larger than for the spectroscopic sample, but more importantly it more closely follows the distribution of the photometric catalog. Primarily, these redshifts include many more faint objects and galaxies at high ( $z > 1$ ) redshifts. Furthermore, the grism redshifts include vastly better sampling of the quiescent galaxy population improving by more than an order of magnitude on the number



**Figure 3.** Grism vs. Spectroscopic redshift, Photometric vs. Spectroscopic redshift, and Photometric vs. Grism redshift for all galaxies with spectroscopic redshifts in the 3D-HST catalogs. The scatter is lower between spectroscopic and grism redshifts ( $\sigma_{\text{NMAD}} = 0.0038$ ) than with photometric redshifts ( $\sigma_{\text{NMAD}} = 0.0158$  between  $z_{\text{spec}}$  and  $z_{\text{phot}}$  and  $\sigma_{\text{NMAD}} = 0.0156$  between  $z_{\text{grism}}$  and  $z_{\text{phot}}$ ); however, the outlier fraction is similar for grism and photometric redshifts.

of quiescent galaxies at faint magnitudes and high redshifts, although these numbers are still small.

### 3. PHOTOMETRIC REDSHIFT ACCURACY: QUANTIFYING SCATTER AND FAILURE RATES

The strongest test of photometric redshift performance given a fitting methodology can be obtained by comparing photometric redshifts to *true* redshifts for a subset of detected objects that reflects the parameter space spanned by the photometric catalogs themselves. Spectroscopic surveys provide excellent data sets with which to perform these tests, but are often quite biased either due to selection criteria or measurement failures. Due to its untargeted nature, redshifts determined from the 3D-HST grism spectra are not susceptible to these selection biases. In fact, the distribution of galaxies in the grism sample very closely follows that of the full photometric sample down to  $JH_{\text{IR}} \leq 24$  with a slight offset due to the smaller footprints (see Figures 1 and 2). Furthermore, we note that spectroscopic redshifts do not always represent the true redshift, either due to errors in spectroscopic analysis or misidentification of photometric counterparts.

In order to test the accuracy of the photometric redshifts in the 3D-HST/CANDELS fields, particularly for faint, high-redshift, and/or quiescent galaxies we benefit significantly by using grism redshifts, instead of those from higher resolution spectroscopy, as a proxy for the true redshifts of galaxies in the catalogs. In this section we demonstrate the feasibility of using the grism redshifts in this way and test the photometric redshift performance in the 3D-HST catalogs.

#### 3.1. Spectroscopic Sample

We begin by identifying a subset of 3278 galaxies in the 3D-HST catalogs with photometric, grism, and spectroscopic redshifts. Taking the spectroscopic redshift to be the true value, the scatter between redshift estimates is indicative of the errors in the photometric and grism redshifts. For the following tests, we compare all three redshifts for the full spectroscopic sample. Comparisons with the spectroscopic redshifts may yield the best estimate of redshift measurement errors, since these are more precise measurements of  $z_{\text{true}}$  and the grism and photometric redshift measurements may be correlated.

However the spectroscopic sample will always be smaller and more biased than the grism redshifts.

In Figure 3 we show the photometric, grism, spectroscopic redshift comparisons. Outlier thresholds of  $|\Delta z|/(1+z) > 0.1$  are indicated by dotted lines in each panel, where  $\Delta z$  is the difference between redshift measurements (either  $\Delta z = z_{\text{spec}} - z_{\text{phot}}$ ,  $\Delta z = z_{\text{spec}} - z_{\text{grism}}$ , or  $\Delta z = z_{\text{grism}} - z_{\text{phot}}$ ). Qualitatively, the left panel (grism versus spectroscopic redshift) exhibits less scatter than the center (photometric versus spectroscopic redshift) panel.

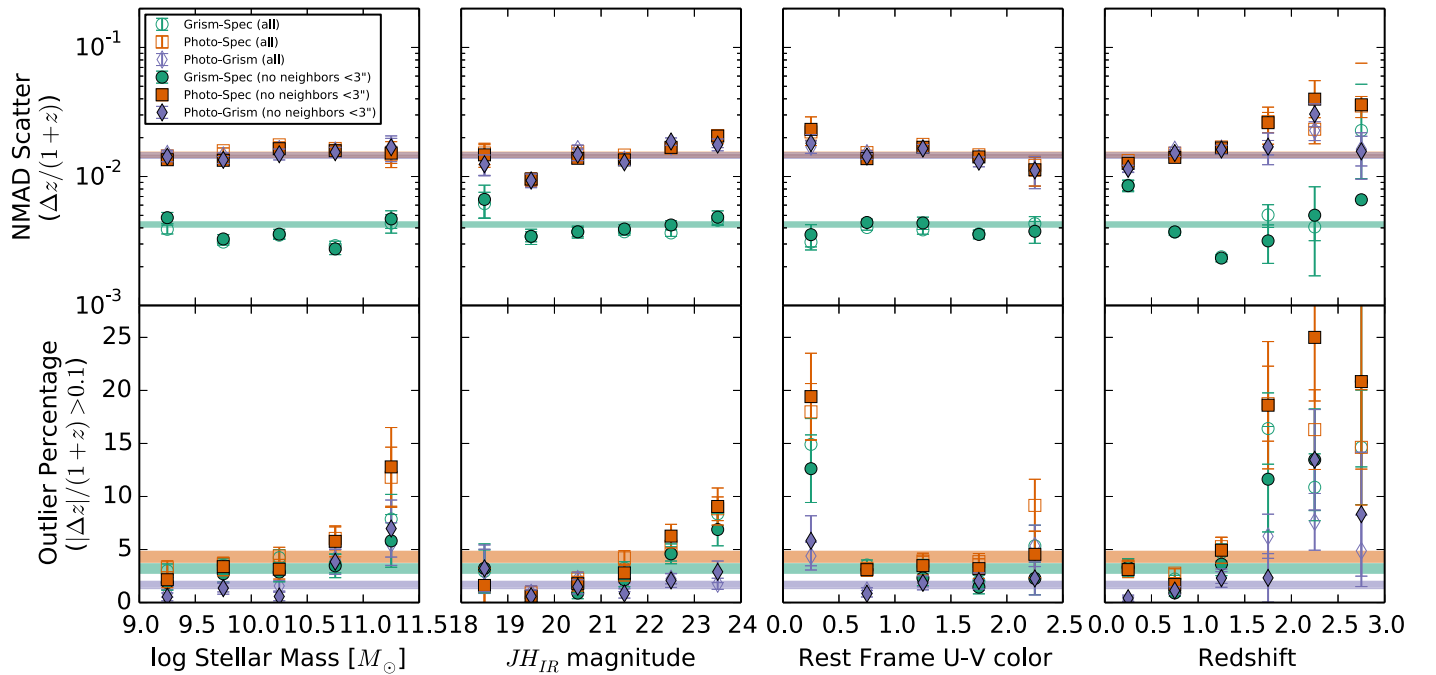
We quantify the scatter in  $\Delta z/(1+z)$  using the normalized median absolute deviation (NMAD) as

$$\sigma_{\text{NMAD}} = 1.48 \times \text{median}(|\Delta z|/(1+z)). \quad (1)$$

This measure of scatter is sensitive to the median deviations but less sensitive to catastrophic redshift failures than an rms scatter. The outlier fraction is defined as the fraction of galaxies with  $|\Delta z|/(1+z) > 0.1$ , although we find similar results with different definitions of this quantity. We emphasize that this definition of outliers does not include formal errors; we return to evaluating the redshift accuracy with respect to photometric redshift error estimates in Section 5. In this and subsequent sections, we only calculate scatter and outlier fraction for subsamples with more than 10 galaxies. Figure 4 shows the scatter and outlier fraction as a function of mass,  $JH_{\text{IR}}$  magnitude, rest-frame  $U - V$  color, and redshift for the spectroscopic sample. All comparisons are made for the same sample of galaxies: photometric versus spectroscopic redshifts in orange, grism versus spectroscopic in green, and photometric versus grism in purple. Errors in each measurement are estimated via bootstrap resampling of the full sample. The average value for each sample is indicated by the colored horizontal band in each panel and average scatter and outlier fractions are reported in Table 2.

One concern in interpreting accuracies derived from comparisons with spectroscopic redshifts is the possibility that published spectroscopic redshifts can also be erroneous. The spectroscopic redshift catalog contains only high-quality redshifts, as assessed by each independent study; however, there is still the possibility that the spectroscopic measurement was not assigned to the correct object in the 3D-HST catalogs.





**Figure 4.** Redshift accuracy for 3D-HST galaxies with spectroscopic, photometric, and grism redshifts. Each column includes NMAD scatter (top row) and outlier fraction (bottom row) for this sample as a function of stellar mass (left column),  $JH_{IR}$  magnitude (second column), rest-frame  $U - V$  color (third column), and redshift (fourth column). Comparison between spectroscopic and grism redshifts is included with green symbols, spectroscopic and photometric redshifts in orange, and grism and photometric redshifts in purple. Filled symbols include only galaxies without neighbors within  $3''$  to eliminate possible spec- $z$  misidentifications; this does not significantly decrease outlier fractions. Scatter between grism and spectroscopic redshifts is much lower than for photometric redshifts, but the outlier fraction is similar. Scatter between photometric redshifts and grism redshifts is extremely similar to scatter with spectroscopic redshifts, suggesting grism redshifts can also be used as a proxy for true redshift.

**Table 2**  
Scatter and Outlier Fraction in Each Sample

S1	S2	$\sigma_{\text{NMAD}}$	Outlier%
Full Spectroscopic Sample <sup>a</sup>			
Phot	Spec	$0.0159 \pm 0.0005$	$4.9\% \pm 0.4$
Gris	Spec	$0.0039 \pm 0.0001$	$4.2\% \pm 0.4$
Phot	Gris	$0.0154 \pm 0.0005$	$1.9\% \pm 0.3$
Spectroscopic Sample without possible mis-IDs <sup>b</sup>			
Phot	Spec	$0.0148 \pm 0.0006$	$4.3\% \pm 0.5$
Gris	Spec	$0.0043 \pm 0.0002$	$3.2\% \pm 0.4$
Phot	Gris	$0.0145 \pm 0.0007$	$1.7\% \pm 0.3$
Grism Sample with Narrowed PDFs <sup>c</sup>			
Phot	Gris	$0.0193 \pm 0.0003$	$3.5\% \pm 0.2$
Grism Sample with Narrowed PDFs (Star-Forming) <sup>c</sup>			
Phot	Gris	$0.0198 \pm 0.0003$	$3.6\% \pm 0.2$
Grism Sample with Narrowed PDFs (Quiescent) <sup>c</sup>			
Phot	Gris	$0.0161 \pm 0.0008$	$2.7\% \pm 0.5$

**Notes.** Average scatter and outlier fraction between photometric, grism, and spectroscopic redshifts in the 3D-HST survey.

<sup>a</sup> Sample selection:  $z_{\text{spec}} > 0$ ,  $\text{use\_phot} = 1$ ,  $\text{use\_zgrism} = 1$ .

<sup>b</sup> Sample selection:  $z_{\text{spec}} > 0$ ,  $\text{use\_phot} = 1$ ,  $\text{use\_zgrism} = 1$ , no neighbors within  $3''$  for which  $z_{\text{spec}}$  falls within 95% confidence interval for  $z_{\text{phot}}$ .

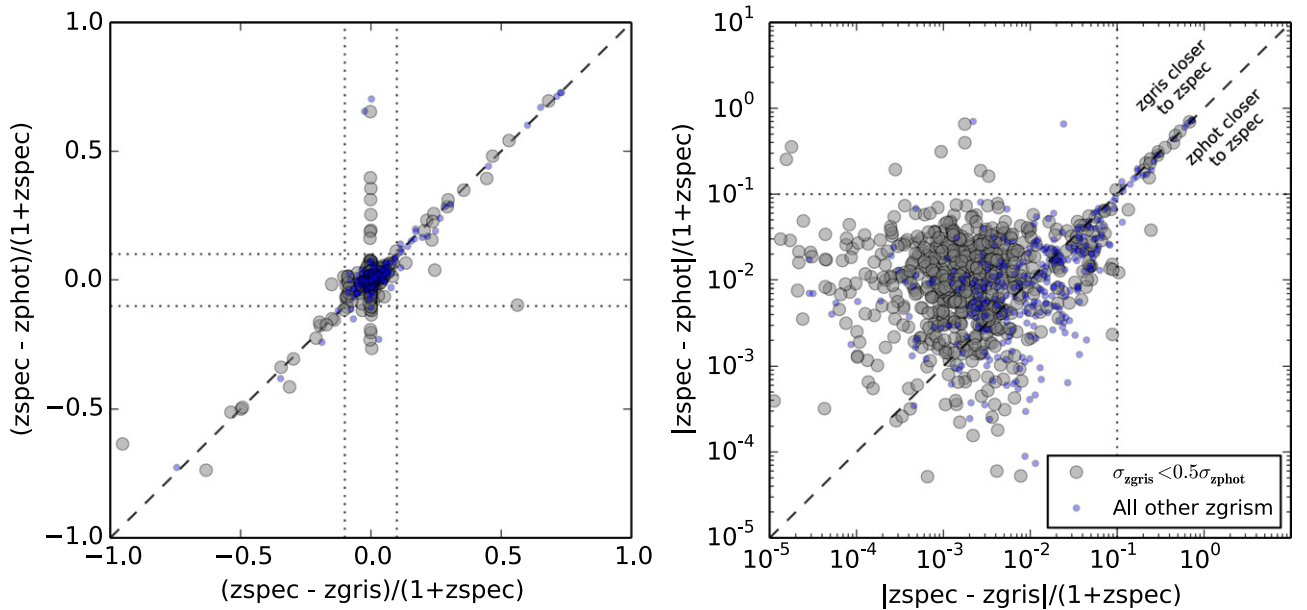
<sup>c</sup> Sample selection:  $\text{use\_phot} = 1$ ,  $\text{use\_zgrism} = 1$ , 68% confidence interval for  $z_{\text{grism}}$  less than or equal to half that of  $z_{\text{phot}}$ .

Spectroscopic counterparts in the 3D-HST catalogs were matched within a radius of  $0''.5$  (Skelton et al. 2014). Although this is a conservative matching aperture, misidentification of

photometric counterparts due to faulty astrometry or close neighbors, could falsely boost the measured rate of catastrophic failures in photometric redshift estimations. We can minimize this possibility by only including spectroscopic redshifts for galaxies with a unique counterpart in the 3D-HST photometric catalogs, removing galaxies from the sample for which there was at least one neighboring galaxy within  $3''$  for which the spectroscopic redshift falls inside of the 95% confidence interval of the photometric  $P(z)$ . The scatter and outlier fractions for this sample are included as filled symbols in Figure 4. This aggressive cut decreases the sample to 1807 galaxies. However, the effect on scatter and outlier fractions is extremely subtle. Therefore, the catastrophic redshift failures cannot be explained simply by incorrect comparisons, but note that there additional errors in spectroscopic redshift identification could also contribute these outliers.

A number of overall trends appear in each column. Scatter between grism and spectroscopic redshifts is much lower than that for photometric redshifts, but the outlier fraction is comparable. The outlier fractions are lower between grism and photometric redshifts, suggesting that the two measurements are correlated. The NMAD scatter between photometric redshifts and grism or spectroscopic redshifts is strikingly similar, both on average and as a function of galaxy properties. In most cases the measurements completely overlap. This suggests that if grism redshifts are used to evaluate the accuracy of photometric redshifts,  $\sigma_{\text{NMAD}}$  will be a robust indication of the scatter about  $z_{\text{true}}$ .

The outlier fraction is  $\sim 2$  times lower for the grism redshifts compared to the spectroscopic redshifts, which suggests the existence of correlated errors if the grism catastrophic redshift failures are a subset of photometric failures. We investigate



**Figure 5.** Correlations between photometric and grism redshift errors with respect to spectroscopic redshifts in linear scale (left panel) and logarithmic scale (right panel). Large gray symbols indicate galaxies for which the 68% confidence interval for  $z_{\text{grism}}$  is  $\leq 0.5$  that of the photometric redshift, small blue points mark galaxies with untightened PDFs. The vertical axes show residuals in photometric redshifts and horizontal axes show the residuals in grism redshifts about the known spectroscopic redshifts. Dotted lines indicate the  $|\Delta z|/(1+z) > 0.1$  outlier threshold. Diagonal trends (dashed lines) indicate correlated errors between photometric and grism redshifts for a small subset of the complete sample. The vertical trend in the left panel highlights a subsample of galaxies for which the grism redshifts catch the true (spectroscopic) redshift, while the photometric redshifts exhibit a fair amount of scatter. Minimal correlated residuals between photometric and grism redshifts and the spectroscopic redshifts suggest that grism redshifts provide an independent measurement of an object’s true redshift.

how much of this is driven by cases where the spectroscopic redshifts are not accurately identifying  $z_{\text{true}}$ . We visually inspect the spectral energy distributions and images of the 58 outliers ( $|\Delta z_{\text{spec}} - z_{\text{grism}}|/(1+z_{\text{spec}}) > 0.1$ ), 6 of which are  $|\Delta z_{\text{phot}} - z_{\text{grism}}| > 0.1$  outliers. First, we find that  $\sim 36\%$  (21) of the galaxies are below  $z_{\text{grism}} = 0.7$ , where the G141 grism provides little additional information due to a lack of spectral features. Additionally, many of these outliers (63%, 36) are in the GOODS-N MODS compilation (Kajisawa et al. 2010), which does not have quality flags. Furthermore, the grism spectra caught emission lines for 17 (29%) of these galaxies. Finally, the grism spectra are extracted using the photometric positions and therefore, in the absence of blending in the *HST* imaging, they are not susceptible to misidentification. We conclude that for a significant fraction of catastrophic redshift outliers, the grism provides estimates of true redshifts of the photometric objects in the catalog that are as good or better than the spectroscopic redshifts and this difference could account for the difference in outlier fractions.

