

Simulating JWST/NIRCam Color Selection of High-Redshift Galaxies

KEVIN N. HAINLINE,¹ RAPHAEL E. HVIDING,¹ MARCIA RIEKE,¹ IRENE SHIVAEI,² RYAN ENDSLEY,¹ EMMA CURTIS-LAKE,^{3,4} RENSKE SMIT,⁴ CHRISTINA C. WILLIAMS,^{1,*} STACEY ALBERTS,¹ KRISTAN N K BOYETT,⁵ ANDREW J. BUNKER,⁵ EIICHI EGAMI,¹ MICHAEL V. MASEDA,⁶ SANDRO TACCHELLA,⁷ AND CHRISTOPHER N.A. WILLMER¹

¹Steward Observatory, University of Arizona, 933 North Cherry Avenue, Tucson, AZ 85721, USA

²Hubble Fellow, Steward Observatory, University of Arizona, 933 North Cherry Avenue, Tucson, AZ 85721, USA

³Kavli Institute for Cosmology, Madingley Road, Cambridge CB3 0HA, UK

⁴Cavendish Laboratory, University of Cambridge, 19 JJ Thomson Avenue, Cambridge CB3 0HE, UK

⁵Department of Physics, University of Oxford, Denys Wilkinson Building, Keble Road, OX1 3RH, UK

⁶Leiden Observatory, Leiden University, P.O. Box 9513, NL-2300 RA Leiden, the Netherlands

⁷Center for Astrophysics, Harvard & Smithsonian, 60 Garden Street, Cambridge, MA 02138, USA

(Received 11/25/19; Revised 2/28/20; Accepted 3/7/20)

ABSTRACT

The NIRCam instrument on the upcoming *James Webb Space Telescope* (JWST) will offer an unprecedented view of the most distant galaxies. In preparation for future deep NIRCam extragalactic surveys, it is crucial to understand the color selection of high-redshift galaxies using the Lyman dropout technique. To that end, we have used the JADES extraGalactic UltraDeep Artificial Realizations (JAGUAR) mock catalog to simulate a series of extragalactic surveys with realistic noise estimates. This enables us to explore different color selections and their impact on the number density of recovered high-redshift galaxies and lower-redshift interlopers. We explore how survey depth, detection signal-to-noise ratio, color selection method, detection filter choice, and the presence of the Ly α emission line affects the resulting dropout selected samples. We find that redder selection colors reduce the number of recovered high-redshift galaxies, but the overall accuracy of the final sample is higher. In addition, we find that methods that utilize two or three color cuts have higher accuracy because of their ability to select against low-redshift quiescent and faint dusty interloper galaxies. We also explore the near-IR colors of brown dwarfs and demonstrate that, while they are predicted to have low on-sky densities, they are most likely to be recovered in F090W dropout selection, but there are color cuts which help to mitigate this contamination. Overall, our results provide NIRCam selection methods to aid in the creation of large, pure samples of ultra high-redshift galaxies from photometry alone.

Keywords: galaxies: distances and redshifts – galaxies: high-redshift

1. INTRODUCTION

The discovery and characterization of high-redshift ($z > 6$) galaxies offers fundamental insights into galaxy assembly and star formation, including the creation of dust and metals, in the first billion years of the history of the universe. Deep imaging with the Wide Field Camera 3 (WFC3) instrument on board the *Hubble Space Telescope* (HST) has revealed samples of galaxies at these redshifts (Bouwens et al. 2003, 2004, 2007, 2008; Bunker et al. 2004, 2010; McLure et al. 2010; Wilkins et al. 2011; Lorenzoni et al. 2011, 2013), including an intriguing, if limited, population of ultra-high redshift galaxies at $z > 10$ (Oesch et al. 2014, 2015a, 2018; Zitrin et al. 2014; Infante et al. 2015; Ishigaki et al. 2015; McLeod et al. 2016; Salmon et al. 2018). Assembling larger populations of galaxies at higher redshifts is challenging due

to the lack of infrared coverage of the instruments on HST (the longest wavelength filter on WFC3 is at 1.6 μm), the limited sensitivity and low resolution of observations made at longer wavelengths by the *Spitzer Space Telescope*, and infrared atmospheric absorption for ground-based observations. Overcoming these limitations is fundamental for understanding the evolution of the earliest galaxies (see reviews by Dunlop 2013; Stark 2016).

The selection of high-redshift galaxies is crucial for our understanding of reionization, where the neutral hydrogen that filled the universe after recombination was ionized in a process thought to be driven by early star-forming galaxies between $z \sim 6-10$ (Robertson et al. 2015), although accretion onto supermassive black holes is also thought to be a contributing factor (Giallongo et al. 2015; Madau & Haardt 2015; Onoue et al. 2017). By characterizing the galaxies that comprise the faint end of the UV luminosity function, the exact source and timescale of reionization can be understood. In addition, observations of these galaxies give us insight

* NSF Fellow

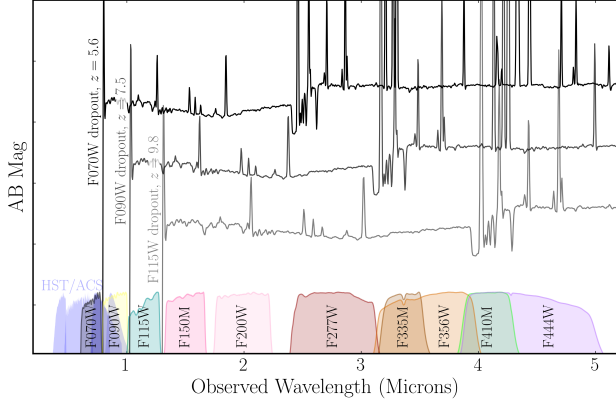


Figure 1. Example JAGUAR mock galaxy SEDs for F070W ($z = 5.6$, black), F090W ($z = 7.5$, grey), and F115W ($z = 9.8$, light grey) dropout galaxies. We also plot the HST/ACS bands we use in this work in light blue, and the *JWST*/NIRCam filters in multiple colors as labelled.

into the evolution of the star-formation rate density in the early universe, which has been observed to increase by almost an order of magnitude in the 170 million years between $8 < z < 10$ (Oesch et al. 2014, 2018; Ishigaki et al. 2018), although observations by McLeod et al. (2016) indicate a shallower evolution during this period. This tension may be due to cosmic variance and small sample sizes, providing a clear impetus to uncover larger samples of ultra-high-redshift galaxies.