Photometric redshift errors increase with both  $JH_{\text{IR}}$  magnitude and redshift in a comparison with either spectroscopic or grism redshifts for the spectroscopic sample, as found by Dahlen et al. (2013). We find very little correlation between redshift accuracy and stellar mass, and a non-monotonic but clear trend of decreasing scatter for the reddest colors. Although it is tempting to interpret trends in redshift accuracy shown in Figure 4, we caution that these are based on a heterogeneous (and biased) spectroscopic sample. In particular, the outlier fraction increases dramatically with redshift at  $z > 1.5$ . However, this is also where the size of the spectroscopic sample dwindles. One takeaway is that for this sample of galaxies for which spectroscopic redshifts are obtainable (and perhaps easy), the grism redshifts are excellent (NMAD scatter is low), but the outlier fraction is similar to that of the

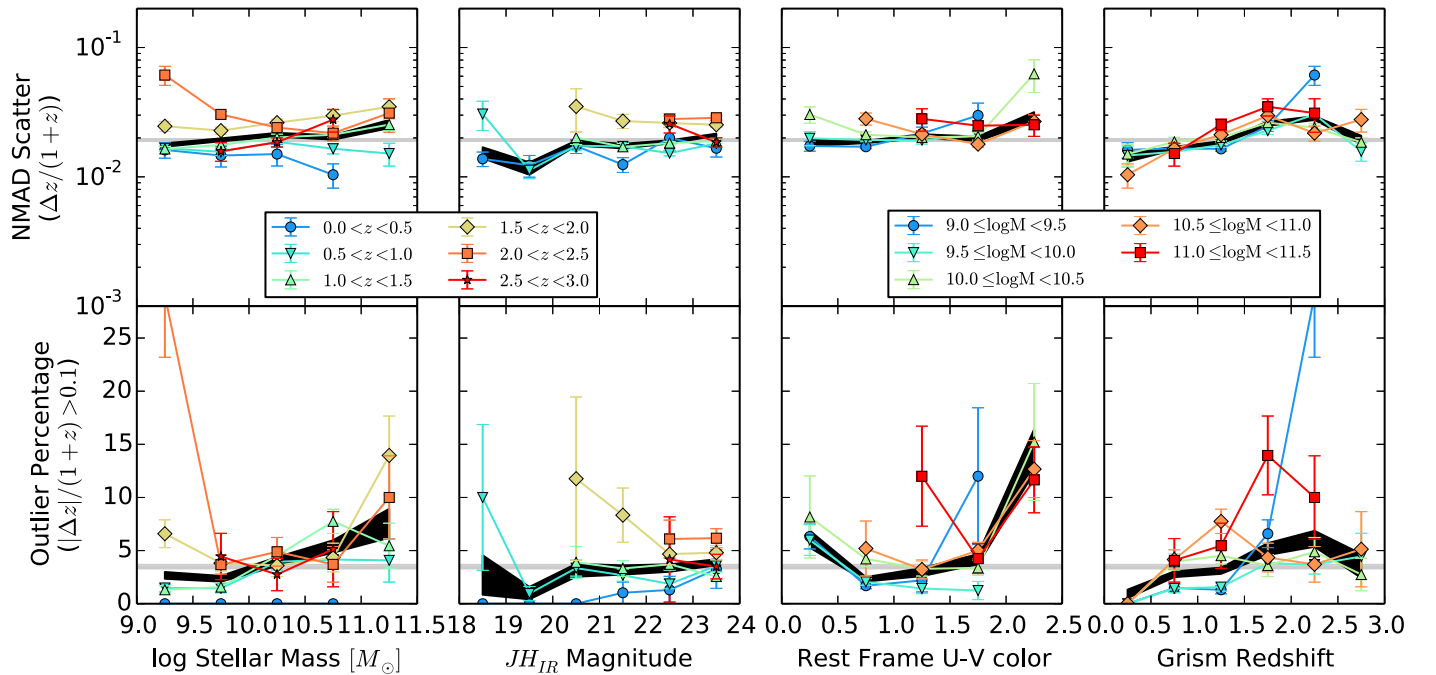
photometric redshifts. The spectroscopic subsample is too small to disentangle trends in both redshift and mass; for this we must utilize grism redshifts for a larger sample.

### 3.2. How Correlated Are Grism and Photometric Redshifts?

The 3D-HST grism redshift fits are made using a joint fit to the photometric catalogs and grism spectra; the resulting redshift estimates may be correlated with purely photometric redshifts. In this section we assess the magnitude of this correlation and therefore the utility of grism redshifts as an independent estimate of *true redshift*. For this test, we include the full sample with spectroscopic and grism redshifts, investigating the residuals between the spectroscopic, or true, redshift of a galaxy and its photometric and grism redshift.

Figure 5 shows the residuals with respect to spectroscopic redshift in photometric versus grism redshift. In each panel, the large symbols indicate galaxies for which the grism redshift  $P(z)$  is tightened with respect to that of the photometric redshift (68% confidence interval of  $z_{\text{grism}}$  is narrower than that of  $z_{\text{phot}}$  by a factor of 0.5), the small symbols show the remainder of the sample. This criterion does not severely impact the demographics of galaxies with grism redshifts (dashed orange lines in Figure 1), but does minimize the effect of correlated residuals. In this section we aim to quantify the effects of correlated errors on this sample; in subsequent sections we will use only grism redshifts to test photometric redshifts. Dotted lines indicate our adopted outlier threshold ( $|\Delta z|/(1+z) = 0.1$ ).

Figure 5(a) shows the residuals with linear scaling. An interesting feature of the left panel is the vertical trend, indicating a subsample of galaxies for which the grism identifies the spectroscopic redshift, but the photometric redshifts exhibit higher residuals. Only a small fraction of galaxies lie along the diagonal trend, on which photometric and



**Figure 6.** NMAD scatter (top row) and outlier percentage (bottom row) for all grism redshifts with narrower  $P(z)$ s than for the photo- $z$  (by a factor of  $\leq 0.5$ ) split by stellar mass,  $JH_{IR}$  magnitude,  $U - V$  color, and redshift are indicated by filled black bands. Mean values are indicated by gray bands. Samples are further split by either redshift (left two panels) or stellar mass (right two panels). Photometric redshift accuracy depends primarily on magnitude and redshift, with non-monotonic variations as a function of galaxy mass or color.

grism redshifts exhibit strongly correlated residuals, particularly for the subset of galaxies with “tightened” PDFs. Only  $\sim 3\%$  of galaxies are outliers in both photometric and grism redshifts for the total sample; however,  $\sim 67\%$  of photometric outliers are also outliers in grism redshifts. For galaxies with tightened PDFs, this correlated outlier rate is lower at  $\sim 2\%$ . For this sample,  $\sim 58\%$  photometric outliers are also grism outliers. The true rate of correlated redshift failure could be even lower. From visual inspections of images, SEDs, and grism spectra, we find that 33% of the galaxies with tightened PDFs and correlated residuals have possible neighbors that could contribute to  $z_{\text{spec}}$  misidentification and 55% of the grism spectra include an identified emission line, suggesting that the grism redshift is the true redshift.

The Figure 5(b) shows the absolute value of the photometric and grism redshift residuals. In logarithmic scaling, galaxies preferentially lie above the diagonal line, indicating that grism redshifts have smaller residuals than photometric redshifts. The scatter is higher for the photometric residuals ( $\sim 0.037$ ) than grism redshifts ( $\sim 0.0145$ ), when correlated residuals ( $> 0.05$  in both) are excluded. The cut in grism redshift uncertainty eliminates a large fraction of high residual and correlated objects in this projection. Only a small fraction (2.7% of tightened sample) of all galaxies lie on the diagonal trend of correlated errors. Therefore, scatter and outlier fractions between photometric and grism redshifts will be dominated by the independent accuracy of each redshift estimate, but will not be artificially reduced by correlated errors.

### 3.3. Beyond Spec-zs: Trends in Photometric Redshift Accuracy with Mass, Magnitude, Color, Redshift, and Size

#### 3.3.1. Testing Photometric Redshifts with Grism Redshifts

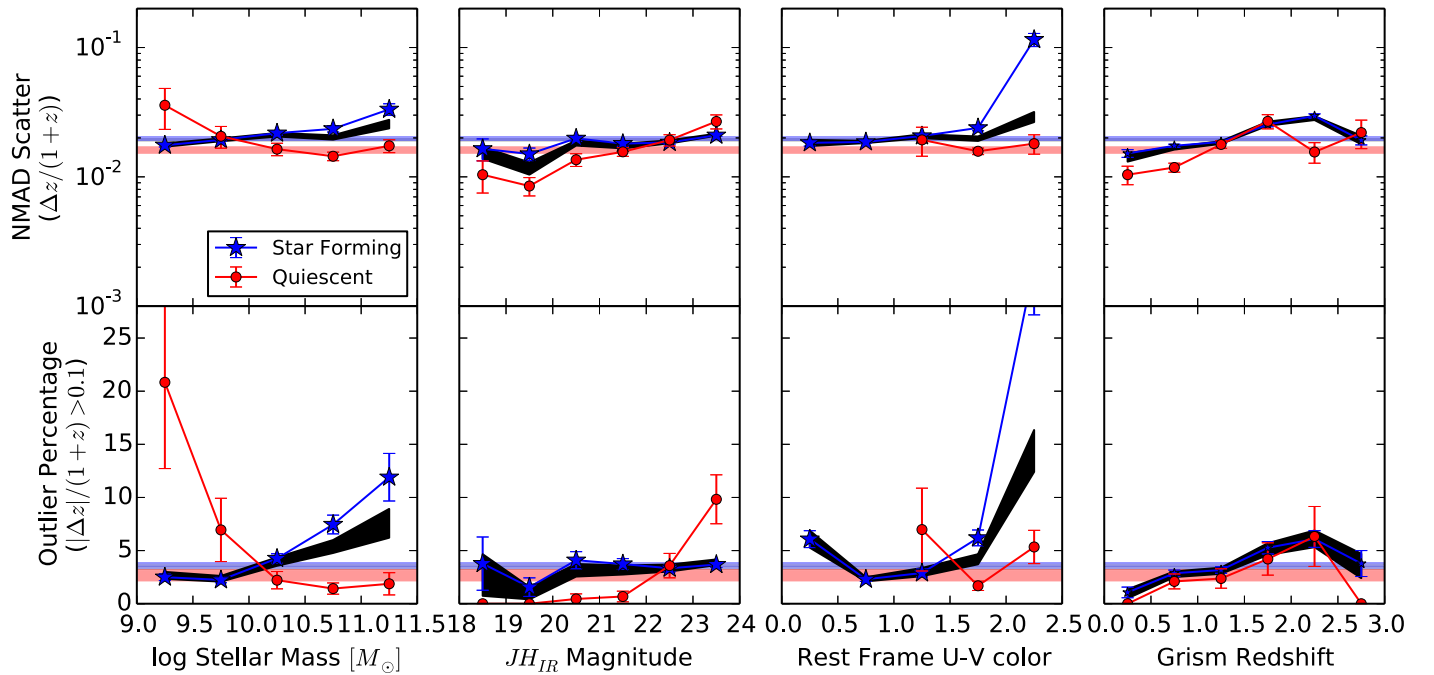
We have demonstrated that 3D-HST grism redshifts can be used to provide a measurement of  $z_{\text{true}}$  and assess photometric

redshift quality, improving upon the severe biases inherent with using spectroscopic redshifts. In this section we utilize the full sample of grism redshifts to investigate the variation in photometric redshift performance. For this test, we include all galaxies with good photometry and grism redshifts ( $\text{use\_phot} = 1$ ,  $\text{use\_zgrism} = 1$ ) and narrowed PDFs (as defined in the previous section). The uniformity and size of the sample of galaxies with grism redshifts, as opposed to a spectroscopic sample (see Figures 1 and 2), allows us to dissect trends photometric redshift accuracy in mass, apparent magnitude, galaxy color, and redshift.

Figure 6 demonstrates trends in NMAD scatter (top row) and outlier fraction (bottom row) as a function of stellar mass and magnitude in the  $JH_{IR}$  imaging (first and second columns: split into redshift ranges) and redshift and  $U - V$  color (third and fourth columns: split by stellar mass). The average scatter and fraction are indicated by the gray band in each panel. On average, the scatter between  $z_{\text{phot}}$  and  $z_{\text{grism}}$  is slightly higher than that of the spectroscopic sample ( $\sigma_{\text{NMAD}} = 0.0198 \pm 0.0003$  versus  $\sigma_{\text{NMAD}} = 0.0148 \pm 0.0006$  for the  $z_{\text{spec}}$  comparison). There are certain mass and redshift ranges for which the outlier fraction increases dramatically, but part of this seems to be driven by uncertain grism redshifts or small subsample size.

The NMAD scatter does not depend strongly on stellar mass or UV color, with the minor exception of galaxy populations such as extremely red high-redshift galaxies that are likely to be ill fit (lower left panel). We emphasize that many of the derived galaxy properties including stellar masses and colors are derived based on the grism redshifts. This will likely lead to correlated trends such as poorly modeled galaxies with  $2 < z < 2.5$  that are assigned very low stellar masses. However, by using this unique data set, it is apparent that the fraction of photometric redshifts that will catastrophically fail in estimating the true redshift of a galaxy depends strongly on the properties of and distance to the galaxy. For example, the





**Figure 7.** NMAD scatter (top row) and outlier percentage (bottom row) compared to grism redshifts (with narrower  $P(z)$ s than for the photo- $z$  by a factor of  $\leq 0.5$ ) split by stellar mass,  $JH_{\text{IR}}$  magnitude,  $U - V$  color, and redshift are indicated by filled black band. The sample is split into star-forming (blue stars) and quiescent (red circles) galaxies based on their  $UV$  and  $VJ$  colors. Mean values are indicated by blue and red bands. Quiescent galaxies have more accurate photometric redshifts than star-forming galaxies; however, this accuracy is strongly dependent on magnitude and redshift.

outlier fraction for low-mass galaxies is extremely low ( $\lesssim 5\%$ ) at low redshift ( $z < 1.5$ ) and for those with blue colors, but increases by a factor of  $\sim 2-3$  at higher redshifts. The outlier fraction of massive galaxies ( $\log(M_*/M_\odot) > 10.5$ ) is a factor of  $\sim 2$  higher than average at all but the highest and lowest redshift bins.

We note that increased scatter or outlier fractions in this sample could indicate regimes in which either photometric or grism redshifts, or both, are less accurate. For example low-mass galaxies at  $z \sim 2$  exhibit large outlier fractions, even though the  $\sigma_{\text{NMAD}}$  is less dependent on these properties.

Another key question is how photometric redshift performance depends on galaxy type. Directly testing this is uniquely possible with the 3D-HST data set. Figure 7 includes the same trends in scatter and outlier fraction, but now indicates the trends for  $U - V$  and  $V - J$  identified star-forming (blue stars) and quiescent (red circles) galaxies. The errors on rest-frame colors are negligible ( $< 0.1$ ) with respect to this color separation (see e.g., Whitaker et al. 2011), although catastrophic redshift failures could lead to systematic errors in rest-frame color estimates and contaminate the samples, we expect this effect to be minimal due to low outlier rates. Overall, photometric redshifts are more accurate for quiescent galaxies than star-forming galaxies in scatter and outlier fractions, as indicated by the red and blue bands. This can be readily understood because quiescent galaxies have stronger Balmer/4000 Å breaks, which are easily identified in broad or medium-band photometry. Above  $z \sim 2.5$ , the Lyman break for star-forming galaxies begins to fall into the optical photometric bands and improves photometric redshift accuracies (see also e.g., Whitaker et al. 2011).

Additionally, trends in these panels are extremely helpful in interpreting Figure 6. For example, although the scatter of the full sample does not depend on stellar mass, scatter increases to

$\sim 0.03(1+z)$  for star-forming galaxies above  $M_* > 10^{11}M_\odot$ . Similarly, the increasing outlier fraction (up to  $\sim 10\%$ ) is due to star-forming galaxies alone; photometric redshift accuracy does not appear to depend on stellar mass for quiescent galaxies. On the other hand, photometric redshift accuracy decreases more dramatically with magnitude for quiescent galaxies (rising from  $\sigma_{\text{NMAD}} \sim 0.01$  at  $JH_{\text{IR}} \sim 18$  to  $\sim 0.03$  at the faint end versus star-forming galaxies, which exhibit  $\sigma_{\text{NMAD}} \sim 0.02$  at all magnitudes).

Perhaps the most striking trend is with rest-frame color, where photometric redshift error and outlier fractions dramatically increase to  $\sigma_{\text{NMAD}} = 0.11\%$  and  $\sim 37\%$  outliers at the reddest  $U - V$  colors. This trend was also apparent in Figure 6, but it is now apparent that only star-forming galaxies are contributing to the increase in scatter and outlier fraction. These galaxies must be extremely dusty to explain their red colors and they appear to have highly degenerate redshifts with the current template set (G. B. Brammer et al., in preparation), despite the inclusion of the old, dusty template. It is noteworthy that this trend does not exist in the spectroscopic sample, highlighting the importance of the 3D-HST grism redshifts in fully characterizing photometric redshift performance.

These red, dusty star-forming galaxies are an increasingly prevalent population at high redshift (e.g., Marchesini et al. 2010, 2014; Muzzin et al. 2013). Estimating the photometric redshifts for galaxies that are both red in  $U - V$  and  $V - J$  colors is helped by including an appropriate dusty starburst template (Marchesini et al. 2010), but in general estimating their photometric redshifts becomes more difficult as the dust degrades the prominence of the break. Not accounting for this growing population of galaxies can systematically place them at the wrong photometric redshifts (Marchesini et al. 2010), significantly influence the observed evolution of the stellar mass function for star-forming galaxies (Muzzin

et al. 2013), and underestimate star formation rates (e.g., Fumagalli et al. 2014). However, the Skelton et al. (2014) photometric redshift fits already include a dusty and old template in the EAZY template set. In this case, the scatter and outlier fraction of the reddest galaxies points to a subset of extremely red star-forming galaxies for which photometric redshifts still consistently fail.

Another observational property that may contribute to the accuracy of photometric redshifts is the apparent size and surface brightness of each individual galaxy. We now investigate the scatter and outlier fraction as a function of size at fixed magnitude and redshift as follows. First, we interpolate between effective radii measured in *F125W* and *F160W* to best match the central wavelength of the grism spectra and direct imaging. In addition to cuts for good photometry and tightened grism redshift confidence intervals, we include only galaxies with reliable size estimates based on the van der Wel et al. (2014) flags (flag  $\leq 1$  in both *F125W* and *F160W* fits). Figure 8 shows the distribution of sizes (top row), trends in NMAD scatter (middle row), and photometric redshift outlier fractions (bottom row). The first column indicates trends in photometric redshift accuracy with magnitude for larger and smaller galaxies. Except for the brightest (and overall largest) and faintest (and smallest) galaxies, there is a clear trend that the most extended galaxies at fixed magnitude exhibit higher scatter (up to  $\sim 0.005$ ) and larger outlier fractions (by up to a factor of five). This effect is likely due to a combination of lower signal-to-noise ratio (S/N) for these galaxies and the subtle effects of point-spread function (PSF) matching and aperture photometry. We note that S/N due to surface brightness will also effect the quality of spectroscopic or grism redshift estimates. However, the 3D-HST grism redshifts are fit in two dimensions, using the full galaxy morphology, and will not depend on aperture effects like the photometric redshifts.

The second column of Figure 8 demonstrates the trends in photometric redshift precision with size at fixed redshift. Below  $z = 2$ , the photometric redshifts of the largest galaxies exhibit more scatter with the grism redshifts than average or compact galaxies and exhibit larger outlier fractions. However, above  $z > 2$  the photometric redshift accuracy is statistically consistent, at least partially due to the fact that nearly all galaxies lie significantly within the photometric apertures. At fixed redshift, the size effect is less clear than at fixed magnitude, suggesting that surface brightness plays a dominant role in photometric redshift performance.

Star-forming and quiescent galaxies exhibit distinct structural properties in the local universe (e.g., Shen et al. 2003) and at high redshift (e.g., van der Wel et al. 2014), such that star-forming galaxies are more extended at fixed mass. Therefore we evaluate the dependence of photometric redshift accuracy on size for the two populations separately in the third and fourth columns of Figure 8. Because star-forming galaxies dominate the number of galaxies at all epochs, the dependence of photometric redshift accuracy on redshift looks extremely similar for this population as for the full sample: larger galaxies have less accurate photometric redshifts than smaller galaxies below  $z = 2$ . Conversely, quiescent galaxies do not show strong trends with size, although between  $1.5 < z < 2$  there is a statistically significant increased fraction of outliers. This lack of sensitivity to galaxy size can be attributed to the compactness of quiescent galaxies, which primarily lie well within the photometric apertures.

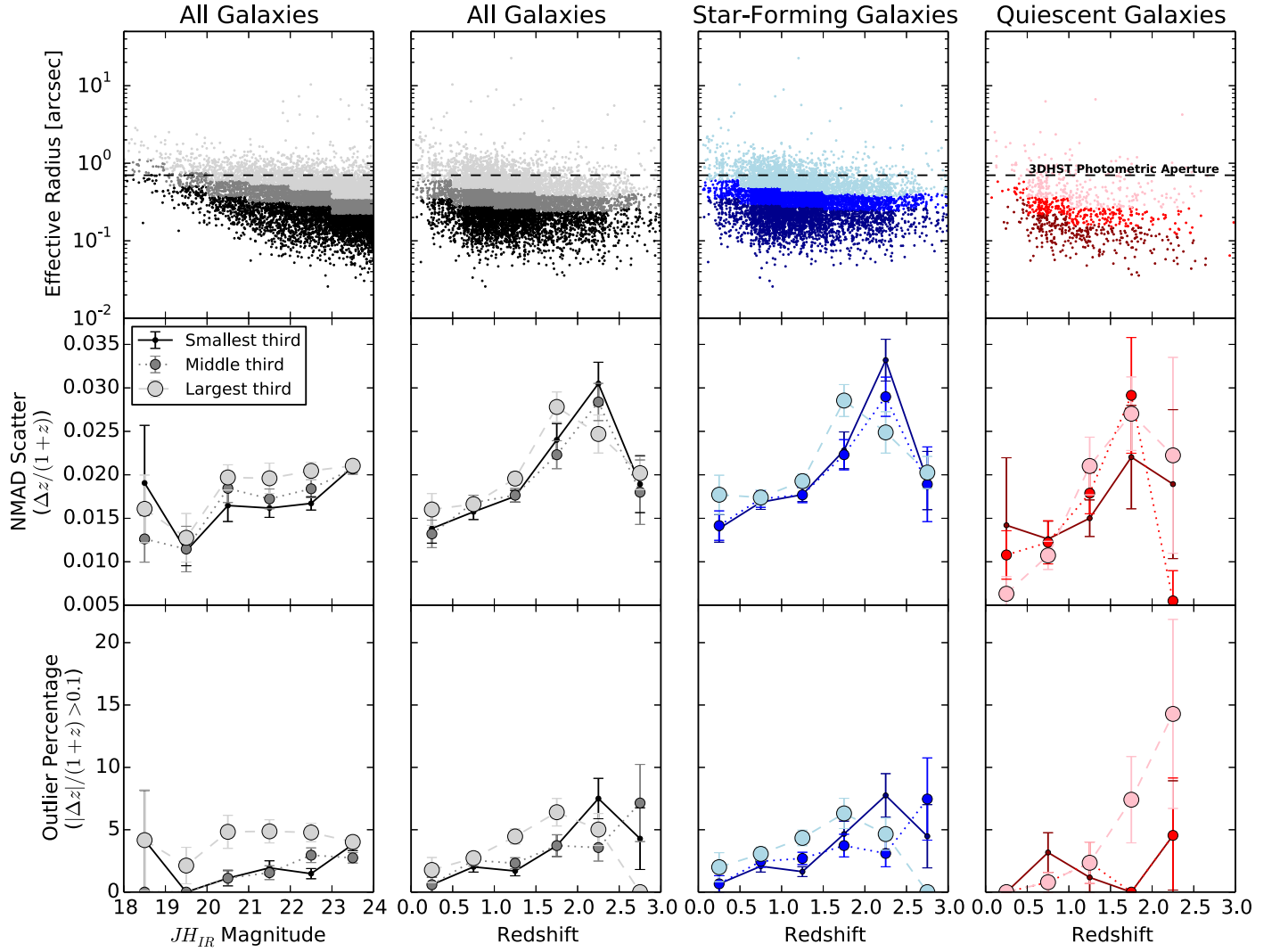
### 3.3.2. Photometric Redshift Accuracy from Close Pairs

The analysis in the previous subsection depended on the use of 3D-HST grism redshifts to estimate  $z_{\text{true}}$  for galaxies in the photometric catalogs. We perform an independent test of the photometric redshift accuracy by following the close pairs analysis described by Quadri & Williams (2010). Due to the clustered nature of galaxies throughout space, galaxies which appear very near to one another projected on the sky are likely to be physically associated. If both galaxies are at the same true redshift, then differences in measured redshift will be due to the measurement and template/fitting errors. Although true physical pairs cannot be individually identified, Quadri & Williams (2010) defined a statistical method to utilize the distribution of these redshift differences, subtract out the contribution of chance alignment to the sample of close pairs in a survey, and calculate the photometric redshift errors. Figure 9 shows the redshift distributions of close pairs, after subtracting the contribution of random superpositions, for the entire photometric sample (use\_phot = 1) in each of the 3D-HST fields (in the *F140W* footprint) down to the approximate magnitude limit of the grism redshift sample ( $JH_{\text{IR}} \leq 24$ ). Pairs are selected within 2–30 arcsec of separation, with the lower limit to avoid erroneous correlations due to blending in the IRAC images. Errors on the photometric redshifts of the pairs in the sample are related to the Gaussian distribution of their redshift deviations as

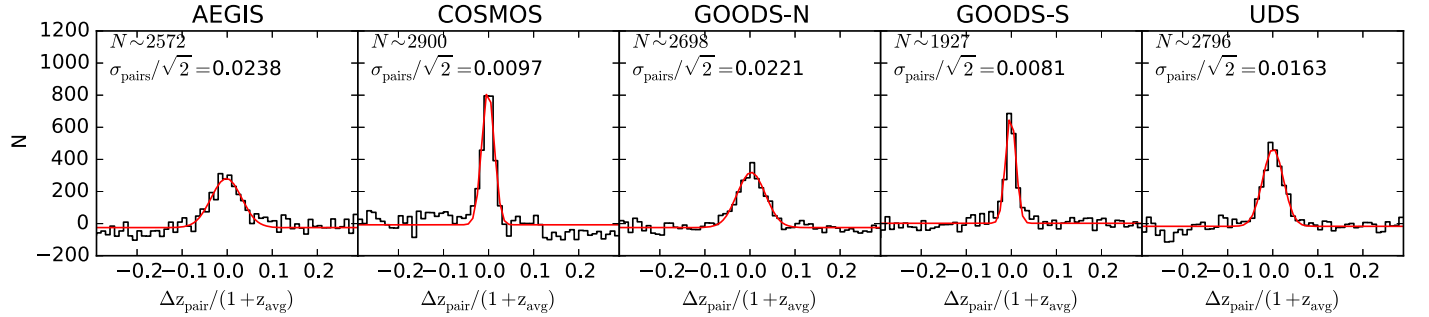
$$\sigma(\text{phot}) = \sigma \left( \frac{\Delta z_{\text{pair}}}{1 + z_{\text{avg}}} \right) \frac{1}{\sqrt{2}}, \quad (2)$$

where  $\sigma(\text{phot})$  is the representative photo- $z$  scatter,  $\Delta z_{\text{pair}}$  is the difference between the photometric redshifts ( $z_1$  and  $z_2$ ) of pair members,  $z_{\text{avg}}$  is their average redshift ( $z_{\text{avg}} = z_1 + z_2$ ), and  $\sigma \left( \frac{\Delta z_{\text{pair}}}{1 + z_{\text{avg}}} \right)$  is the width of the best fit Gaussian to the distribution of normalized redshift deviations. As with the grism redshift comparisons, the measured photometric redshift uncertainty exhibits clear field-to-field variation (see Section 4 for a more detailed discussion of inhomogeneous photometric data). We calculate the overall scatter as the average of all five fields and utilize jackknife resampling to estimate errors given that the formal errors on each individual Gaussian fit are negligible with respect to field-to-field variation. We calculate the photometric redshift errors following this method for galaxy pairs in bins of mass, magnitude, color, and redshift, only including pairs for which both galaxies are included in the selection. Figure 10 shows the measured photometric redshift errors as a function of stellar mass (left panel),  $JH_{\text{IR}}$  magnitude (second panel), rest-frame  $U - V$  color (third panel), and redshift (right panel). There is significant variation among fields (see Section 4); however, we find reasonable agreement between the average error estimated by this methodology (gray bands) and the  $\sigma_{\text{NMAD}}$  from direct comparison between grism and photometric redshifts (shown as gray dashed lines in each panel for comparison). For an evaluation of grism redshift accuracy via close pairs analysis, we refer the reader to Momcheva et al. (2015).

Again these estimates of photometric redshift errors do not exhibit strong trends with stellar mass or  $U - V$  color (as in Figure 10 from grism redshift tests), although there is a clear

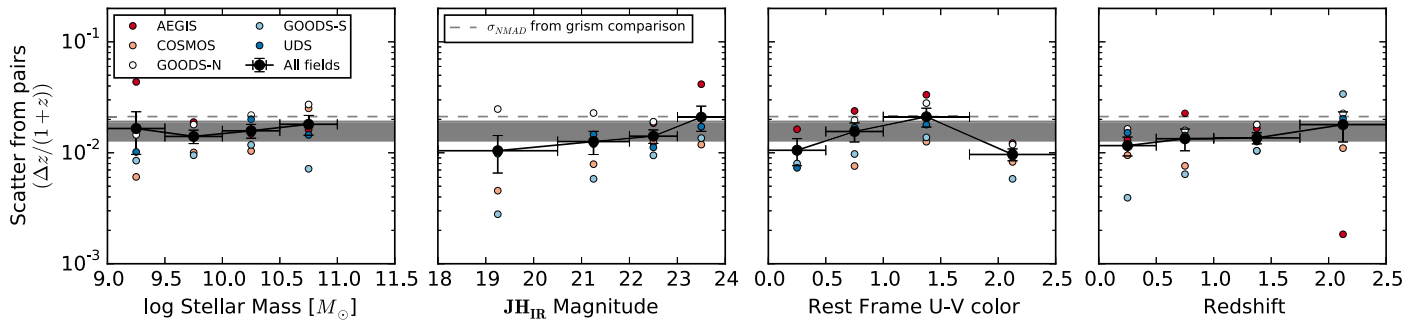


**Figure 8.** Dependence of photometric redshift accuracy on galaxy size at fixed magnitude (left column) and fixed redshift (second, third, and fourth columns). The photometric aperture used in creating the 3D-HST photometric catalogs is indicated by a horizontal dashed line in all of the top panels. The top row shows the size distribution of all galaxies (first and second columns) and for star-forming (third column) and quiescent galaxies (fourth column). The second row includes trends in NMAD scatter between photometric and grism redshifts and the bottom row includes trends in redshift outlier percentages. In each panel, the measured scatter or outlier percentage for the subdivided galaxy samples are indicated by the symbol size: largest filled circles for the largest galaxies and the smallest circles for the most compact galaxies at fixed mass or redshift. Photometric redshift accuracy depends most strongly on size at fixed magnitude, with extended galaxies exhibiting up to  $\sim 30\%$  higher scatter and  $\sim 5$  times higher outlier rates than small galaxies. Photometric redshift accuracy exhibits a weaker trend with size at fixed redshift, but this is primarily due to star-forming galaxies; the photometric redshift accuracy does not significantly depend on galaxy size for quiescent galaxies.



**Figure 9.** Distribution of photometric redshift differences  $(z_1 - z_2)/(1 + (z_1 + z_2)/2)$  for close pairs in each 3D-HST fields. The characteristic  $1\sigma$  error in photometric redshift is approximated by  $\sigma/\sqrt{2}$ , where  $\sigma$  is the Gaussian width of the distribution of redshift differences (Quadri & Williams 2010). The approximate number of pairs, calculated by integrating the Gaussian fits, and measured photometric redshift errors are indicated in the upper left corner of each panel. Photometric redshift accuracy varies significantly among fields, at least in part by differing photometric coverage.





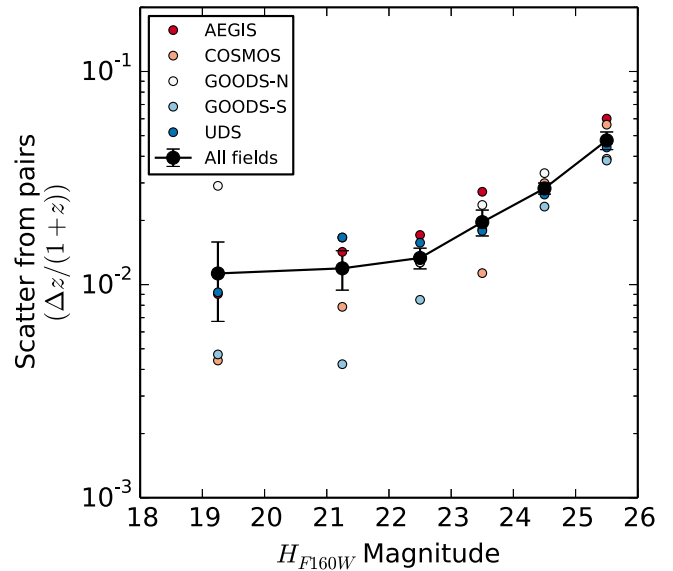
**Figure 10.** Photometric redshift accuracy from close pairs analysis as a function of stellar mass (left panel),  $JH_{\text{IR}}$  magnitude (second panel), rest-frame  $U - V$  color (third panel), and redshift (right panel). Average error derived from pairs analysis is indicated by gray shaded band and average NMAD scatter between photometric and grism redshifts by the gray dashed horizontal line. Redshift errors derived from close pairs are slightly lower than those derived from direct grism redshift comparisons; however, trends of increasing errors with magnitude and redshift persist.

decrease in scatter at the reddest colors. We note that galaxy pairs may sample the galaxy mass distribution differently than that of the full grism sample, which could lead to subtle differences in the estimates of redshift accuracies. On the other hand, the pairs-derived photometric errors depend slightly more strongly on  $JH_{\text{IR}}$  magnitude. Part of this is due to a small number of bright ( $<20.5$ ) pairs of galaxies, however the trend extends to faint magnitudes. This may in part be due to the fact that the grism redshift sample is less complete below  $z \sim 0.7$  (see Figure 1). Although we only use grism redshifts with tightened PDFs to test photometric redshift accuracy, we note that if residual correlations between photometric and grism redshifts depend on galaxy properties, this could alter trends in redshift performance. In this case, trends evaluated by the analysis of close pairs could be stronger than for the grism comparisons. We expect this correlation to be higher for galaxies without emission lines, particularly those with fainter continuum. This effect could contribute to the differences in scatter with magnitude. On the other hand, the photometric redshift errors estimated using close pairs could be artificially diminished by redshift “attractors” in photometric redshift space that artificially place galaxies at the same redshifts (Quadri & Williams 2010).

Finally, one benefit of photometric redshift accuracy based on statistical pairs analysis is that it is limited to the photometric depths, not those of the grism redshifts. In the case of the 3D-HST grism catalogs, redshift fits can be made to an arbitrary limit; however, we expect that these will only be valuable for galaxies with emission lines at fainter magnitudes. Therefore, we have limited our analysis to galaxies with visually inspected redshift fits at  $JH_{\text{IR}} < 24$ . We now extend the study of close pairs down to a fainter limiting magnitude ( $H_{F160W} < 26$ ), beyond the limits of the grism redshift estimates. Figure 11 shows the measured scatter as a function of  $H_{F160W}$  magnitude. Indeed photometric redshift accuracy diminishes significantly beyond the limits of the grism redshifts, with scatter increasing by over a factor of two below  $JH_{\text{IR}} = 24$  and a factor of four across the whole magnitude range.

#### 4. THE DEPENDENCE OF REDSHIFT ACCURACY ON FILTERS

All of the five extragalactic fields have extraordinary photometric coverage; the excellent photometric redshift accuracy ( $\langle\sigma_{\text{NMAD}}\rangle \lesssim 0.02$ ) results from analysis of the fully sampled galaxy SEDs. In this section, we investigate the



**Figure 11.** Scatter as a function of magnitude as determined by pairs analysis of galaxies below the magnitude limits of the grism redshifts. Photometric redshift accuracy continues to diminish with redshift, with increased scatter of  $0.046 \pm 0.005(1+z)$  for the faintest objects with  $25 \leq H_{F160W} < 26$ .

importance of various categories of photometric data in determining photometric redshifts. For this analysis, we rerun EAZY for the 3D-HST photometric catalogs, including only specific subsets of filters. The EAZY code includes a redshift prior in the fitting, for which we use the  $K$  band magnitude in the default fitting. In cases where a  $K$  or  $Ks$  filter is not included in the subset, we use an  $R$  band magnitude prior. Filter combinations for each field are included in Table 3, including appropriate references.

For this test, we compare derived photometric redshifts with the spectroscopic and grism redshift measurements of  $z_{\text{true}}$ . Calculated scatter and outlier fractions between photometric and spectroscopic and photometric and grism redshifts are included in Table 4. In the latter comparison, we again only include grism redshifts with tightened PDFs to minimize the effects of correlated errors on the measured scatter. Full comparisons are included in the Appendix. Although there are still variations in the details of photometry in each field, we classify subsets of photometry into the following categories: (1) all filters, (2) *HST* imaging ( $F606W-F160W$ ), (3) *HST* imaging ( $F435W-F160W$ ), (4) *HST* and IRAC, (5) *HST*, IRAC, and broadband, ground-based near-IR imaging, (6)

**Table 3**  
3D-HST/CANDELS Field Filter Subsets

Field	Subset	Descriptive Label	Filters	References
AEGIS	<i>HST</i>	<i>HST(F606W–F160W)</i>	<i>F606W, F814W, F125W, F140W, F160W</i>	1, 2, 3
	IRAC	IRAC	3.6, 4.5, 5.8, 8.0 $\mu\text{m}$	4, 5
	NMBS	medium-band NIR	<i>J1, J2, J3, H1, H2, K</i>	6
	CFHTLS	broadband optical	<i>u, g, r, i, z</i>	7, 8
COSMOS	<i>HST</i>	<i>HST(F606W–F160W)</i>	<i>F606W, F814W, F125W, F140W, F160W</i>	1, 2, 3
	IRAC	IRAC	3.6, 4.5, 5.8, 8.0 $\mu\text{m}$	4, 9
	UltraVISTA	broadband NIR	<i>Y, J, H, K</i>	10
	NMBS	medium-band NIR	<i>J1, J2, J3, H1, H2, K</i>	6
	CFHTLS	broadband optical	<i>u, g, r, i, z</i>	7, 8
	Subaru	broad and medium-band optical	<i>B, V, r', i', z', IA427, IA464, IA484, IA505, IA527, IA624, IA679, IA709, IA738, IA767, IA827</i>	11
	GOODS-N	<i>HST</i>	<i>HST(F435W–F160W)</i>	<i>F435W, F606W, F775W, F850LP, F125W, F140W, F160W</i>
IRAC		IRAC	3.6, 4.5, 5.8, 8.0 $\mu\text{m}$	4, 13
HDFN		broadband optical	<i>U, B, V, r_c, i_c, z'</i>	14
MODS		broadband NIR	<i>J, H, Ks</i>	15
GOODS-S	<i>HST</i>	<i>HST(F435W–F160W)</i>	<i>F435W, F606W, F775W, F814W, F850W, F125W, F140W, F160W</i>	1, 2, 3, 12
	IRAC	IRAC	3.6, 4.5, 5.8, 8.0 $\mu\text{m}$	4, 13
	GaBoDs	broadband optical	<i>U38, B, V, R_c, I</i>	16, 17
	MUSYC	medium-band optical	<i>IA427, IA445, IA505, IA527, IA550, IA574, IA598, IA624, IA651, IA679, IA738, IA767, IA797, IA856</i>	18
	FIREWORKS	broadband NIR	<i>J, H, Ks</i>	19, 20
	UDS	<i>HST</i>	<i>HST(F606W–F160W)</i>	<i>F606W, F814W, F125W, F140W, F160W</i>
IRAC		IRAC	3.6, 4.5, 5.8, 8.0 $\mu\text{m}$	4, 21
SXDS		broadband optical	<i>B, V, R_c, i', z'</i>	22
UKIDSS		broadband NIR	<i>J, H, Ks</i>	23

**References.** (1) Grogin et al. (2011), (2) Koekemoer et al. (2011), (3) Brammer et al. (2012), (4) Ashby et al. (2013), (5) Barmby et al. (2008), (6) Whitaker et al. (2011), (7) Erben et al. (2009), (8) Hildebrandt et al. (2009), (9) Sanders et al. (2007), (10) McCracken et al. (2012), (11) Taniguchi et al. (2007), (12) Giavalisco et al. (2004), (13) Dickinson et al. (2003), (14) Capak et al. (2004), (15) Kajisawa et al. (2011), (16) Hildebrandt et al. (2006), (17) Erben et al. (2005), (18) Cardamone et al. (2010), (19) Wuyts et al. (2008), (20) Retzlaff et al. (2010), (21) J. Dunlop et al. (in preparation), (22) Furusawa et al. (2008), (23) O. Almaini et al. (in preparation).