A widely used method for selecting galaxies at high redshift involves searching for their redshifted Lyman break, a feature in their spectrum caused by the absorption of extreme ultraviolet radiation by neutral hydrogen in the intergalactic medium along the line of sight and surrounding a given galaxy. In this technique, a galaxy observed in a filter that probes a wavelength range bluewards of the Lyman break will have reduced flux compared to a filter that lies to the red of the break. As a result, by selecting for galaxies with extreme red colors in adjacent bands, a rough estimate of the redshift of the galaxy can be obtained (Guhathakurta et al. 1990). Galaxies selected in this way are referred to as “dropouts.” This method was used to assemble a large sample of galaxies at $z = 2 - 4$ using ground-based observations in the optical U , G , and R filters, which was subsequently observed spectroscopically to confirm individual galaxy redshifts (Steidel et al. 1996, 1999, 2003). This technique has subsequently been supported with spectroscopic observations of galaxies out to $z \sim 8$ (Bunker et al. 2003; Stanway et al. 2004; Vanzella et al. 2009, 2011; Stark et al. 2010; Ono et al. 2012; Schenker et al. 2012; Shibuya et al. 2012; Cassata et al. 2015; Oesch et al. 2015b; Roberts-Borsani et al. 2016; Song et al. 2016; Tasca et al. 2017).

An alternate method for estimating accurate photometric redshifts relies on modeling a galaxy’s full spectral energy distribution (SED). The use of this method requires additional observed photometry over what is often needed for Ly-

man dropout selection, as well as a diverse suite of observed galaxy templates or stellar population synthesis models. In addition, it is less straightforward to understand the sample selection and survey completeness for SED modeling techniques than for color selection methods, and color selection is significantly quicker than full template fitting. For these reasons, in this paper we will focus on dropout selection of high-redshift galaxies.

The near-infrared wavelength coverage of *HST* and *Spitzer* has been used to select dropout galaxies out to the current redshift frontier of $z = 9 - 11$ (Ellis et al. 2013; Oesch et al. 2013; McLure et al. 2013; Bouwens et al. 2015). At higher redshifts, the Lyman break is shifted further into the infrared, and this technique is therefore limited by the lack of HST WFC3 filters at wavelengths longer than $1.6 \mu\text{m}$. The infrared wavelength coverage and sensitivity of the Near Infrared Camera (NIRCam) instrument on the James Webb Space Telescope (*JWST*, Gardner et al. 2006) will enable the discovery of galaxies out to $z > 15$. Following the projected launch of *JWST* in 2021, NIRCam will provide $0.7 \mu\text{m}$ to $5 \mu\text{m}$ imaging over a 9.7 arcmin^2 field of view at resolutions of $0''.04 - 0''.1$. NIRCam offers excellent sensitivity in this wavelength range, with 10σ point source depths of 28 magnitude (AB) achievable in only 2 ksec at $2 \mu\text{m}$. As *JWST* is designed for only a nominal 5 - 10 year mission, it is imperative that we explore the ways in which NIRCam observations can be quickly and efficiently leveraged to assemble large samples of high-redshift galaxies.

To that end, in this study we use a catalog of mock galaxies to explore the relationship between various color selection methods and the properties of recovered high-redshift dropout galaxies. We use the JADES extraGalactic Ultra-deep Artificial Realizations (JAGUAR) mock catalog (Williams et al. 2018), which was developed by members of the joint NIRCam and NIRSpec Guaranteed Time Observation (GTO) teams to aid in preparing for the early observations that will be made with *JWST*, with a focus on the *JWST* Deep Extragalactic Survey (JADES) GTO program. JAGUAR offers a catalog of photometry and spectra for mock galaxies along with self-consistent modeling of strong UV and optical emission lines, and was created using current observations of the number counts of galaxies as a function of UV luminosity and mass. To prepare for future deep *JWST*/NIRCam surveys, we simulate NIRCam noise at various observational depths to explore how color cuts affect the number densities, redshift distributions, and intrinsic properties of recovered mock galaxies. We explore dropout selection using both *JWST*/NIRCam filters alone as well as selection with NIRCam + *HST*/ACS filters which are helpful for imaging below the Lyman break and rejecting low-redshift interlopers. The goal of this present study is not to provide canonical color cuts, but rather to demonstrate the types of color cut selection scenarios that can be employed to assemble galaxy samples at multiple redshift ranges.

We begin by outlining the creation of our photometric catalogs with realistic noise properties in Section 2. There we provide an overview of JAGUAR, describe the methods by

which we generate estimates of NIRC*am* noise, and discuss the overall design of the surveys we used to explore NIRC*am* color space. We outline our results in Section 3 for both simple single color selection and more complex, multi-color, selection. Additionally, we look at how dropout galaxy recovery is impacted by filter choice with an eye towards designing surveys that best utilize the unique dichroic beam splitter on board NIRC*am*. We also explore other common statistics used to separate interlopers, and the NIRC*am* colors of brown dwarfs. Finally, we discuss these results in Section 5, and conclude in Section 6. Throughout we adopt a cosmology with $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$, $\Omega_M = 0.3$, and $\Omega_\Lambda = 0.7$. All magnitudes are presented in the AB system (Oke & Gunn 1983).

2. METHODS

To explore the impact that color selection choices can make on recovered galaxy samples, we require a mock catalog that is diverse in star formation properties, redshifts, stellar masses, and dust attenuation as well as simulated observational noise at multiple depths. We also depend on statistical measures of how successful a given set of color cuts is at recovering high-redshift galaxy samples. In this section, we outline the JAGUAR catalog and describe our method for adding photometric noise to the JAGUAR photometry to produce mock observational catalogs at different simulated exposure times. We then describe the figures of merit we will use to compare the results from changing color selection methods, and finally, we discuss how we use these noisy photometric data to explore the NIRC*am* color space.