*HST*, IRAC, and medium-band optical imaging, (7) *HST*, IRAC, and medium-band near-IR imaging, (8) broadband, ground-based optical, and near-IR imaging, (9) medium-band, ground-based optical, and near-IR imaging. However, we note that the data included in each category will still vary in specific filter sets, photometric depths, and data quality. Scatter and outlier fractions between photometric and spectroscopic redshifts in each of these categories and fields is included in Figure 12 and for photometric and grism redshifts in Figure 13.

The first thing to notice in Figure 12 is that the scatter between photometric and spectroscopic redshifts varies significantly from field to field. This is due to a number of different factors, including heterogeneity in both the available photometry and spectroscopic follow-up in addition to cosmic variance. Overall scatter is lowest in the COSMOS field ( $\sigma_{\text{NMAD}} = 0.008$ ) and highest in GOODS-N ( $\sigma_{\text{NMAD}} = 0.027$ ). This may in part be due to the optical and near-IR medium-band photometry in the COSMOS field. Scatter is also low in GOODS-S, where photometry includes medium-band filters in the optical from the MUSYC survey, whereas the scatter in AEGIS is somewhat higher ( $\sigma_{\text{NMAD}} = 0.023$ ) even though the Newfirm Medium Band Survey (NMBS) near-IR medium-band filters are included. The relative importance of optical medium-band filters is partially due to the redshift distribution of the spectroscopic comparison sample, most of which are at low redshift. Using the grism redshifts, we can overcome this bias and assess the importance of filters in setting the photometric redshift accuracy for a more representative sample.

Comparing to grism redshifts has the effect of normalizing variable spectroscopic redshift quality and quantity (Figure 13). In this case, the scatter and outlier fractions are somewhat more uniform across the five fields when all filters are included. In fields where *HST* ACS *F435W* imaging is not available (AEGIS, COSMOS, UDS), the scatter and outlier fraction of photometric redshifts are significantly higher when only *HST* imaging is fit. The inclusion of the blue filter in GOODS-N and GOODS-S introduces a decrease in scatter when only *HST* imaging is used to estimate photometric redshifts, although the effect is weaker than in the spectroscopic comparison. This emphasizes the importance of including blue wavelengths to identify spectral features such as the Lyman Break and the stellar bump. Improvements in identifying these features can be made by including deep optical imaging, in particular medium-band imaging. Another striking improvement in the photometric redshift accuracy is gained with the addition of *Spitzer*-IRAC imaging, in some cases decreasing the scatter by over a factor of two.

It is interesting to note that the photometric redshift accuracy is uniformly worse when compared to the grism redshift sample than relative to the spectroscopic redshift sample, especially in the outlier fractions. This is likely because the grism redshifts probe regions of parameter space where photometric redshifts are harder to measure: fainter galaxies, higher redshifts, and different galaxy types. Even with *HST* imaging and *Spitzer*-IRAC photometry, the outlier fractions range from  $\sim 14\%$  to  $\sim 21\%$  without the bluer *HST* imaging. Systematics in these redshift measurements are examined in individual fields in the

**Table 4**  
Photometric Redshift Accuracy with Filters

Field	Subsets	Spectroscopic Redshift Comparison		Grism Redshift Comparison	
		$\sigma_{\text{NMAD}}$	Outlier %	$\sigma_{\text{NMAD}}$	Outlier %
AEGIS	all	$0.024 \pm 0.001$	$2.5\% \pm 0.6$	$0.027 \pm 0.001$	$5.1\% \pm 0.5$
	<i>HST</i>	$0.089 \pm 0.006$	$31.6\% \pm 2.4$	$0.090 \pm 0.003$	$31.2\% \pm 1.0$
	<i>HST</i> , IRAC	$0.042 \pm 0.002$	$5.4\% \pm 0.8$	$0.050 \pm 0.001$	$14.4\% \pm 0.7$
	<i>HST</i> , IRAC, NMBS	$0.037 \pm 0.002$	$6.5\% \pm 0.9$	$0.040 \pm 0.001$	$11.8\% \pm 0.7$
	CFHTLS, NMBS	$0.026 \pm 0.001$	$13.2\% \pm 1.3$	$0.050 \pm 0.002$	$21.2\% \pm 0.9$
COSMOS	all	$0.007 \pm 0.001$	$1.2\% \pm 0.6$	$0.012 \pm 0.000$	$1.9\% \pm 0.3$
	<i>HST</i>	$0.098 \pm 0.005$	$29.8\% \pm 3.1$	$0.102 \pm 0.003$	$35.2\% \pm 1.2$
	<i>HST</i> , IRAC	$0.038 \pm 0.002$	$5.0\% \pm 1.2$	$0.057 \pm 0.002$	$17.7\% \pm 1.0$
	<i>HST</i> , IRAC, UVISTA	$0.029 \pm 0.003$	$5.0\% \pm 1.2$	$0.027 \pm 0.001$	$8.3\% \pm 0.7$
	<i>HST</i> , IRAC, NMBS	$0.028 \pm 0.003$	$3.8\% \pm 1.0$	$0.032 \pm 0.001$	$9.6\% \pm 0.7$
	<i>HST</i> , IRAC, Subaru	$0.008 \pm 0.001$	$1.2\% \pm 0.6$	$0.017 \pm 0.001$	$3.4\% \pm 0.4$
	CFHTLS, UVISTA	$0.014 \pm 0.001$	$0.9\% \pm 0.5$	$0.026 \pm 0.001$	$4.2\% \pm 0.5$
	NMBS, Subaru	$0.008 \pm 0.001$	$1.5\% \pm 0.6$	$0.017 \pm 0.001$	$3.3\% \pm 0.4$
GOODS-N	all	$0.026 \pm 0.001$	$8.5\% \pm 0.9$	$0.023 \pm 0.001$	$3.7\% \pm 0.4$
	<i>HST</i>	$0.039 \pm 0.002$	$15.2\% \pm 1.2$	$0.036 \pm 0.001$	$8.1\% \pm 0.6$
	<i>HST</i> , 600, 160	$0.089 \pm 0.005$	$33.1\% \pm 1.5$	$0.054 \pm 0.002$	$16.5\% \pm 0.8$
	<i>HST</i> , IRAC	$0.031 \pm 0.001$	$10.8\% \pm 1.0$	$0.027 \pm 0.001$	$5.0\% \pm 0.5$
	<i>HST</i> , IRAC, HDFN	$0.032 \pm 0.001$	$8.9\% \pm 0.9$	$0.029 \pm 0.001$	$4.4\% \pm 0.4$
	<i>HST</i> , IRAC, MOIRCS	$0.029 \pm 0.002$	$9.8\% \pm 1.0$	$0.024 \pm 0.001$	$4.5\% \pm 0.4$
	HDFN, MOIRCS	$0.044 \pm 0.002$	$10.1\% \pm 1.0$	$0.050 \pm 0.001$	$11.3\% \pm 0.6$
	all	$0.010 \pm 0.001$	$5.2\% \pm 0.8$	$0.013 \pm 0.000$	$2.8\% \pm 0.4$
GOODS-S	<i>HST</i>	$0.036 \pm 0.002$	$10.2\% \pm 1.2$	$0.034 \pm 0.001$	$10.2\% \pm 0.6$
	<i>HST</i> , 600, 160	$0.063 \pm 0.005$	$26.3\% \pm 1.7$	$0.048 \pm 0.002$	$20.7\% \pm 0.9$
	<i>HST</i> , IRAC	$0.026 \pm 0.002$	$5.8\% \pm 0.9$	$0.028 \pm 0.001$	$8.0\% \pm 0.6$
	<i>HST</i> , IRAC, MUSYC	$0.013 \pm 0.001$	$6.5\% \pm 0.9$	$0.017 \pm 0.001$	$5.0\% \pm 0.5$
	GaBoDs, FIREWORKS	$0.033 \pm 0.002$	$10.0\% \pm 1.2$	$0.053 \pm 0.002$	$19.4\% \pm 0.8$
	MUSYC, FIREWORKS	$0.012 \pm 0.001$	$8.8\% \pm 1.1$	$0.024 \pm 0.001$	$14.5\% \pm 0.8$
	all	$0.021 \pm 0.002$	$6.3\% \pm 2.2$	$0.024 \pm 0.001$	$3.7\% \pm 0.4$
	<i>HST</i>	$0.075 \pm 0.009$	$24.7\% \pm 4.5$	$0.071 \pm 0.002$	$19.8\% \pm 0.9$
UDS	<i>HST</i> , IRAC	$0.044 \pm 0.007$	$9.9\% \pm 2.8$	$0.050 \pm 0.002$	$12.8\% \pm 0.7$
	<i>HST</i> , IRAC, SXDS	$0.026 \pm 0.003$	$8.1\% \pm 2.7$	$0.028 \pm 0.001$	$4.7\% \pm 0.4$
	<i>HST</i> , IRAC, UKIDSS	$0.034 \pm 0.006$	$9.9\% \pm 2.8$	$0.033 \pm 0.001$	$10.0\% \pm 0.6$
	SXDS, UKIDSS	$0.036 \pm 0.003$	$7.2\% \pm 2.5$	$0.032 \pm 0.001$	$5.4\% \pm 0.5$

[Appendix](#). This highlights the utility of the 3D-HST sample, but also emphasizes the importance of assessing the breadth of any spectroscopic sample used to evaluate photometric redshift performance.

We emphasize that collapsing the redshift accuracy into these three measures disguises systematics introduced by different filter combinations. For example, low redshift ( $z \lesssim 0.5-1$ ) galaxies are particularly driving the increased scatter in photometric redshift when only *HST* imaging (*F616W-F160W*) is included. We included detailed figures including the redshift scatter in each individual field in the [Appendix](#).

## 5. PHOTOMETRIC REDSHIFT PROBABILITY DISTRIBUTION FUNCTIONS

In addition to fitting a single-valued photometric redshift estimate, the EAZY code produces individual probability distribution functions. These PDFs provide an estimate of the likelihood that the galaxy lies at a given true redshift. Until this point, we have adopted the redshift with the maximum likelihood ( $z_{\text{peak}}$ ) as the photometric redshift for each galaxy. For certain applications, we would like to incorporate the uncertainty on photometric redshift and ideally utilize the entire PDF function to describe the probability that the galaxy lies at a given  $z_{\text{true}}$ . This technique has been proven to significantly improve measurement uncertainties, for example it can increase

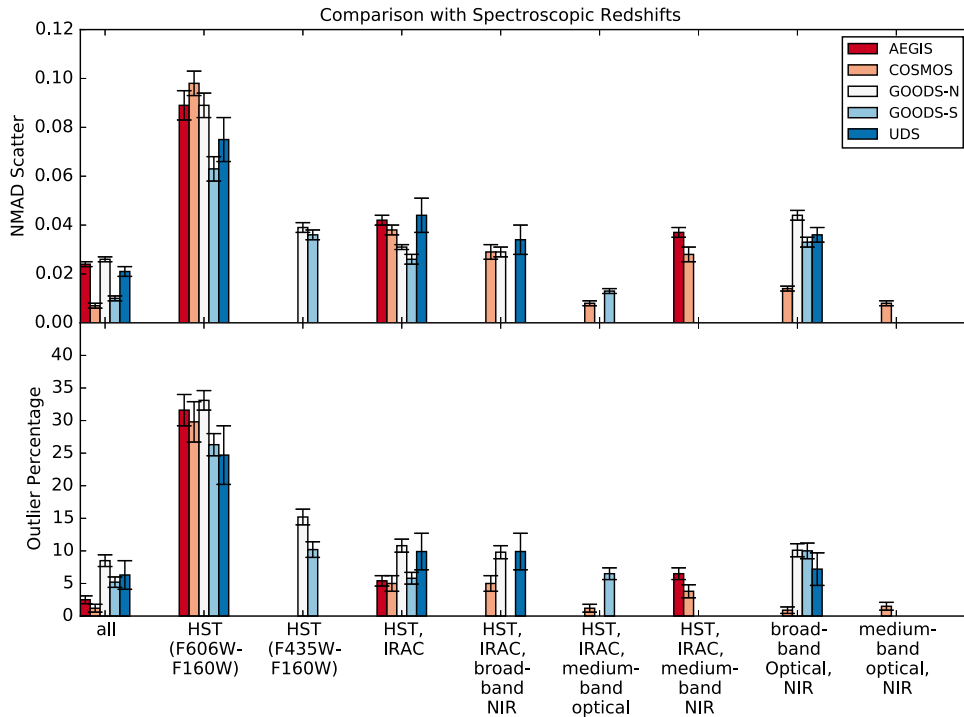
the S/N of clustering measurements by a factor equivalent to an increase in survey size of  $\sim 2-3$  (Myers et al. 2009). In this section, we investigate the ability of the EAZY-generated PDFs to predict the  $z_{\text{true}}$  values for the ensemble of 3D-HST galaxies.

### 5.1. Photometric Redshift Confidence Intervals

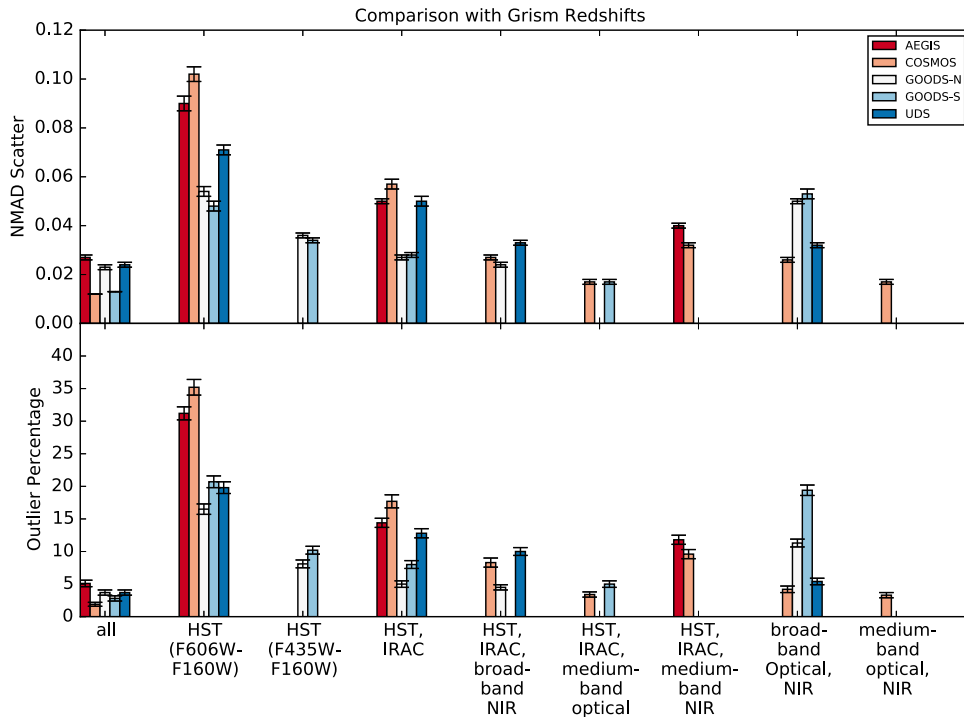
A key question regarding photometric redshift performance is whether the scatter between measured and true redshifts is primarily driven by uncertainties in the photometric redshift estimates. For this, we test the redshift deviations for individual galaxies relative to their estimated confidence intervals. In Figure 14 we show the distribution of deviations ( $z_{\text{phot}} - z_{\text{true}}$ ) normalized by the 68% (left panel), 95% (center panel), and 99.7% photometric redshift confidence intervals. Each panel includes two samples: the orange histogram indicates the distribution for a comparison with independent grism redshifts (with narrowed PDFs) and a purple histogram for the spectroscopic sample. Although the normalization between the two samples and confidence intervals differs, each is characterized by a Gaussian central peak in addition to broader wings. Best-fit Gaussians are fit to each distribution and standard deviations are indicated in the legend of each panel. These Gaussian distributions contain roughly 90% of galaxies in each panel.

For fully representative photometric redshift errors, we would expect a Gaussian of width  $\sigma = 1.0$  for the 68%





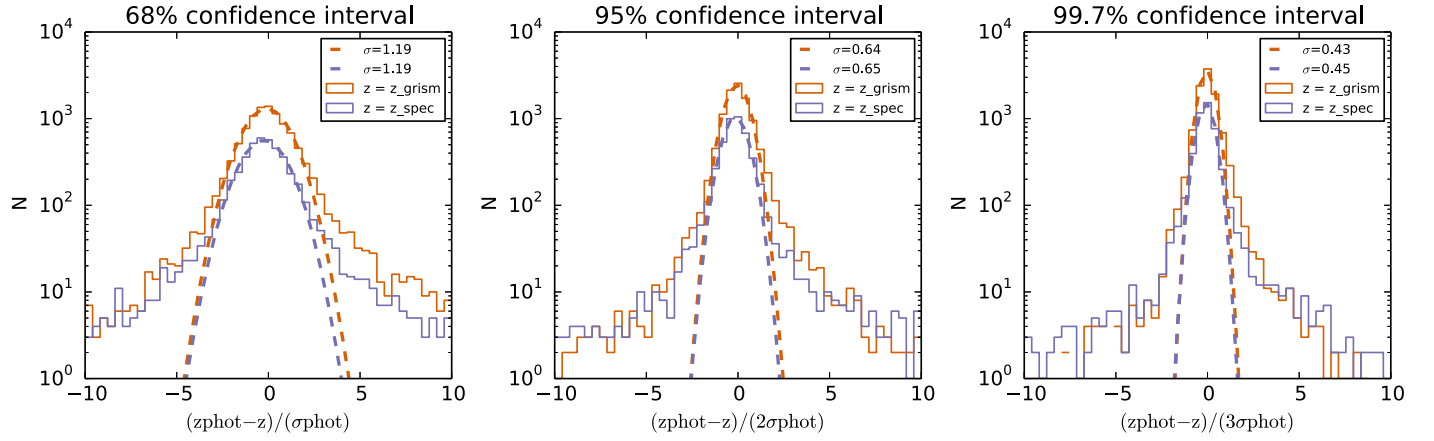
**Figure 12.** Photometric redshift accuracy with different filter combinations compared to spectroscopic redshifts in each 3D-HST field. Photometric redshift accuracy depends strongly on the photometric bandpasses included in redshift fitting; at these magnitudes and redshifts blue optical imaging is crucial. Redshift accuracy varies strongly among fields, even when similar data sets are included. Some of this variation may be due to heterogeneous spectroscopic redshift samples across the fields.