2.1. The JAGUAR Catalog

The JAGUAR mock catalog consists of a series of 11' by 11' photometric and spectroscopic catalogs, as described in Williams et al. (2018). JAGUAR includes both quiescent and star-forming mock galaxies using as the base catalog the observations of the galaxy stellar mass function from Tomczak et al. (2014) at $z < 4$ and the UV luminosity function from Bouwens et al. (2015) and Oesch et al. (2018) at $z > 4$. These mass and luminosity functions are joined at $z = 4$ by modeling the evolution of the relationship between observed galaxy stellar mass and M_{UV} , the absolute magnitude of each galaxy in the ultraviolet, in agreement with measurements in the 3D-HST survey (Skelton et al. 2014). JAGUAR mock galaxies were generated such that they followed the evolution of the mass and luminosity functions, and each object was then assigned a spectrum using BEAGLE, a tool designed to model and interpret galaxy SEDs (Chevallard & Charlot 2016). This code allows for the creation of realistic mock galaxy SEDs with self-consistent nebular continuum and line emission. A large quantity of BEAGLE galaxy realizations was constructed across a wide parameter space, including fits to existing 3D-HST objects, and each galaxy in the mock catalog was matched to an individual SED from these realizations. For each object, simple Sérsic profiles were assigned following observations of high-redshift galaxies in van der Wel et al. (2014), which have been shown to agree with low-

redshift results from Sloan Digital Sky Survey observations (Shen et al. 2003; Guo et al. 2009). We plot an example JAGUAR F070W dropout (at $z = 5.6$), F090W dropout (at $z = 7.5$), and F115W dropout (at $z = 9.8$) with the HST/ACS and JWST/NIRC*am* filters in Figure 1.

The JAGUAR catalogs span a stellar mass range of $\log(M_*/M_\odot) = 6 - 12$ and a redshift range of $z = 0.2 - 15$. For the lowest mass mock galaxies, the catalog requires significant extrapolation of existing mass and luminosity functions. We refer the reader to Williams et al. (2018) for a description of how the JAGUAR catalog agrees with current observations of the evolution of quiescent and star-forming galaxy properties, the cosmic star formation rate density, specific star-formation rate, and mass-metallicity relationship. The effects of IGM absorption in JAGUAR mock galaxies follow the prescription from Inoue et al. (2014). Dust attenuation of both the stars and the photoionized gas in the JAGUAR mock galaxies is described using a two-component model of Charlot & Fall (2000) and parameterized using $\hat{\tau}_V$, the total attenuation optical depth which is allowed to vary between 0 and 4, and the fraction of attenuation arising in the diffuse ISM μ , which is fixed at 0.4. While this range is motivated from observational relations (Schaerer & de Barros 2010), current samples of high-redshift galaxies that form the basis for these relations are likely missing a population of extremely dusty star-forming galaxies which may be observed with NIRC*am* (Casey et al. 2014; Spilker et al. 2016; Williams et al. 2019; Wang et al. 2019). We further discuss these sources in Section 3.8.

We plot the NIRC*am* color space for a 10'x10' JAGUAR realization in Figure 2, with mock galaxy points colored by their redshift. As can be seen in each panel, at specific redshifts where the filters span the Lyman break (plotted on the y-axis), the mock galaxies are observed to have redder colors. In each color selection scenario, there are also lower redshift interlopers with red colors, a mixture of those with strong 4000Å +Balmer breaks, star-forming galaxies with heavy dust obscuration, and quiescent galaxies. As would be expected with the evolution of the galaxy luminosity function to higher redshifts, the density of high-redshift dropout candidates decreases from F070W dropouts to F115W dropouts. We also overlay an example two-color dropout selection box in each panel to illustrate how objects lying inside the region at the top-left of each panel could be selected as dropout candidates.

2.2. Generating NIRC*am* Noise Estimates

While such plots as Figure 2 can be very helpful for choosing color criteria for selecting galaxies at specific redshift ranges, these plots do not incorporate any noise, which will preferentially affect fainter (and often lower mass) galaxies, moving them both into and out of color selection regions. To simulate noise we wrote a suite of custom scripts for use with the JAGUAR catalog, NIRCPrepareMock¹. These

¹ <https://github.com/kevinhainline/NIRCPrepareMock>

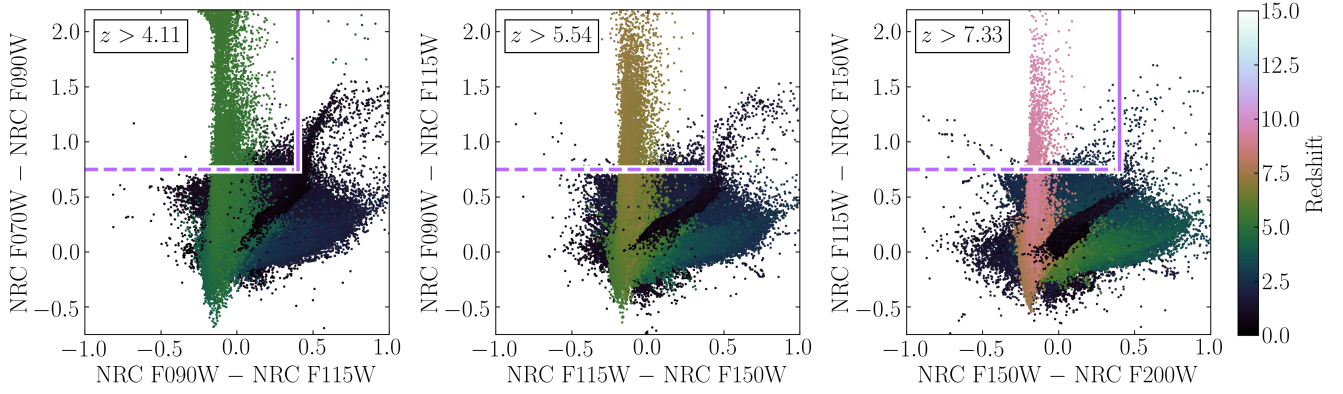


Figure 2. NIRC color-color plots for a $10' \times 10'$ section of the JAGUAR catalog, with mock galaxies at $z = 0.2 - 15$, without adding noise, with points colored by catalog redshift values, as given by the colorbar on the right side of the figure. The left panel shows F070W dropouts at $z \sim 5.5$, the center panel shows F090W dropouts at $z \sim 7.3$, and the right panel shows F115W dropouts at $z \sim 9.7$. In addition, in each panel, populations of lower redshift mock galaxies have red colors on both axes, and most selection methods at these redshifts will deliberately exclude these objects. We plot an example two-color selection method in lavender. In the absence of photometric noise, these selection boxes would return relatively pure samples of mock galaxies above a given redshift limit.

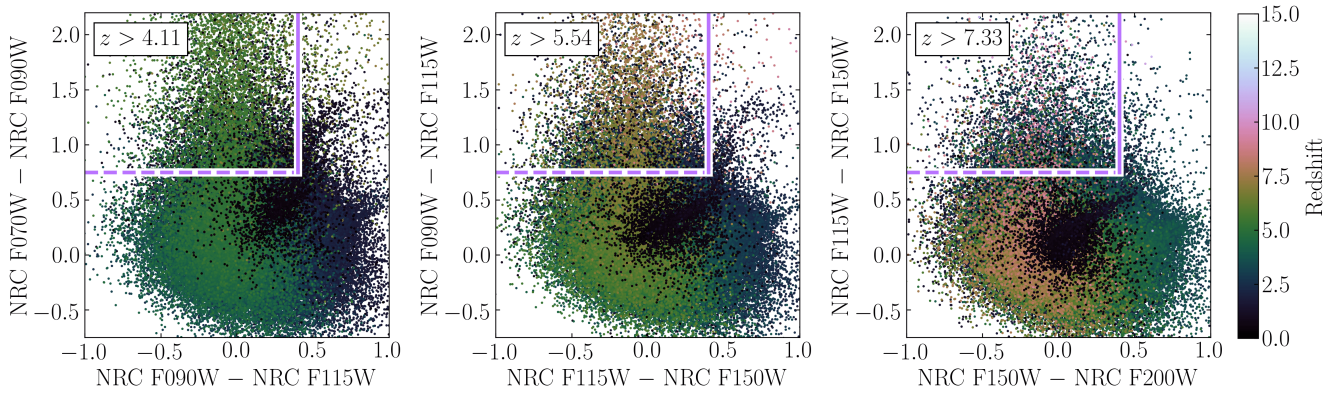


Figure 3. NIRC color-color plots for a $10' \times 10'$ section of the JAGUAR catalog, with mock galaxies at $z = 0.2 - 15$, created by simulating noise from images with a total of 98.9 ksec exposure time in each filter. In each panel, we only plot mock galaxies with detections in filters at wavelengths longer than the Lyman break with a $\text{SNR} > 3$. We show an example two-color selection method used throughout this work in lavender. In Section 3, we will discuss the properties of mock galaxies selected using this selection method, where we fix the Lyman break cut (dashed line) and vary the UV continuum cut (solid line).

scripts generate artificial noise directly from the JAGUAR photometry, which can be used when assessing the efficacy of photometric redshift or SED fitting codes.

We estimate noise for the mock galaxies in each filter separately, starting with the JAGUAR flux in that filter, as well as the morphology of the mock galaxy. The code selects the smallest circular aperture from a series of fixed radii ($0''.16$, $0''.24$, $0''.32$, and $0''.64$) that would encompass the semi-major axis half-light radius of each mock galaxy. At this point, the script calculates the total flux of each mock galaxy through that circular aperture taking into account its Sérsic index. Because we are not extracting flux from mock PSF-convolved NIRC images and extracting fluxes directly, we do not correct for aperture losses. To simulate the

sky background, we use estimates for the zodiacal light emission in the GOODS-S region for each filter², and add this to the flux of each mock galaxy through the aperture to produce the final flux in a given exposure. The uncertainty on the flux for an individual exposure is the Poisson noise summed in quadrature with the instrument read noise (summed over the pixels in the aperture).

When using NIRC, individual frames will be co-added to create a final deep image from which flux will be measured. To simulate this process, the code co-adds exposures

² following the *JWST* background model described here: <https://jwst-docs.stsci.edu/observatory-functionality/jwst-background-model>

(with a base exposure time), while allowing the user to input the total number of summed frames for a final exposure. To estimate the noise, we randomly sample from a Gaussian with the RMS width set to the exposure noise to produce the frame uncertainty, and then sum the noise in quadrature for each co-added frame.

There are some caveats to this approach to estimating uncertainty. Because we are not using mock images, we do not account for the instrumental point spread function, or change the size of the aperture in different filters to account for the change in instrument resolution as a function of wavelength. Accounting for the PSF would serve to decrease the flux that falls inside a given aperture at longer wavelengths, as the 50% encircled energy (defined as the fraction of light contained inside a circular aperture) increases from $0''.068$ for the F070W filter to $0''.092$ for the F444W filter. This effect would serve to make objects artificially more blue when using longer wavelength filters, although PSF-matching can account for this effect. For the majority of the analysis presented here, we focus on the NIRCcam short-wavelength filters (F070W, F090W, F115W, F150W, and F200W), where this effect is minimized. In addition, we don't simulate projection effects which would serve to reduce the number of high-redshift galaxies which are blocked by foreground galaxies. The full treatment of estimating noises from mock images is beyond the scope of the current work. While there are more rigorous methods of measuring fluxes, the flux uncertainties produced by our code agree with the predicted uncertainties from the STScI *JWST* Exposure Time Calculator³, and our code can be run quickly on large samples.

In Figure 3, we plot the same NIRCcam color spaces as in Figure 2, but with fluxes from a 100 square arcminute noisy catalog with images at 49.5 ksec total exposure time. We only plot mock galaxies detected with $\text{SNR} > 3$ in the filters on the x-axis in each panel (we do not set a limit on the SNR for the dropout filter). By comparing the noise-free to the noisy photometry, we can observe how rare dropout candidates are at $z > 8$, even in a 100 square arcminute field, both because of the faint observed fluxes of these objects (less mock galaxies satisfy the $\text{SNR} > 3$ criterion), as well as their low on-sky density based on the observed UV luminosity functions used to constrain JAGUAR. We also plot the same selection boxes as in Figure 2, demonstrating the difficulty in separating high redshift targets and lower redshift interlopers with noisy photometry.