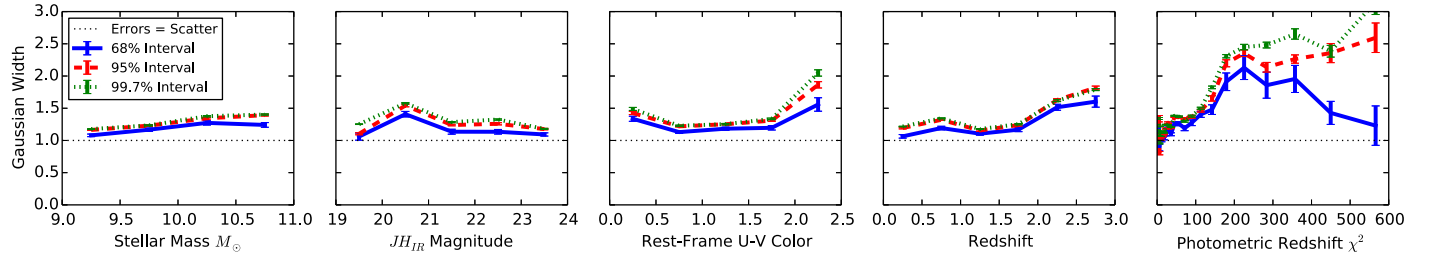
**Figure 13.** Photometric redshift accuracy with different filter combinations compared to tightened grism redshifts. Again this figure demonstrates the strong bandpass dependence of photometric redshift accuracy, with uniformly distributed grism redshifts. When all photometry is included in the fit, there is less variation in redshift quality between fields than in spectroscopic comparison (Figure 12), but the comparison sample is insufficient to explain all field-to-field variation.

confidence interval and  $\sigma = 1/2, 1/3$  when redshift deviations are normalized by the  $2\sigma$  and  $3\sigma$  errors. Like Dahlen et al. (2013), we find that the photometric redshift PDFs are too

narrow in each confidence interval. However, the factor by which the PDFs would need to be broadened differs for each test. At the 68% confidence level, the photometric PDFs are a



**Figure 14.** Redshift deviations normalized by 68%, 95%, and 99.7% confidence intervals for photometric redshifts  $((z_{\text{phot}} - z_{\text{true}})/(\sigma_{\text{phot}}))$  in the left panel,  $(z_{\text{phot}} - z_{\text{true}})/(2\sigma_{\text{phot}})$  in the center panel, and  $(z_{\text{phot}} - z_{\text{true}})/(3\sigma_{\text{phot}})$  in the right panel. Comparisons between photometric redshifts and spectroscopic redshifts (purple) or narrowed grism redshifts (orange) yield similar distribution shapes. Both exhibit roughly Gaussian distributions (fits are indicated with dashed lines) but  $\sigma \sim 1.2, 0.6, 0.4$  suggesting that the redshift PDFs are narrower than the observed scatter in redshift by a factor of  $\sim 1.2$ .



**Figure 15.** Underestimate factor for photometric redshift errors as a function of galaxy properties: stellar mass, apparent magnitude,  $U - V$  color, redshift, and photo- $z$   $\chi^2$  as measured by fitting Gaussians to the scatter in redshift deviations from grism redshifts normalized by photometric error as in Figure 14. The dotted black line at unity indicates the value at which errors explain the observed redshift scatter. Solid lines indicate scatter normalized by  $1\sigma$  error, which is always slightly lower than the  $2\sigma(3\sigma)$  width multiplied by a factor of two(three), indicating that the PDFs are too narrow to explain the photometric redshift scatter, particularly in the tails of the distribution.

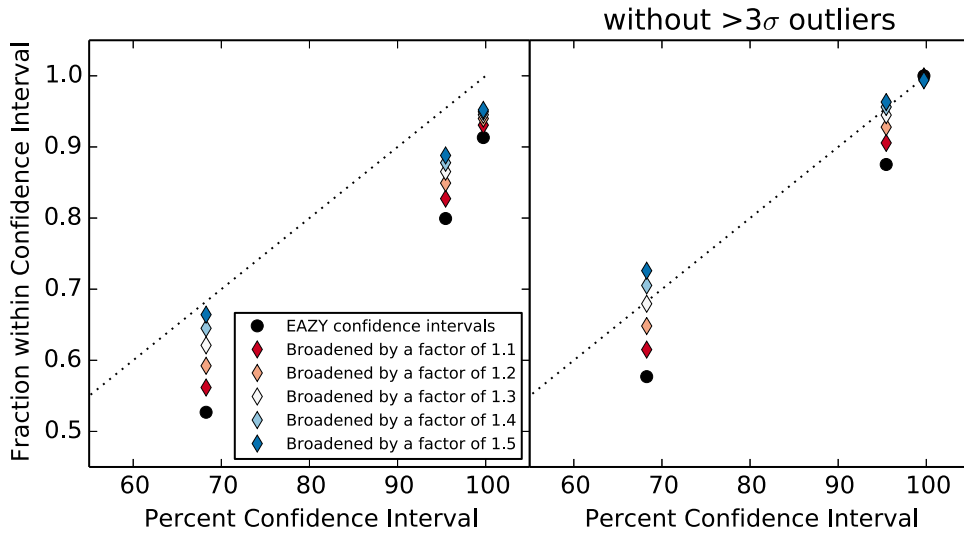
factor of  $\sim 1.2$  too narrow, whereas the tails of the PDFs are further underestimated, requiring a factor of  $\sim 1.6$  to explain the observed scatter between photometric and true redshift.

To complicate the situation, this discrepancy is not uniform among galaxy types. Figure 15 further dissects the trends in uncertainty underestimation by galaxy stellar mass, apparent magnitude, rest-frame  $U - V$  color, redshift, and  $\chi^2$  from the photometric redshift fit using the sample of galaxies with narrowed grism PDFs. As for the total sample, the underestimation of the 68% confidence interval for photometric redshift errors is less than for the 95% and 99.7% confidence intervals. Furthermore, we see clear trends that this depends on galaxy properties. In contrast with the measured scatter in photometric redshifts with stellar mass, photometric redshift errors are decreasingly well calibrated with increasing mass. Aside from the brightest galaxies, which appear to have appropriate error estimates, the normalized scatter does not depend strongly on apparent magnitude or galaxy color. However, the uncertainties are underestimated by an increasing  $\sim 1.1$  at low redshifts to  $\sim 1.6$  at  $z \sim 2.5$ . Finally, for very poorly fit photometric redshifts ( $\chi^2 \gtrsim 100$ ), the scatter in redshifts is vastly underpredicted by EAZY by up to a factor of  $\sim 2.5$ . The right panel indicates strong correlation between photometric redshift scatter, as normalized by the redshift confidence intervals; at the largest  $\chi^2$  values this normalization may include the entire allowed redshift range, driving Gaussian widths back to 1.0.

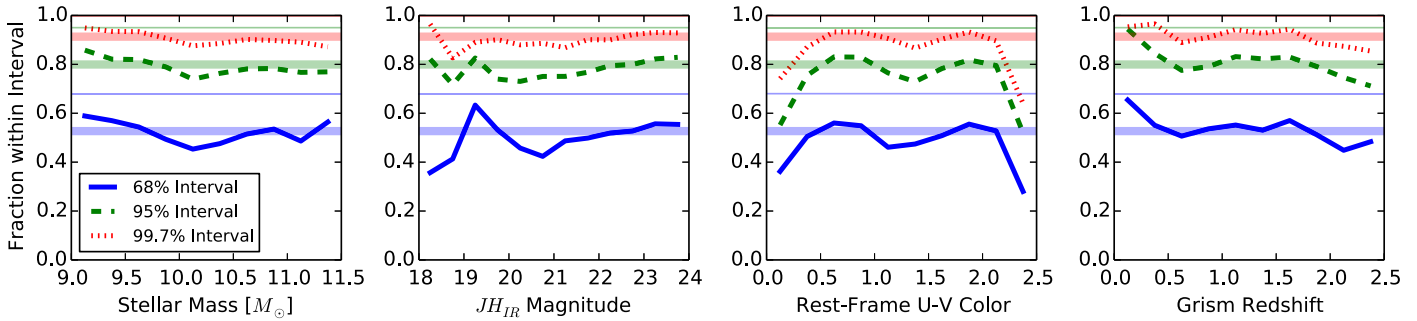
## 5.2. PDF Width and Quantifying Catastrophic Outliers

In the previous section, we demonstrated that the error estimates for the majority of galaxies are underestimated by approximately a factor of 1.2–1.3 by looking at the distribution of scatter between photometric and grism redshifts. However, there are tails of the distribution of redshift scatter for which the errors cannot be described by a Gaussian distribution. We now investigate the properties of these outliers. Following Dahlen et al. (2013), if the redshift error estimates are accurate for the entire population of galaxies,  $\sim 68\%$  of galaxies will have 68% confidence intervals that include  $z_{\text{true}}$  and likewise for the 95% and 99.7% confidence intervals.

Figure 16 indicates the fraction of galaxies within the 68%, 95%, and 99.7% confidence intervals as black circles for the entire sample (left panel) and the sample for which  $z_{\text{grism}}$  falls within the 99.7% photometric redshift confidence interval (right panel). Colored diamonds demonstrate the fractions measured by artificially broadening the confidence intervals. One-to-one correspondence relations are included as black dotted lines. Clearly, the fraction of galaxies within a given confidence interval is well below the predicted value, partially due to underestimated errors. Even by extending the confidence intervals by a range of factors, there is always a fraction of galaxies for after cropping catastrophic outliers, defined such that  $z_{\text{grism}}$  lies well outside the  $3\sigma$  error estimates. These catastrophic fitting failures drive the overall fractions lower than can be explained by inflating error bars alone. In the right



**Figure 16.** Fraction of galaxies with  $z_{\text{grism}}$  within photometric redshift confidence intervals. The left panel includes the full sample and the right panel only includes galaxies for which  $z_{\text{grism}}$  falls within the 99.7% confidence intervals. Overall  $\sim 10\%$  of galaxies will have grossly underestimated photometric redshift uncertainties, and confidence intervals are too narrow by a factor of  $\sim 1.2$ .



**Figure 17.** Trends in the fraction of galaxies for which grism redshifts fall within photometric redshift confidence intervals (blue solid, green dashed, red dotted lines correspond to  $1\sigma$ ,  $2\sigma$ , and  $3\sigma$ . Average values are indicated by thick horizontal lines and expected 68%, 95%, and 99.7% values are indicated by thin colored lines.

panel, we demonstrate that by excluding these outliers (approximately 10% of galaxies) and then broadening the error estimates by a factor of 1.2–1.3 found in the previous section, producing general agreement between the confidence intervals and redshift distributions of galaxies.

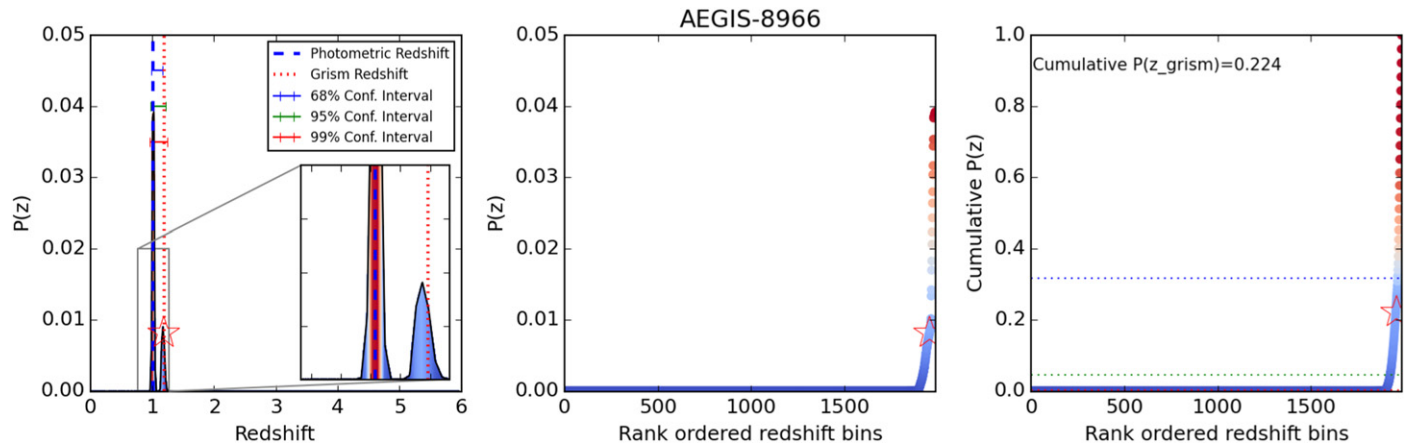
These fractions do not depend strongly on galaxy properties. In Figure 17 we show trends in these fractions as a function of stellar mass,  $JH_{\text{IR}}$  magnitude, rest-frame  $U - V$  color, and grism redshift. Average values are indicated as solid, horizontal lines and trends with galaxy properties are shown in blue for the  $1\sigma$  confidence interval, dashed green for  $2\sigma$ , and dotted red for  $3\sigma$ . In each case, roughly 10% of galaxies lie outside of the  $3\sigma$  confidence intervals. When using a purely photometric sample of galaxies, this will correspond to noise in galaxy counts that will not be accounted for by photometric redshift uncertainties. This outlier rate is significantly higher than the outlier rates in  $z_{\text{phot}}$  versus  $z_{\text{grism}}$  comparisons, but may also be important to include for studies that include photometric redshift error estimates.

### 5.3. How Well Do Photometric PDFs Predict True Redshifts?

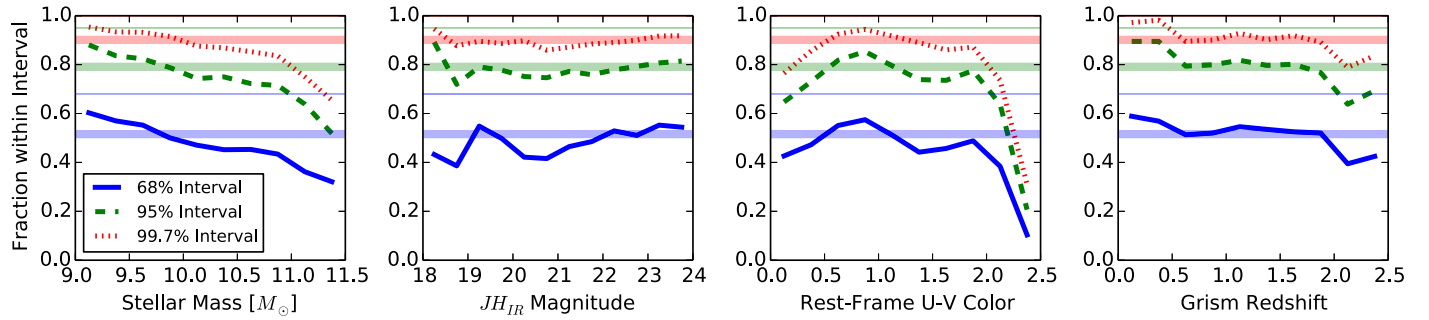
In this section we investigate the use of the full photometric redshift PDF as opposed to a single-valued photometric redshift with errorbars. In particular, this could be important for galaxies that have multi-peaked PDFs. We show an example

galaxy from the catalog in Figure 18. The  $P(z)$  for the galaxy in redshift bins is included in the left panel, along labeled photometric redshift (blue) and grism redshift (red dotted line and star). Confidence intervals are indicated by blue (68%), green (95%), and red (99.7%) errorbars. For this galaxy, the photometric redshift is assigned at the center of the most dominant peak of the PDF however the grism redshift reveals that the second peak is the location of the true redshift for this galaxy. In this specific case, the full PDF gives a clearer understanding of the uncertainties on the photometric redshift. Although the errorbars are somewhat broad, the actual redshift is well constrained between two narrower ranges.

To test the impact of such multi-peaked PDFs, we rank redshift bins by  $P(z)$  in the PDF of every galaxy (as in the center panel of Figure 18) and estimate the cumulative probability that corresponds to the redshift bin in which the grism redshift lies (right panel). In this specific example, grism redshift lies on a second redshift peak, outside of the 68% confidence interval; calculated in this way the  $P(z_{\text{grism}}) = 0.224$ . With this computed for each galaxy in the tightened grism sample, we repeat the test of enclosed fractions as a function of galaxy properties. Figure 19 presents the fraction of galaxies for which  $z_{\text{grism}}$  falls within the three defined confidence intervals. The average values for these fractions are remarkably similar to those derived via standard uncertainties. One can imagine problematic



**Figure 18.** Photometric PDF for sample galaxy AEGIS-8966. Photometric redshift is labeled with vertical, dashed blue line, grism redshift with red dotted line, and red star and color-coding in each panel corresponds to the  $P(z)$  value. Note that the photometric redshift lies in the most likely peak of the PDF, but the true redshift lies on the second peak. Second and third panels include rank-order distribution of  $P(z)$  bins (and cumulative  $P(z)$ ) with the location of  $z_{\text{grism}}$  indicated again by the star.



**Figure 19.** Trends in the fraction of galaxies for which grism redshifts fall within photometric redshift confidence intervals as measured from rank-ordered photometric PDFs. Color coding and resulting average fractions are extremely similar to those measured from confidence intervals in Figure 17, which suggests that confidence intervals estimated by EAZY from the CDFs are sufficient.

PDF with widely separated multiple peaks for which the confidence intervals defined by the EAZY code could overestimate the uncertainty, even though the errors on the whole are narrow with respect to the observed scatter. In practice, we do not find evidence that this is a dominant effect, likely because the subset of these significantly multi-peaked PDFs is small and therefore the average fraction of true redshifts that fall within a given confidence interval is largely independent of the method used to determine confidence intervals. Although for the most part we find similar trends in fractions with galaxy properties as in Figure 17, we do find a fairly strong trend of decreasing fractions of correctly estimated errors within each interval with increasing stellar mass. This is consistent with the trend in uncertainty normalized scatter which increases with stellar mass, as shown in Figure 15.