2.3. Mock Survey Design

To explore high-redshift dropout selection with NIRCcam and *HST*+NIRCcam, we generated multiple sets of mock catalogs with realistic noise estimates. Because NIRCcam may target regions of the sky that do not have adequate deep *HST* coverage, we produced noisy data sets with only NIRCcam coverage over a region of 100 square arcminutes, with three different depths. In each case, we simulated a

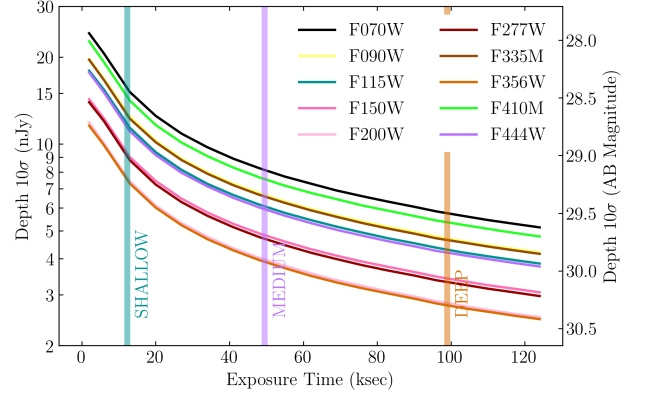


Figure 4. Simulated 10σ depths plotted against total exposure time for the NIRCcam filters used in this work. These values were estimated using the NIRCPrepareMock package. We also plot the exposure times for the SHALLOW (teal), MEDIUM (lavender), and DEEP (orange) surveys with vertical lines.

JWST/NIRCcam observational strategy for observing high-redshift galaxies which utilizes the DEEP8 readout pattern, with 7 groups per integration, for a pixel integration time of 1374.3 seconds. For each depth we assumed a 9-point dither pattern, which samples 3 times the pixel resolution, and we then varied the number of integrations per exposure:

1. A "SHALLOW" mock survey with 1 integration per exposure resulting in an integrated exposure time of 12.3 ksec per filter.
2. A "MEDIUM" mock survey with 4 integration per exposure resulting in an integrated exposure time of 49.5 ksec per filter.
3. A "DEEP" mock survey with 8 integration per exposure resulting in an integrated exposure time of 98.9 ksec per filter.

We plot the median 10σ depths in each of the NIRCcam filters we will use in this study as a function of total exposure time in Figure 4. These depths were calculated from our simulated noisy photometry and are appropriate for extended sources. Future *JWST*/NIRCcam deep surveys will likely be designed with longer exposures in less sensitive bands in order to balance the observational depth, and interested readers can explore the impact of such changes with the NIRCPrepareMock code we make publicly available. Deeper observations at bluer NIRCcam bands will preferentially affect the ability for a given survey to remove low redshift interlopers, while deeper observations in the detection bands for a given selection criterion will lead to a larger number of recovered high-redshift objects. More exposure time in longer-wavelength NIRCcam bands will be important for SED fitting, as these bands cover the rest-frame optical and a suite of strong emission lines in high-redshift galaxies.

In addition, in order to explore how NIRCcam observations can complement those made at shorter wavelengths by *HST*,

³ <http://jwst.etc.stsci.edu/>

Table 1. Simulated NIRCам 10 σ depths for the SHALLOW, MEDIUM, and DEEP surveys.

Filter	10 σ Depth (nJy)		
	SHALLOW	MEDIUM	DEEP
F070W	15.71	8.13	5.76
F090W	12.93	6.67	4.71
F115W	11.79	6.08	4.30
F150W	9.35	4.83	3.42
F200W	7.69	3.98	2.81
F277W	9.11	4.70	3.32
F335M	12.69	6.58	4.65
F356W	7.56	3.91	2.76
F410M	14.67	7.56	5.76
F444W	11.48	5.95	4.20

we simulate a region of the sky of 10.8 square arcminutes at the XDF ACS depth given by [Illingworth et al. \(2013\)](#). We simulate observations in the *HST*/ACS filters F435W (152.4 ksec, 7.06 nJy 10 σ depth), F606W (174.4 ksec, 5.00 nJy 10 σ depth), F775W (377.8 ksec, 5.99 nJy 10 σ depth), F814W (50.8 ksec, 21.93 nJy 10 σ depth), and F850LP (421.6 ksec, 10.61 nJy 10 σ depth), and generate NIRCам fluxes with the same depths as described in the previous paragraph, but over the smaller XDF area. For both mock surveys, we produced 500 noisy samples to explore how our noise estimates affect the uncertainties on the overall density of objects selected by a set of NIRCам color cuts.

2.4. Color Cut Figures of Merit

Because of the large variety of observed galaxy SEDs and photometric noise, there is no single ideal set of color selection criteria that will result in a clean sample of high-redshift galaxies. Our goal in this paper is to estimate statistics on the recovered population of simulated high-redshift galaxies as a function of our color cuts in order to aid in future NIRCам observations. For the purposes of this study, we require a definition of a "high-redshift object" and an "interloper" for a given dropout selection filter. While the Lyman limit is found at 912Å, absorption due to the Ly α forest causes the exact wavelength of the Lyman break to shift to longer wavelengths at higher redshifts, which is simulated within the JAGUAR catalog. At $z > 6$, this absorption is thick enough that the break occurs at 1216Å, the wavelength of Ly α . We define a high-redshift object as one that is above the redshift where the Ly α emission line crosses the half-power response of the blue side of the dropout band, and an interloper is any object that satisfies a given color selection criteria, but is below this redshift.

There are three primary statistics that we explore for choosing a given color selection criterion and assembling a high-redshift dropout sample:

1. The first statistic we report is selection "accuracy," defined as the ratio between the number of high-redshift objects selected to the *total* number of objects selected by a given color selection criterion. This is sometimes referred to as sample "purity" in the literature.
2. Extremely red selection limits will result in more accurate, but smaller total samples, so we also report the on-sky density of high-redshift objects under a given selection criterion, which we refer to as "true positive density", or TPD.
3. The final statistic we provide is selection "completeness," defined as the ratio between the number of high-redshift objects selected to the number of high-redshift galaxies that satisfy the SNR criteria (both red detections and blue non-detections).

Defining the optimal selection criteria will be determined by the trade-off between a more accurate sample, and one that has a higher number of high-redshift objects selected and a higher sample completeness.

2.5. Selecting High-Redshift Galaxies

High-redshift dropout candidates are often selected by observing flux at a given significance in multiple photometric filters at wavelengths longer than the break, with flux below a given significance at wavelengths shorter than the break. In this paper, we select mock galaxies based on a set of color criteria, and require objects to be selected in at least two filters to the red of the Lyman break above a signal-to-noise ratio (SNR) of 3 (although we will describe how our statistics change if we instead select above a SNR of 5, or 10). In addition, because of IGM absorption at rest wavelengths shorter than the Lyman break, we require a non-detection in the bands to the blue of the dropout filter at a SNR less than 2, as is commonly used in the literature (e.g., [Bouwens et al. 2015](#)). We should note that for the NIRCам-only simulations we will describe, F070W dropouts will not have a rejection band, while for F090W dropouts we will use F070W to help reject interlopers, and for F115W dropouts, we require non-detection fluxes in both F070W and F090W. For the *HST*+NIRCам simulations we will also use the *HST* bands for this rejection, highlighting the importance of using shorter wavelength data for selecting more pure samples of objects with fewer lower redshift interlopers. For F070W dropouts, we require non-detections at *HST*/ACS F435W. For F090W dropouts, we require non-detections at *HST*/ACS F435W, F606W, and NIRCам F070W. For F115W dropouts, we require non-detections at *HST*/ACS F435W, F606W, F775W, F814W, and NIRCам F070W and F090W.

3. RESULTS

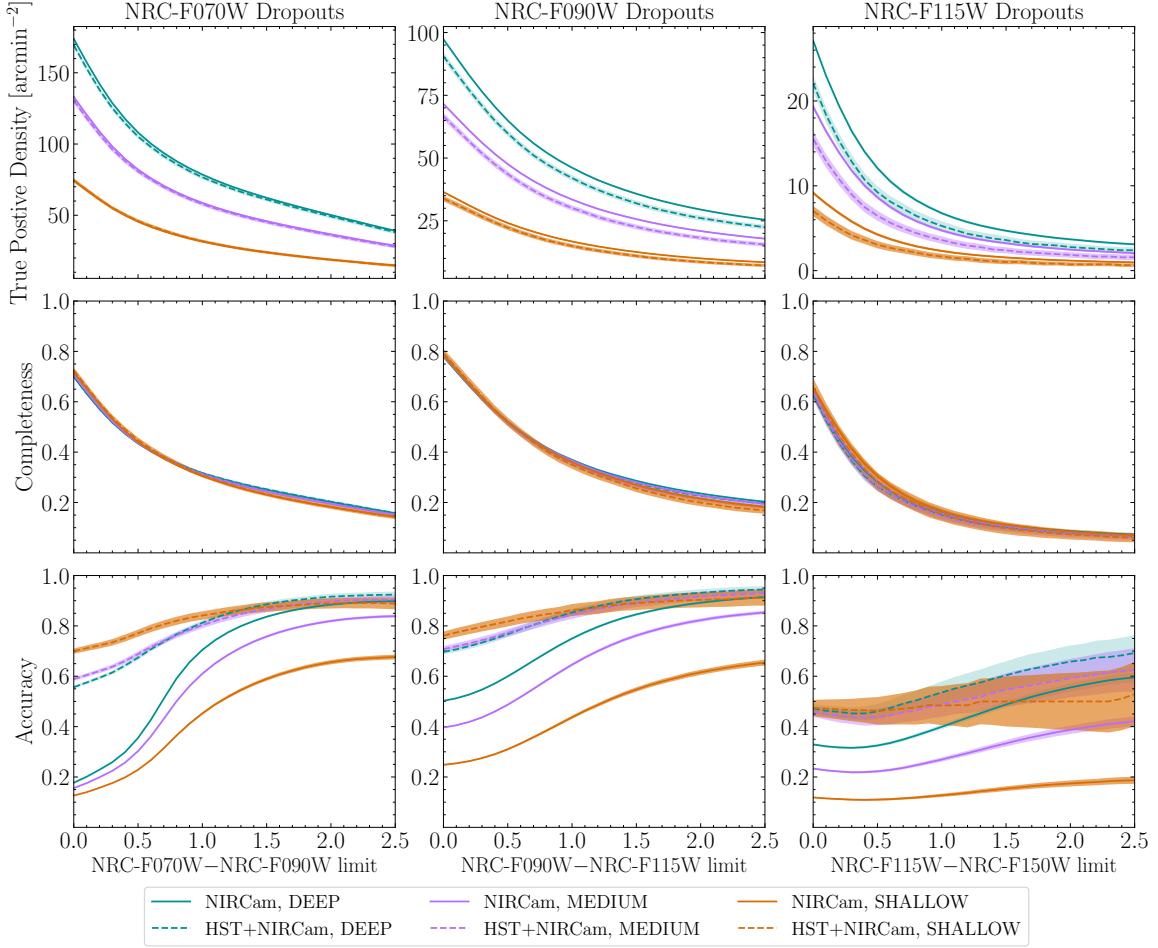


Figure 5. TPD (top row), completeness (middle row), and accuracy (bottom row) as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right), for the DEEP (teal), MEDIUM (lavender), and SHALLOW (orange) NIRC*am* (solid) and *HST*+NIRC*am* (dashed) surveys requiring a SNR > 3 in both detection bands. We additionally require a second color cut of F090W - F115W < 0.4 (left), F115W - F150W < 0.4 (middle), and F150W - F200W < 0.4 (right). The shaded regions indicate the 1σ range on the TPD and accuracy values calculated using the 500 noisy mock catalogs. In each set of panels, using a redder color cut results in lower TPD and completeness at a higher level of accuracy, and requiring a higher SNR limit reduces the overall TPD while increasing the accuracy at a given color cut.

In this section we discuss our number density, completeness, and accuracy results as a function of multiple factors, including survey depth, detection SNR, and survey design. For the majority of this analysis, we will adopt a simple two-color cut selection method, as illustrated in Figures 2 and 3, where we vary the color cut for the filters that straddle the Lyman break (the “Lyman break cut”, represented by a dashed line in these figures), and we fix the color requirement for the filters redward of the Lyman break (the “UV continuum cut” represented by a solid line in these figures). After testing the effects of varying the UV continuum cut on TPD and accuracy, we require F090W - F115W < 0.4 (magnitudes, for F070W dropouts), F115W - F150W < 0.4 (for F090W dropouts), and F150W - F200W < 0.4 (for F150W dropouts). We will be discussing the use of single color cuts or more complicated color selection methods further in Section 3.2.

3.1. Survey Depth and Detection SNR

The design of a survey, and especially the observational depth in the chosen filters, will have a strong impact on the number of high-redshift objects that are recovered with a given selection method. In Figure 5, we plot TPD (top panels), completeness (middle panels), and accuracy (bottom panels) against the Lyman break cut for our DEEP, MEDIUM, and SHALLOW survey depths. In each set of panels we utilize a detection SNR of 3.0 (for at least two filters to the red of the Lyman break), and ensure non-detections in the filters to the blue of the Lyman break as previously described. We plot the NIRC*am* only selection with solid lines, and the *HST*+NIRC*am* selection with dashed lines. We plot the 1σ range on the distribution of these values calculated using the 500 mocks with a shaded region.