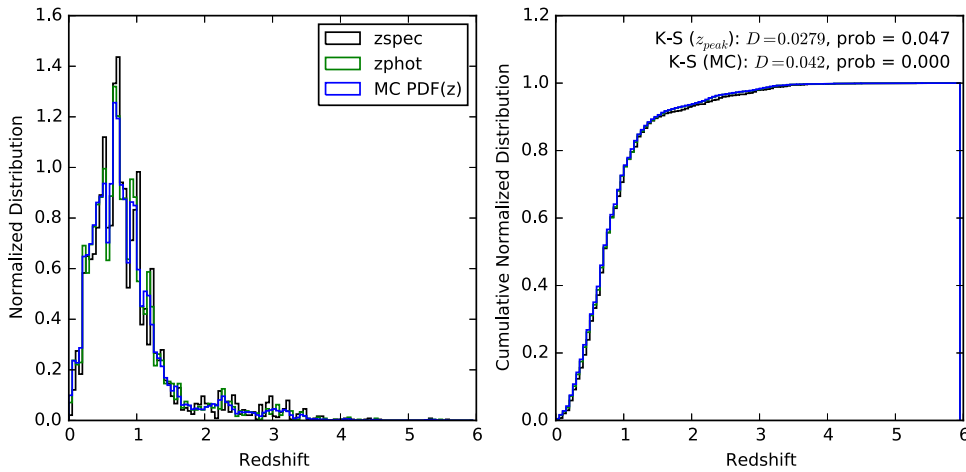
We now shift to a related issue, and investigate how well the photometric redshifts can recover redshift distributions of a sample of galaxies. We compare the overall redshift distribution (both of the spectroscopic sample and of the tightened grism sample) to the distribution of single-valued photometric redshifts and those derived by bootstrap resampling of the individual photometric PDFs. The redshift histogram and cumulative distributions for the spectroscopic sample is presented in Figure 20 and for the grism sample in Figure 21. We emphasize that while we expect the photometric and true redshift distributions to be similar, measurement errors and catastrophic outliers will broaden the photometric redshift

distributions. Therefore, these distributions should fail a simple Kolmogorov–Smirnov (K-S) test from a statistical standpoint; however, the K-S statistic  $D$ , or the maximum distance between each cumulative distribution functions, is still an informative metric of the relative similarity of the two distributions.

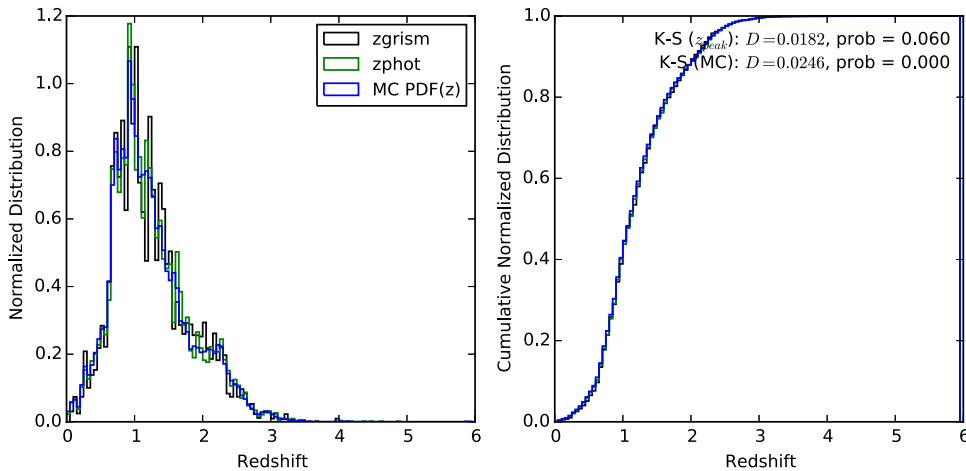
In Figure 20, the two-sided K-S test between the spectroscopic and  $z_{\text{peak}}$  redshift distributions indeed suggests a very low probability that the two samples are drawn from the same distribution. However, we find excellent agreement between the two histograms, with maximum deviations of less than 3%. The distribution of Monte-Carlo sampled redshifts exhibits a larger deviation, but still agrees to within 4% with the spectroscopic redshifts. We interpret this as due to the additional scatter to each galaxy due to the PDF resampling. The comparison with grism redshifts (Figure 21) yields similar results: 2% deviation for the comparison with  $z_{\text{peak}}$  and 2.5% deviations for the resampled redshifts.

The full redshift distributions are not statistically consistent in either case; however, given their similarity there are benefits to using the full redshift PDF. We note that although the distributions are more similar for the comparison with the single-valued  $z_{\text{peak}}$  photometric redshift estimates, there are instances in which the full redshift PDF carries additional useful information. For example, in the case of a multi-peaked PDF the  $z_{\text{peak}}$  carries only information about the most likely redshift, even when the secondary peak may be nearly as significant. Sampling the entire PDF will scatter the redshifts





**Figure 20.** Distribution and cumulative distribution of spectroscopic redshifts (black histograms), photometric redshifts (green), and Monte-Carlo sampling of photometric PDFs (blue). Full distributions agree extremely well and although a KS test between the spectroscopic distribution and that of the photometric samples do not suggest that they are drawn from the same distribution, the distributions deviate by less than 3% for the photometric redshifts and  $\sim 4\%$  for the MC sampled PDFs.



**Figure 21.** Distribution and cumulative distribution of redshifts, as in Figure 20, for comparison with grism redshifts. Again the cumulative distributions of grism and photometric redshifts agree to within  $\sim 2\%$  and MC sampled PDFs to within  $\sim 3\%$ .

for all galaxies, hence the increased deviations, but for multip peaked PDFs will capture all possible solutions, given the fitting methodologies and templates. We conclude that photometric redshifts can reproduce the true redshift distribution of galaxies in a sample to a few percent accuracy, but emphasize that because these deviations may depend on the galaxy properties the effects of using photometric redshifts should be carefully modeled.

## 6. PHOTOMETRIC REDSHIFT ACCURACY IN SIMULATED SURVEYS

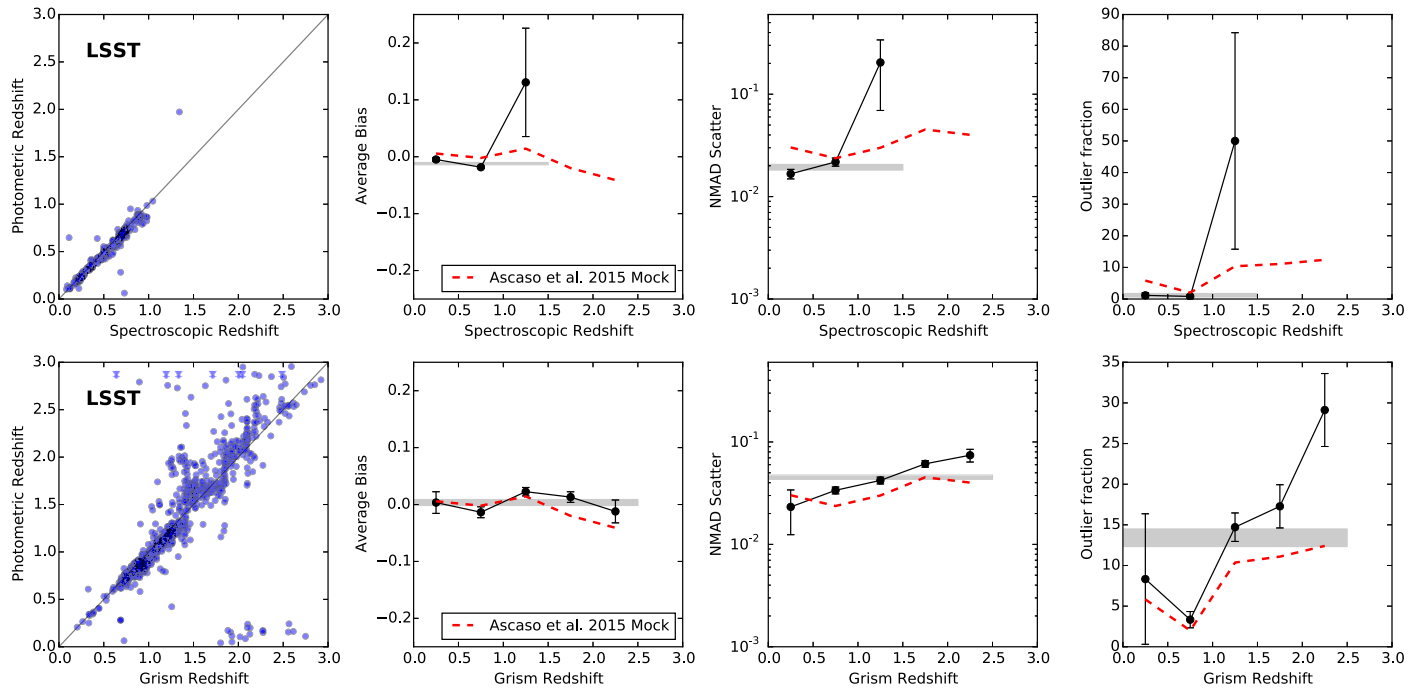
In addition to characterizing the performance of photometric redshifts in the 3D-HST survey, this vast data set can be used to predict or estimate the redshift accuracy in other surveys with similar photometric data. In this section, we extend our analysis to predict the photometric redshift accuracy in three major planned data sets using the 3D-HST catalogs: the LSST, the DES, and DES combined with the VHS. Because both of these surveys are planned to include a  $Y$  band filter, we limit this exercise to the COSMOS field where ground-based  $Y$  band imaging is included from the UltraVista Survey (McCracken et al. 2012). In each case we add noise to the 3D-HST catalogs to match the cited target depths of these surveys. Although the

catalogs are based on real data, we refer to them as data simulations. We use the EAZY code to fit photometric redshifts to the resulting catalogs and analyze the photometric redshift accuracy for each simulated survey.

### 6.1. LSST Survey

The LSST survey is planned to image  $>10,000 \text{ deg}^2$  in  $u$ ,  $g$ ,  $r$ ,  $i$ ,  $z$ , and  $y$  filters down to 26.1, 27.4, 27.5, 26.8, 26.1, and 24.9  $5\text{-}\sigma$  limiting magnitudes in the final coadded images (Ivezic et al. 2008). The 3D-HST catalog in the COSMOS field includes imaging in each of these bands; however, the  $g$ ,  $r$  and  $z$  band imaging from the CFHTLS survey (Erben et al. 2005; Hildebrandt et al. 2009) is slightly shallower than the proposed LSST depths. The differences are extremely small, less than  $\sim 0.3 \text{ mag}$  in each case. We add noise to the other catalog fluxes ( $u$ ,  $i$ ,  $Y$ ) to the planned depths prior to fitting photometric redshifts.

Results of the photometric redshift performance for the simulated LSST survey are shown in Figure 22. The top row includes a comparison of photometric and spectroscopic redshifts, the bottom row of photometric and grism redshifts. Each row includes a  $z_{\text{phot}}$  versus  $z_{\text{true}}$  scatter plot, average bias  $\langle \Delta z \rangle = (z_{\text{phot}} - z_{\text{true}})/(1 + z_{\text{true}})$ , scatter, and outlier



**Figure 22.** Simulated photometric redshift performance for the LSST survey with  $u$ ,  $g$ ,  $r$ ,  $i$ ,  $z$ , and  $y$  filters created by adding noise to the 3D-HST photometric catalog in COSMOS. Scatter predicted by mock galaxy catalogs is indicated by a dotted red line (Ascaso et al. 2015). Below  $z \sim 1$ , the accuracy is quite good when compared to the spectroscopic redshift sample  $\Delta z / (1+z) \sim 0.02$ . However, when the fainter and more representative grism redshift comparison sample is included, the scatter clearly depends strongly with redshift, increasing to  $\sim 0.04$  by  $z \sim 1.25$  and  $0 \sim 0.07$  by  $z \sim 2.5$ . Additionally, while the outlier fraction is excellent below  $z \sim 1$  for brighter spectroscopic redshifts, the percentage of  $|\Delta z| > 0.1$  outliers increases from  $\sim 5\%$  to  $\sim 30\%$  at  $z \sim 2.5$ .

fraction as a function of redshift. Ascaso et al. (2015) also conducted a simulation of the LSST survey using mock redshift catalogs based on dark matter halos from the Millennium Simulation (Springel et al. 2005) and GALFORM semi-analytic models (Cole et al. 2000; Bower et al. 2006). Average values from the current study are indicated by gray bands and predictions from Ascaso et al. (2015) study are indicated by red dotted lines.

We find that a comparison with only spectroscopic redshifts yields a fairly optimistic view of  $z < 1$  photometric redshifts in the LSST survey, predicting  $\sigma_{\text{NMAD}} \sim 0.02$  and very few outliers. However, when the full sample of grism redshifts is included in the test, the measured scatter increases significantly, spanning from  $\sim 2\%$  at  $z \sim 0.25$  to  $\sim 7\%$  at  $z \sim 2$ . Additionally, the number of catastrophic outliers increases dramatically with redshift from 5% to 30%. Ascaso et al. (2015) performed similar tests using mock catalogs and found slightly more optimistic results, with NMAD scatter spanning  $\sim 0.03$ – $0.04$  (red dashed lines in NMAD and outlier panels). Aside from the lowest redshift bin, where the COSMOS field is small and the grism adds little to the redshift determination, photometric redshift accuracy predicted from the mock catalogs is optimistic. Estimates from mock catalogs are lower by up to a factor of two compared to simulations leveraging real data.

## 6.2. DES and VHS Surveys

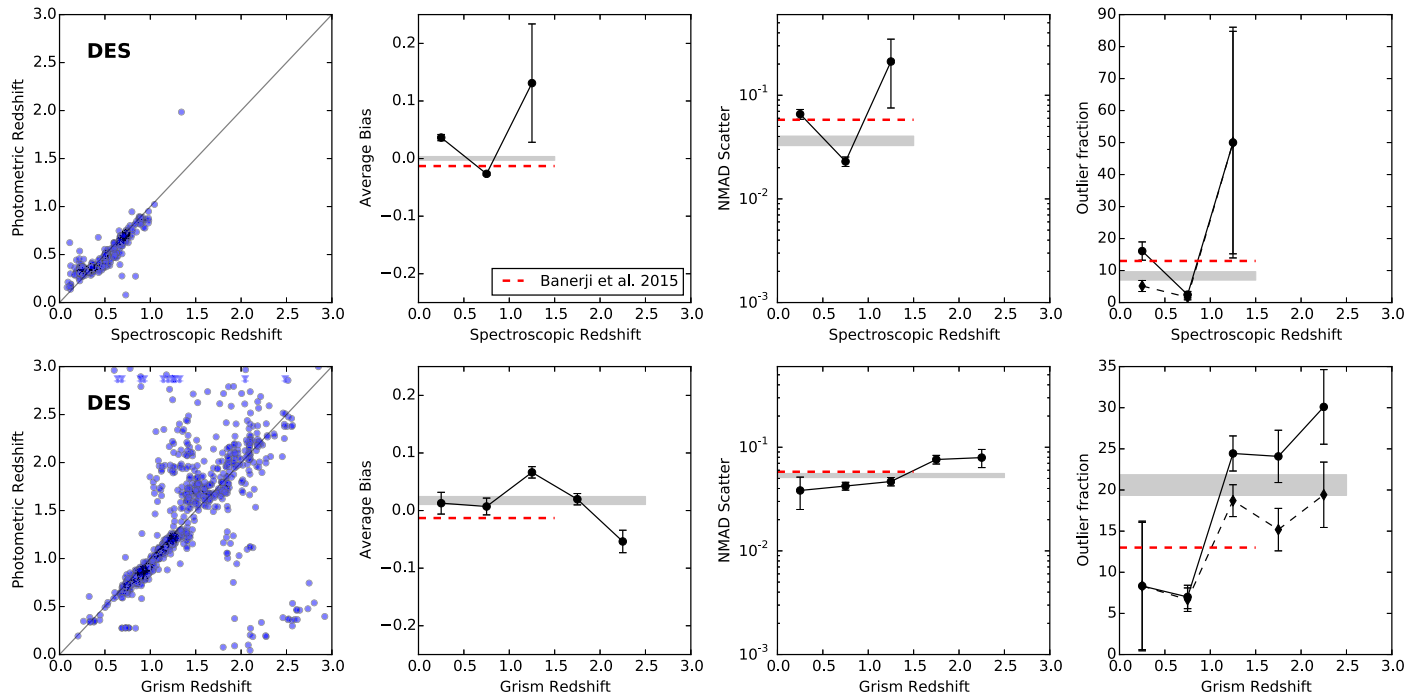
The DES survey is a photometric survey of  $5000 \text{ deg}^2$  of the southern sky that includes *grizY* imaging using the Dark Energy Camera (Flaugher 2005; Diehl & For Dark Energy Survey Collaboration 2012). According to the survey description document,<sup>11</sup> the target  $5\text{-}\sigma$  limiting magnitudes for point

sources in the DES survey will be 26.5, 26.0, 25.4, 24.7, and 23.0. As each of these limits is shallower than the imaging in COSMOS, we are able to accurately create a simulated catalog. Additionally, the DES footprint overlaps with the VISTA Hemisphere Survey (VHS), which can complement the DES optical data with near-IR  $JHK$  imaging over  $\sim 20,000 \text{ deg}^2$  down to limiting magnitudes of 21.5, 21.16, 20.3 (McMahon et al. 2013).