For all three color selection criteria, at redder color cuts the density and completeness of recovered true high-redshift

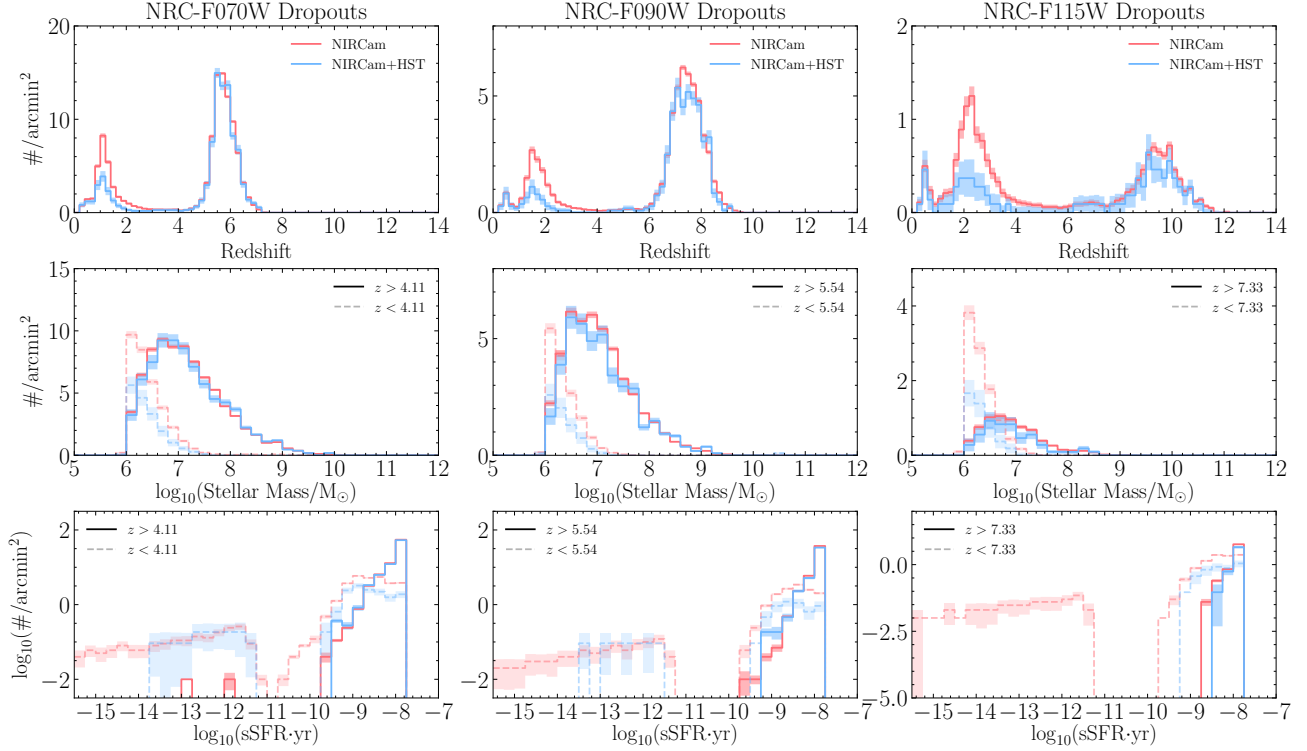


Figure 6. Histograms of mock galaxies at the DEEP survey depth with F070W - F090W > 1.0 (left), F090W - F115W > 1.0 (middle), and F115W - F150W > 1.0 (right), as well as a second color cut as described in Figure 5, with a detection SNR > 3 for all three plots. (Top) The spectroscopic redshifts of the dropouts, where in red we plot the number density of objects with NIRCам data alone, and in blue we plot those objects with *HST*+NIRCам data. The addition of deep *HST* data for constraining blue non-detections has a significant effect in removing interlopers. (Middle) The JAGUAR stellar masses of these objects are plotted, with red and blue as in the top panel, but now the dashed lines correspond to a subsample of low-redshift interlopers in each panel, while we plot the mass distribution of the true high-redshift objects with a solid line. (Bottom) The JAGUAR specific star formation rates (sSFR) are plotted similar to the middle row, but using a logarithmic scale on both axes. Interlopers are primarily found at low stellar masses, although for the F070W and F090W dropouts, a number of higher-mass quiescent mock galaxies are selected as interlopers with these color cuts, which are also found at lower sSFR values.

sources decreases, but the accuracy of the sample increases. The total number of recovered sources, as well as the overall accuracy, increases at deeper survey depths. The completeness, however, does not depend strongly on survey depth, as this statistic is a ratio between two values that depend on depth in roughly the same manner. At a detection SNR > 3, it is only possible to reach high levels of accuracy with extremely red color cuts. In all survey depths and dropout criteria, the accuracy plateaus to a value less than 1.0 owing to contamination by mock galaxies at low redshifts and low SNR with non-detections in the bluer filter of the Lyman break color cut. As a result, these objects have extremely red Lyman break colors, and would be contaminants at any choice of cut. We should also note that the 1σ distributions are much larger for the NRC-F115W dropouts in the SHALLOW depth survey because of the small number of objects recovered at this survey depth.

While the usage of *HST* blue filter non-detections results in overall lower densities of actual high-redshift objects, as would be expected, it has a much larger effect on the accuracy. For F090W - F115W > 1.0, in the DEEP survey,

with NIRCам observations only, the density of sources at $z > 5.5$ is 50 arcmin⁻², at an accuracy of 0.70, while with *HST*+NIRCам observations, the density is 10% smaller, but at an increased accuracy of 0.80. Interestingly, when using *HST* fluxes, the measured accuracy at blue color cuts for the SHALLOW depth survey is higher than for the MEDIUM or DEEP surveys. The addition of an *HST* SNR cut has a strong effect on reducing the total number of galaxies selected by a set of color cuts (and thereby increasing accuracy) which is more significant in the SHALLOW survey due to the larger flux uncertainties. The discrepancy between the NIRCам and the *HST*+NIRCам TPD values is larger for higher-redshift dropouts, because of the additional blue filters that are used to reject low-redshift interlopers. The recovered completeness is not significantly different between *HST*+NIRCам and NIRCам observations only.