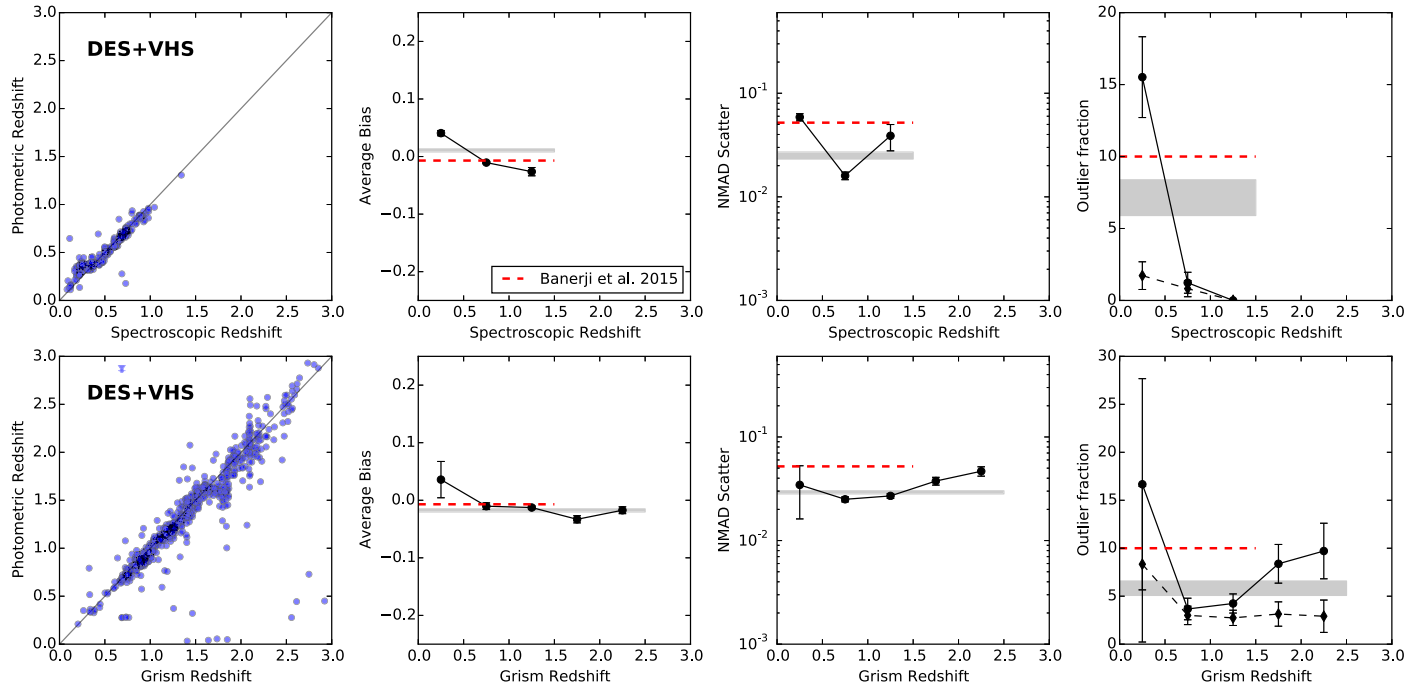
Figure 23 demonstrates the photometric redshift performance for the spectroscopic and grism samples in the simulated DES survey and Figure 24 includes both DES and VHS filters. For the DES filters alone, we find that the photometric redshift scatter will be higher than for the LSST ( $\sim 5\%$ ) and will increase to  $\sim 8\%$  by a redshift of 2.5. This is partially due to the omission of the  $u$  filter and shallower depths of the DES survey. Additionally, from the grism sample, we predict that the outlier fraction will be quite high with only DES imaging ( $\sim 20\%$ ). These estimates are somewhat different from those found in a study by Banerji et al. (2015) (red dashed line on Figures 23 and 24), which found very similar scatter, but lower outlier fractions. We note that Banerji et al. (2015) adopted a slightly different definition ( $\Delta z / (1+z) > 0.15$ ). We include outliers defined in this way as a black dashed line in Figures 23 and 24 and find closer agreement between the two studies.

We find that the addition of near-IR photometry with the VHS specifications improves the photometric redshift performance dramatically, in contrast with the findings of Banerji et al. (2015). When these data are included, the mean NMAD scatter decreases to  $\sim 3\%$ , with an increase to  $\sim 5\%$  above  $z \gtrsim 1.5$ . In a comparison with the grism redshifts, which will be less biased than the spectroscopic catalogs used by Banerji et al. (2015), we find a lower outlier fraction of  $\sim 6\%$  versus

<sup>11</sup> <http://www.darkenergysurvey.org/survey/des-description.pdf>



**Figure 23.** Photometric redshift performance in the simulated DES survey ( $g, r, i, z, y$  filters). Photometric redshift accuracy will depend strongly on redshift, with a significant outlier fraction ( $\geq 20\%$ ) except at  $0.5 < z < 1.0$ . These systematics will be improved significantly by including near-IR photometry, e.g., from the VHS survey, as shown in Figure 24. Red dashed lines indicate measured redshift performance with DES science verification data as estimated for a sample of bright objects with spectroscopic redshifts (Banerji et al. 2015).



**Figure 24.** Photometric redshift performance in the simulated combined DES and VHS surveys ( $ugrizyJHKs$  filters) from spectroscopic and grism samples. Errors in photometric redshifts will be  $\sim 3\%$  on average, ranging up to  $\sim 5\%$  at  $z \sim 2.5$  with  $\sim 3\%$ – $13\%$  outliers. Photometric redshift scatter and outlier fraction is lower in this data set than in Banerji et al. (2015), suggesting the importance of including near-IR photometry from the VHS survey for photometric redshift performance.

10%. This average value decreases even further with the less strict outlier threshold ( $\Delta z / (1+z) > 0.15$ )

We emphasize that simulations like those presented in this section are overly simplistic. The filters used in planned or ongoing surveys may not exactly match those used in the 3D-HST catalogs. Furthermore, although we attempt to match

quoted catalog depths, there will naturally be differences in image quality (e.g., seeing) and redshift fitting methodology that will influence photometric redshift performance. In particular, the COSMOS catalogs are *HST*-detected, therefore ground-based photometry will suffer more dramatically from blending. Furthermore, the redshift accuracy in the COSMOS

field is excellent when only the broadband optical and near-IR imaging is included (see e.g., Figure 13), and therefore these estimates could further be a generous estimate of photometric redshift performance in the planned surveys. We emphasize the discrepancies between the mock and empirical predictions presented in this section and suggest the importance of including empirical tests with representative spectroscopic samples in addition to mock simulations in order to robustly predict photometric redshift accuracy.

## 7. SUMMARY

Studies of the high-redshift universe are increasingly reliant on photometric redshifts to probe fainter targets in scope and variety than are inaccessible to even the most ambitious spectroscopic campaigns. The goal of this paper is to assess and quantify the photometric redshift accuracy in the 3D-HST photometric catalogs. We summarize the major findings below.

1. The 3D-HST photometric catalogs consist of PSF-matched aperture photometry across  $\sim 900$  square arcminutes in the CANDELS extragalactic field, including ground- and space-based imaging from 0.3 to 8.0  $\mu\text{m}$ . Overall, photometric redshift quality in the catalogs, calculated using EAZY (Brammer et al. 2008), is excellent, with an overall characteristic scatter of  $\Delta z / (1+z) \sim 0.02$  down to  $JH_{\text{IR}} = 24$ . This result is fairly robust to measurement technique, e.g., comparison with spectroscopic or grism redshifts versus galaxy pair counts, although it does vary among the five fields by  $\pm 0.006$ .
2. The characteristic, or NMAD, scatter does not depend strongly on galaxy stellar mass or  $U - V$  rest-frame color; however, we do find significant variations in the fraction of catastrophic outliers ( $\Delta z / (1+z) > 0.1$ ). Photometric redshift scatter increases by  $\sim 1\% - 2\%$  ( $1+z$ ) with apparent magnitude (down to the limiting magnitude of  $JH_{\text{IR}} = 24$ ) and redshift (out to  $z \sim 2.5 - 3$ ). Analysis of close pairs suggests that redshift accuracy further degrades for fainter objects, reaching  $\sim 0.046$  ( $1+z$ ) at  $H_{F160W} = 26$ .
3. Photometric redshift accuracy depends significantly on galaxy surface brightness; at fixed apparent magnitude, the photometric redshifts for the most extended galaxies can exhibit up to  $\sim 0.005$  more scatter and a  $\sim 5$  times greater fraction of outliers than the most compact galaxies.
4. We confirm that the error estimates and PDFs produced by the EAZY code are narrow with respect to the photometric redshift scatter, but this underestimation cannot be improved by uniformly broadening the PDFs. Furthermore, errors in photometric redshift estimates do not capture the outlier behavior. However, the effect on the derived overall properties of a sample may be subtle; the overall spectroscopic/grism and photometric redshift distributions as probed by single-valued estimates and full PDFs agree to within  $\sim 3\% - 4\%$ . In many specific cases, such as deriving luminosity or mass functions, scatter and outliers can tend to bias the derived properties. Although the size of this bias is not immediately calculable, it must be simulated for any given survey, magnitude limit, and redshift range.
5. Finally, a fraction of the field-to-field variation in photometric redshift quality can be attributed to the heterogenous nature of available imaging bands. We investigate the contribution of various filter combinations on the derived redshift accuracy, highlighting the dramatic impact driven by the inclusion of *Spitzer*-IRAC photometry, blue (*F435W*) *HST* photometry, and medium-band filters particularly in the optical.

The conclusions from this paper extend far beyond the use of the 3D-HST catalogs and can be applied in the interpretation of current surveys for which grism spectroscopy is not available. Furthermore, the systematics in redshift accuracy can be used in the planning of future surveys. To illustrate this possibility, we included simulations of photometric redshift performance in the LSST, DES, and DES plus VHS data sets in Section 6. This type of empirical simulation could more realistically reflect the input galaxy population than a spectroscopic or mock catalog, which could yield overly optimistic estimates for redshift accuracy.

Additionally, the demonstrated filter dependence can influence survey design choices. For example, the inclusion of blue (*F435W*) imaging in the GOODS fields significantly improved both the scatter and outlier fractions. One key question in planning photometric surveys is the balance between depth in broad filters and shallower imaging in narrower filters. The significant improvement in photometric redshift accuracy, especially in the outlier fraction, due to the inclusion of medium-band imaging for the current sample, centered at  $1 < z < 2$ , can inform future studies of the earlier universe. Similar medium-band imaging in the near-IR, as used by NMBS (Whitaker et al. 2011) and the FourStar Galaxy Evolution Survey (ZFOURGE; I. Labbé et al., in preparation) could be crucial for future studies at higher redshift to maximize confidence in redshift estimates for individual galaxies as opposed to their statistical properties.

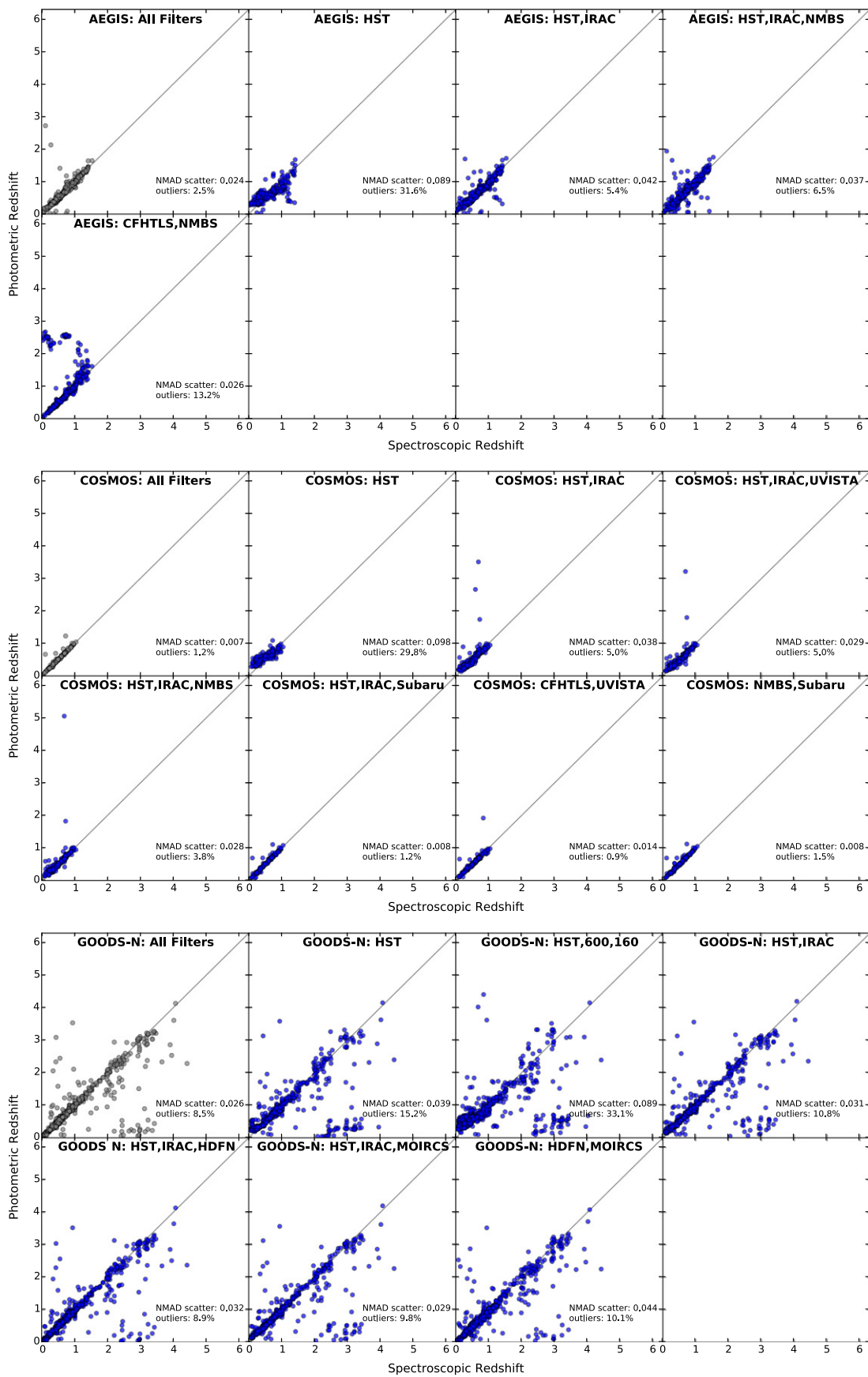
R.B. and K.E.W. gratefully acknowledge support by NASA through Hubble Fellowship grants #HF-51318 and #HF2-51368 awarded by the Space Telescope Science Institute, which is operated by the Association of Universities for Research in Astronomy, Inc., for NASA, under contract NAS5-26555. This research made use of Astropy, a community-developed core Python package for Astronomy (Astropy Collaboration et al. 2013). This work is based on observations taken by the 3D-HST Treasury Program (GO 12177 and 12328) with the NASA/ESA *HST*, which is operated by the Association of Universities for Research in Astronomy, Inc., under NASA contract NAS5-26555. Finally, this work is also based on observations taken by the CANDELS Multi-Cycle Treasury Program with the NASA/ESA *HST*, which is operated by the Association of Universities for Research in Astronomy, Inc., under NASA contract NAS5-26555.

## APPENDIX

### FIELD TO FIELD VARIATION IN PHOTOMETRIC REDSHIFT ACCURACY

For the most part, we have treated the 3D-HST catalogs as a uniform photometric sample. However, aside from the availability of optical and near-IR *HST* imaging and *Spitzer*-IRAC photometry, each field includes a heterogenous collection of photometry and spectroscopic redshifts. In this Appendix, we





**Figure 25.** Photometric vs. spectroscopic redshifts in the 3D-HST/CANDELS fields for different filter combinations. Fits to the full photometric catalogs are included as gray symbols; all other tests are included as blue symbols.

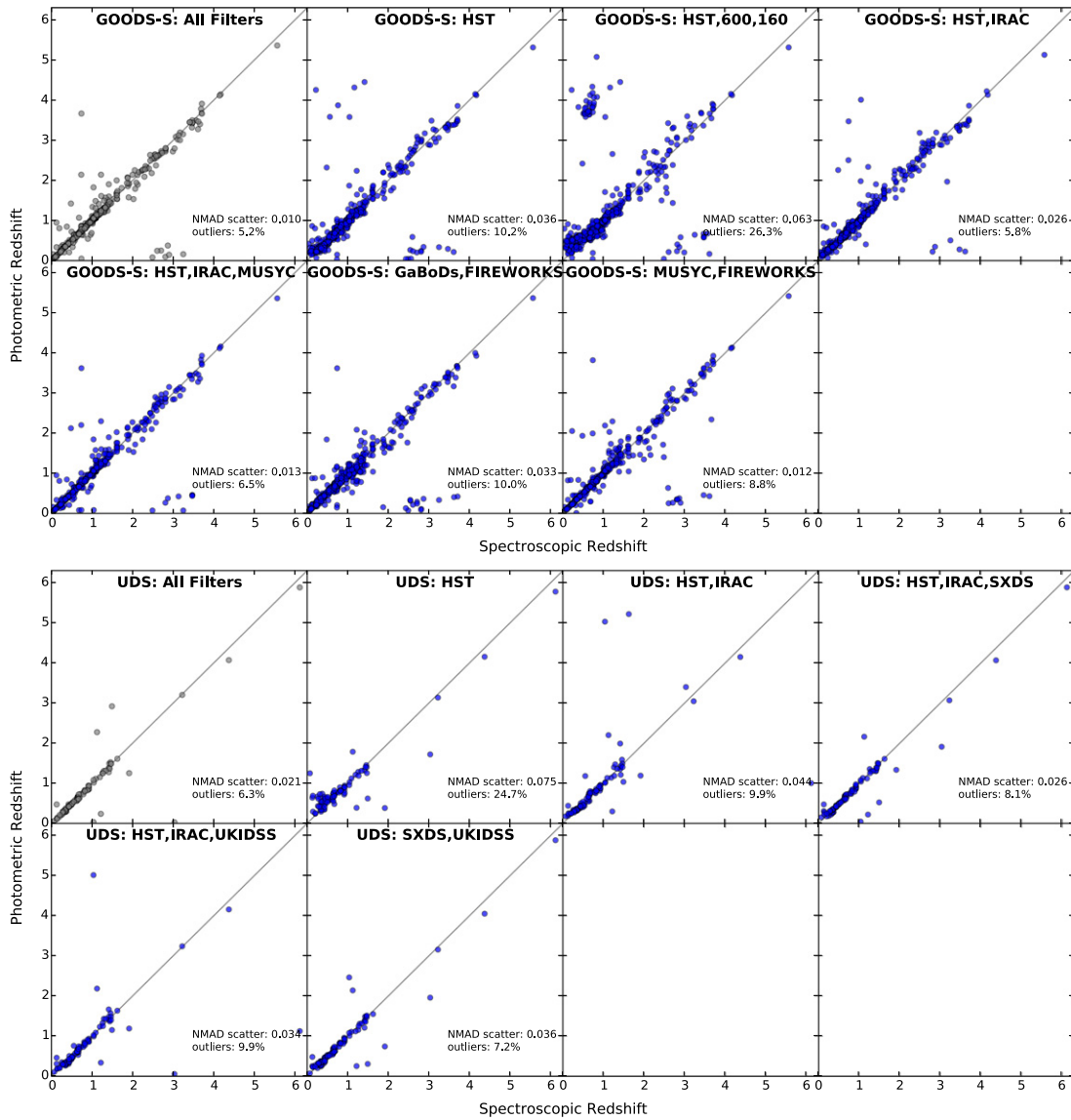


Figure 25. (Continued.)

show the scatter between photometric redshifts and true redshifts (grism and spectroscopic) for each field and subset of photometry (Figures 25 and 26). Panels in each figure are divided into fields, with the redshifts from the full photometric catalog included in the top left panel (gray points) and redshifts measured from subsets of filters in the additional panels (blue points). Figure 25 includes comparisons with spectroscopic redshifts and Figure 26 with grism redshifts. Measured scatter and outlier fractions are included in Table 4.

These figures illuminate some of the reasons for the strong field-to-field variance in photometric redshift scatter and outlier fractions shown in Figures 12 and 13. In Figure 25 it is apparent that the spectroscopic redshift follow-up varies wildly from field-to-field. For example, GOODS-N exhibits significantly more scatter than the other fields, however it also includes much better sampling of  $z > 1$  galaxies e.g., than COSMOS, which has an amazingly tight relationship between photometric and spectroscopic redshifts. On the other hand, GOODS-S also includes a large number of high-redshift galaxies, but much lower NMAD scatter.

Differences in spectroscopic data sets do not explain the full field-to-field variation as the variations persist in comparisons with grism redshifts (Figure 26), where the redshift coverage should be approximately uniform. Still COSMOS and GOODS-S, which both have medium-band optical imaging, have the most accurate photometric redshifts. These optical filters appear to have a greater effect than medium-band filters in the near-IR, which are included in the AEGIS and COSMOS fields from the NMBS Survey (Whitaker et al. 2011). We expect that this is due to the redshift distribution probed by the grism; at higher redshifts, such deep medium-band imaging should be increasingly important.

We find that many systematics are introduced by including various subsets of photometric bands in the redshift fitting. Although we have discussed some of these trends in Section 4, we include all tests in separate panels in these figures to illustrate some of the systematic redshift failures that are behind the increased scatter and catastrophic failure rates. For example, it is clear that IRAC photometry breaks an important degeneracy that systematically narrows the distribution of

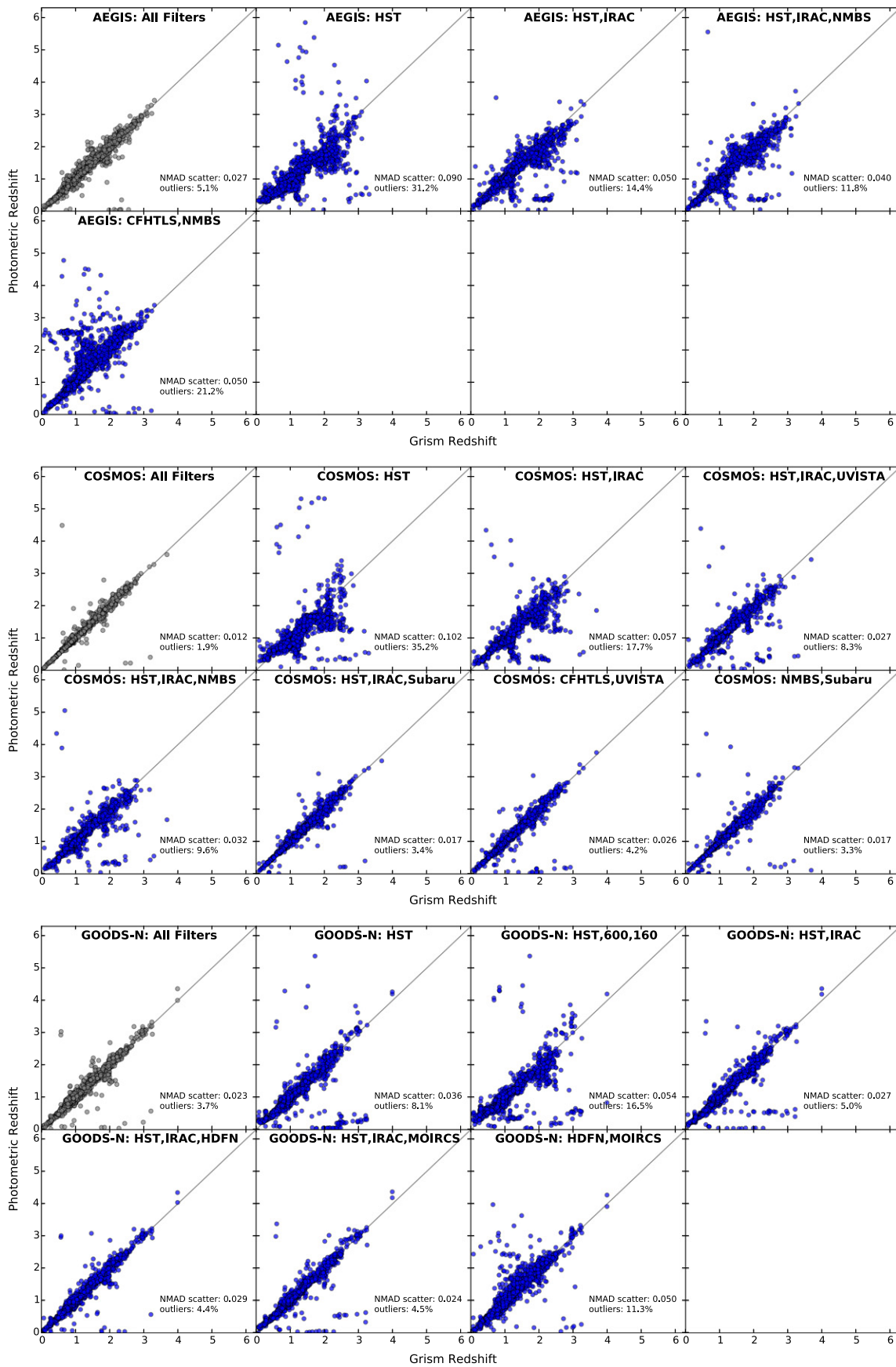


Figure 26. Photometric vs. grism redshifts in the 3D-HST/CANDELS fields for different filter combinations.

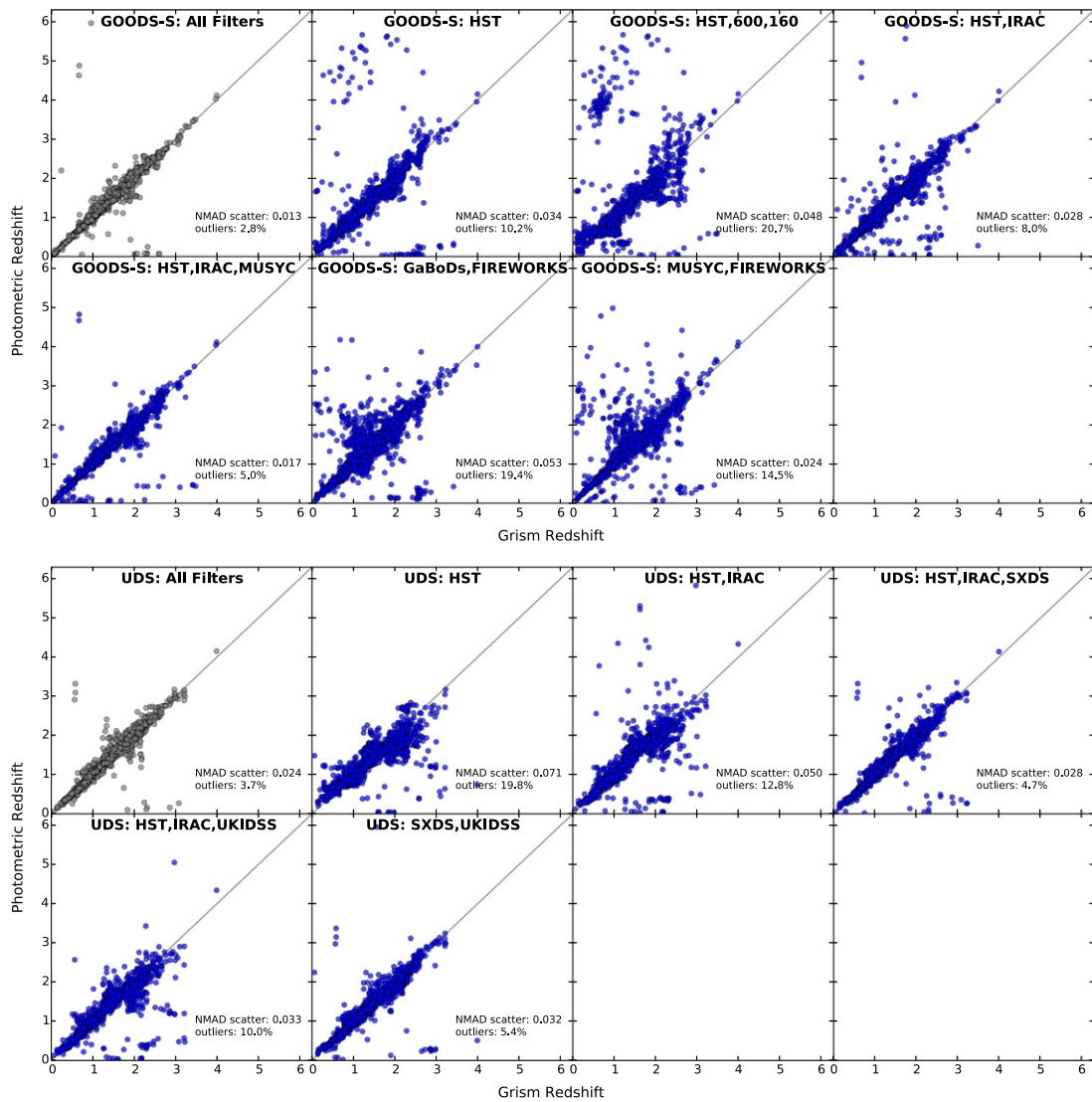


Figure 26. (Continued.)

photometric redshifts toward  $z \sim 0.7$ , improves the accuracy at  $z \sim 2$ , and discriminates between low-redshift ( $z < 1$ ) and very high-redshift galaxies ( $z \sim 4-5$ ).

## REFERENCES

- Abdalla, F. B., Banerji, M., Lahav, O., & Rashkov, V. 2011, *MNRAS*, 417, 1891
- Akiyama, M., Ueda, Y., Watson, M. G., et al. 2015, *PASJ*, 67, 82
- Aragon-Calvo, M. A., Weygaert, R. v. d., Jones, B. J. T., & Mobasher, B. 2015, *MNRAS*, 454, 463
- Ascaso, B., Mei, S., & Benítez, N. 2015, *MNRAS*, 453, 2515
- Ashby, M. L. N., Willner, S. P., Fazio, G. G., et al. 2013, *ApJ*, 769, 80
- Astropy Collaboration, Robitaille, T. P., Tollerud, E. J., et al. 2013, *A&A*, 558, A33
- Banerji, M., Jouvel, S., Lin, H., et al. 2015, *MNRAS*, 446, 2523
- Barger, A. J., Cowie, L. L., & Wang, W.-H. 2008, *ApJ*, 689, 687
- Barnby, P., Huang, J.-S., Ashby, M. L. N., et al. 2008, *ApJS*, 177, 431
- Benjamin, J., van Waerbeke, L., Ménard, B., & Kilbinger, M. 2010, *MNRAS*, 408, 1168
- Bezanson, R., Franx, M., & van Dokkum, P. G. 2015, *ApJ*, 799, 148
- Bezanson, R., van Dokkum, P. G., van de Sande, J., et al. 2013, *ApJL*, 779, L21
- Bielby, R., Hudelot, P., McCracken, H. J., et al. 2012, *A&A*, 545, A23
- Bower, R. G., Benson, A. J., Malbon, R., et al. 2006, *MNRAS*, 370, 645
- Brammer, G. B., van Dokkum, P. G., & Coppi, P. 2008, *ApJ*, 686, 1503
- Brammer, G. B., van Dokkum, P. G., Franx, M., et al. 2012, *ApJS*, 200, 13
- Brammer, G. B., Whitaker, K. E., van Dokkum, P. G., et al. 2011, *ApJ*, 739, 29
- Bruzual, G., & Charlot, S. 2003, *MNRAS*, 344, 1000
- Capak, P., Cowie, L. L., Hu, E. M., et al. 2004, *AJ*, 127, 180
- Cardamone, C. N., van Dokkum, P. G., Urry, C. M., et al. 2010, *ApJS*, 189, 270
- Chabrier, G. 2003, *PASP*, 115, 763
- Chen, H.-W., Marzke, R. O., McCarthy, P. J., et al. 2003, *ApJ*, 586, 745
- Cohen, J. G. 2001, *AJ*, 121, 2895
- Cohen, J. G., Hogg, D. W., Blandford, R., et al. 2000, *ApJ*, 538, 29
- Cole, S., Lacey, C. G., Baugh, C. M., & Frenk, C. S. 2000, *MNRAS*, 319, 168
- Cooper, M. C., Newman, J. A., Davis, M., Finkbeiner, D. P., & Gerke, B. F. 2012, *ascl:1203.003*
- Cowie, L. L., Barger, A. J., Hu, E. M., Capak, P., & Songaila, A. 2004, *AJ*, 127, 3137
- Dahlen, T., Mobasher, B., Faber, S. M., et al. 2013, *ApJ*, 775, 93
- Dawson, S., Stern, D., Bunker, A. J., Spinrad, H., & Dey, A. 2001, *AJ*, 122, 598
- Dickinson, M., Papovich, C., Ferguson, H. C., & Budavári, T. 2003, *ApJ*, 587, 25
- Diehl, T. & For Dark Energy Survey Collaboration 2012, *PhPro*, 37, 1332
- Dobos, L., Csabai, I., Yip, C.-W., et al. 2012, *MNRAS*, 420, 1217
- Erben, T., Hildebrandt, H., Lerchster, M., et al. 2009, *A&A*, 493, 1197
- Erben, T., Schirmer, M., Dietrich, J. P., et al. 2005, *AN*, 326, 432
- Fioc, M., & Rocca-Volmerange, B. 1997, *A&A*, 326, 950



- Flaugher, B. 2005, *IJMPA*, 20, 3121
- Fumagalli, M., Labbé, I., Patel, S. G., et al. 2014, *ApJ*, 796, 35
- Furusawa, H., Kosugi, G., Akiyama, M., et al. 2008, *ApJS*, 176, 1
- Geach, J. E., Simpson, C., Rawlings, S., Read, A. M., & Watson, M. 2007, *MNRAS*, 381, 1369
- Gialalisco, M., Ferguson, H. C., Koekemoer, A. M., et al. 2004, *ApJL*, 600, L93
- Grogin, N. A., Kocevski, D. D., Faber, S. M., et al. 2011, *ApJS*, 197, 35
- Hildebrandt, H., Arnouts, S., Capak, P., et al. 2010, *A&A*, 523, A31
- Hildebrandt, H., Erben, T., Dietrich, J. P., et al. 2006, *A&A*, 452, 1121
- Hildebrandt, H., Pielorz, J., Erben, T., et al. 2009, *A&A*, 498, 725
- Hildebrandt, H., Wolf, C., & Benítez, N. 2008, *A&A*, 480, 703
- Hogg, D. W., Cohen, J. G., Blandford, R., et al. 1998, *AJ*, 115, 1418
- Hsieh, B.-C., Wang, W.-H., Hsieh, C.-C., et al. 2012, *ApJS*, 203, 23
- Ivezic, Z., Tyson, J. A., Abel, B., et al. 2008, arXiv:0805.2366
- Kajisawa, M., Ichikawa, T., Tanaka, I., et al. 2011, *PASJ*, 63, 379
- Kajisawa, M., Ichikawa, T., Yamada, T., et al. 2010, *ApJ*, 723, 129
- Koekemoer, A. M., Faber, S. M., Ferguson, H. C., et al. 2011, *ApJS*, 197, 36
- Kriek, M., van Dokkum, P. G., Labbé, I., et al. 2009, *ApJ*, 700, 221
- Labbé, I., Bouwens, R., Illingworth, G. D., & Franx, M. 2006, *ApJL*, 649, L67
- Lilly, S. J., Le Fèvre, O., Renzini, A., et al. 2007, *ApJS*, 172, 70
- Marchesini, D., Muzzin, A., Stefanon, M., et al. 2014, *ApJ*, 794, 65
- Marchesini, D., van Dokkum, P. G., Förster Schreiber, N. M., et al. 2009, *ApJ*, 701, 1765
- Marchesini, D., Whitaker, K. E., Brammer, G., et al. 2010, *ApJ*, 725, 1277
- McCracken, H. J., Milvang-Jensen, B., Dunlop, J., et al. 2012, *A&A*, 544, A156
- McCracken, H. J., Wolk, M., Colombi, S., et al. 2015, *MNRAS*, 449, 901
- McMahon, R. G., Banerji, M., Gonzalez, E., et al. 2013, *Msngr*, 154, 35
- Momcheva, I. G., Brammer, G. B., van Dokkum, P. G., et al. 2015, arXiv:1510.02106
- Muzzin, A., Marchesini, D., Stefanon, M., et al. 2013, *ApJ*, 777, 18
- Myers, A. D., White, M., & Ball, N. M. 2009, *MNRAS*, 399, 2279
- Newman, J. A. 2008, *ApJ*, 684, 88
- Newman, J. A., Cooper, M. C., Davis, M., et al. 2013, *ApJS*, 208, 5
- Nonino, M., Dickinson, M., Rosati, P., et al. 2009, *ApJS*, 183, 244
- Ono, Y., Ouchi, M., Shimasaku, K., et al. 2010, *MNRAS*, 402, 1580
- Ouchi, M., Shimasaku, K., Akiyama, M., et al. 2008, *ApJS*, 176, 301
- Papovich, C., Momcheva, I., Willmer, C. N. A., et al. 2010, *ApJ*, 716, 1503
- Quadri, R. F., & Williams, R. J. 2010, *ApJ*, 725, 794
- Quadri, R. F., Williams, R. J., Lee, K.-S., et al. 2008, *ApJL*, 685, L1
- Reddy, N. A., Steidel, C. C., Erb, D. K., Shapley, A. E., & Pettini, M. 2006, *ApJ*, 653, 1004
- Retzlaff, J., Rosati, P., Dickinson, M., et al. 2010, *A&A*, 511, A50
- Sanders, D. B., Salvato, M., Aussel, H., et al. 2007, *ApJS*, 172, 86
- Shen, S., Mo, H. J., White, S. D. M., et al. 2003, *MNRAS*, 343, 978
- Simpson, C., Almaini, O., Cirasuolo, M., et al. 2006, *MNRAS*, 373, L21
- Simpson, C., Rawlings, S., Ivison, R., et al. 2012, *MNRAS*, 421, 3060
- Skelton, R. E., Whitaker, K. E., Momcheva, I. G., et al. 2014, *ApJS*, 214, 24
- Smail, I., Sharp, R., Swinbank, A. M., et al. 2008, *MNRAS*, 389, 407
- Sołtan, A. M., & Chodorowski, M. J. 2015, *MNRAS*, 453, 1013
- Springel, V., White, S. D. M., Jenkins, A., et al. 2005, *Natur*, 435, 629
- Steidel, C. C., Adelberger, K. L., Shapley, A. E., et al. 2003, *ApJ*, 592, 728
- Taniguchi, Y., Scoville, N., Murayama, T., et al. 2007, *ApJS*, 172, 9
- Treu, T., Ellis, R. S., Liao, T. X., & van Dokkum, P. G. 2005, *ApJL*, 622, L5
- van Breukelen, C., Cotter, G., Rawlings, S., et al. 2007, *MNRAS*, 382, 971
- van de Sande, J., Kriek, M., Franx, M., et al. 2013, *ApJ*, 771, 85
- van der Wel, A., Franx, M., van Dokkum, P. G., et al. 2014, *ApJ*, 788, 28
- Wake, D. A., Whitaker, K. E., Labbé, I., et al. 2011, *ApJ*, 728, 46
- Whitaker, K. E., Kriek, M., van Dokkum, P. G., et al. 2012, *ApJ*, 745, 179
- Whitaker, K. E., Labbé, I., van Dokkum, P. G., et al. 2011, *ApJ*, 735, 86
- Wirth, G. D., Willmer, C. N. A., Amico, P., et al. 2004, *AJ*, 127, 3121
- Wuyts, S., Labbé, I., Franx, M., et al. 2007, *ApJ*, 655, 51
- Wuyts, S., Labbé, I., Schreiber, N. M. F., et al. 2008, *ApJ*, 682, 985
- Yamada, T., Kodama, T., Akiyama, M., et al. 2005, *ApJ*, 634, 861
- Yoshikawa, T., Akiyama, M., Kajisawa, M., et al. 2010, *ApJ*, 718, 112