We can examine in more detail the properties of the mock galaxies that are recovered by a specific color cut. In Figure 6, we show the redshift (top), stellar mass (middle), and specific star formation rate (sSFR, defined as the mock galaxy star formation rate normalized by the stellar mass, bottom)

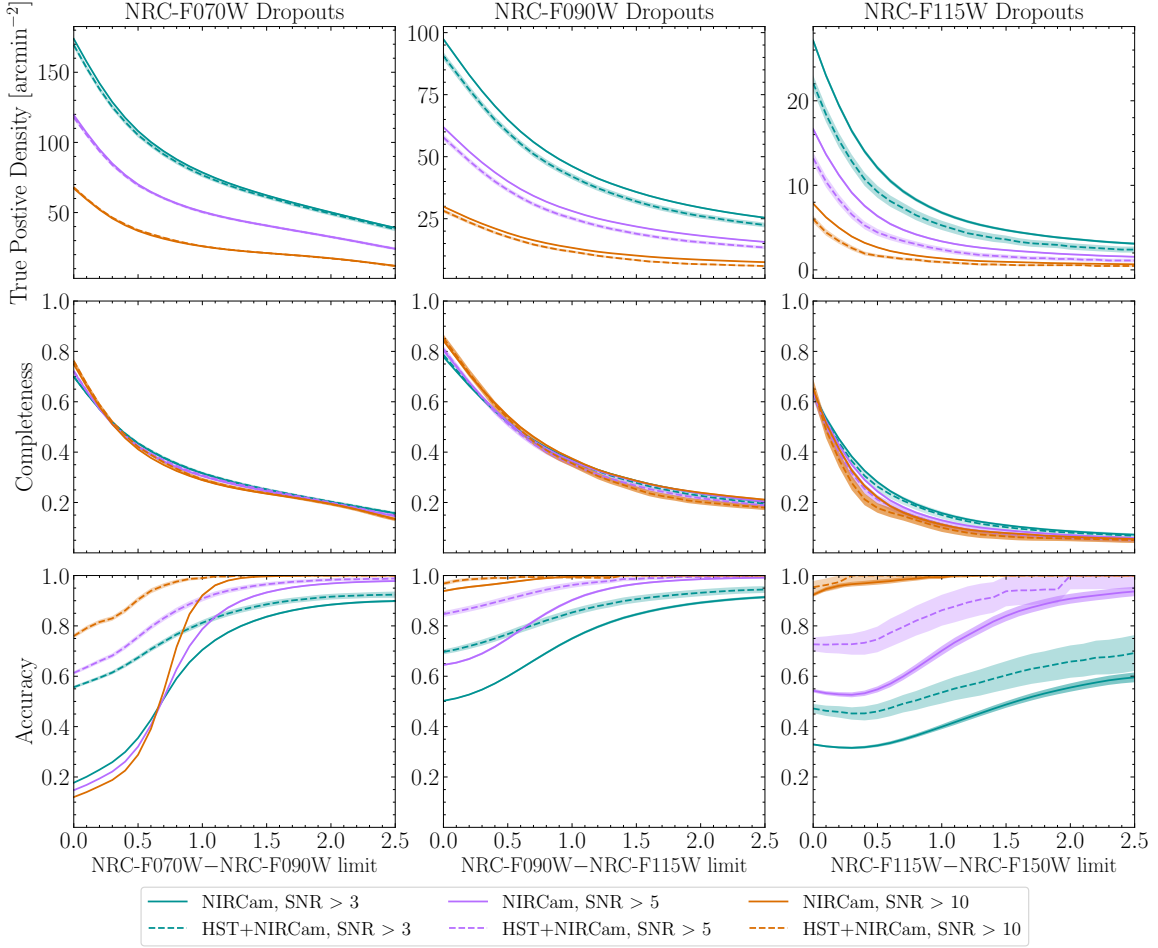


Figure 7. TPD (top row), completeness (middle row), and accuracy (bottom row) as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right), for a detection SNR of > 3 (teal), > 5 (lavender), and > 10 (orange) for NIRCam (solid) and *HST*+NIRCam (dashed) DEEP surveys. We additionally require a second color cut of $F090W - F115W < 0.4$ (left), $F115W - F150W < 0.4$ (middle), and $F150W - F200W < 0.4$ (right). In each set of panels, a higher SNR restriction leads to an increase in the accuracy, but at a significant decrease in TPD.

distributions for the mock galaxies selected by color cuts of $F070W - F090W > 1.0$ (left), $F090W - F115W > 1.0$ (middle), and $F115W - F150W > 1.0$ (right) (in addition to the UV continuum cuts described above) for the DEEP survey. In all three columns, we plot the NIRCam-only selection in red, and the *HST*+NIRCam selection in blue. In each case, we can see how mock galaxies at $z \sim 1-4$ are the primary contaminants, and based on the mass distributions, these objects have masses $10^6 - 10^7 M_\odot$ and lower sSFR values. The addition of the *HST* data helps mitigate the contaminants, but in all cases, red, low-mass, faint mock galaxies are selected as Lyman-break galaxies.

We also explored how the detection SNR affects the dropout selection. For the DEEP survey depth, using the NIRCam and *HST*+NIRCam observations, we calculated the TPD, completeness, and accuracy for detection SNR of 3, 5, and 10, and plot these results in Figure 7 for the F070W, F090W, and F115W dropouts. Changing the SNR has a

strong effect on the accuracy of the recovered samples, such that samples with greater than 90% accuracy can be recovered with a detection SNR of 5 - 10 at redder color cuts. However, this comes at a significant cost to the recovered TPD: almost twice as many objects are detected for a detection SNR of 3 vs. 5, and at 5 vs. 10 at all color limits. While the completeness is similar between the different SNR cuts for F070W and F090W dropouts, we find a slightly higher completeness for the detection SNR of 3 for F115W dropouts.

3.2. Single Color Cut vs. Two Color Cut Selection

Throughout this analysis, we have shown results with a simple two-color selection method, as star-forming mock galaxies at lower redshift can have red colors which mimic dropout galaxies. In this section we additionally explore the recovery of high-redshift galaxies with a method which uses only a single color cut, as well as a compound method which uses two color cuts and an additional angled color selection

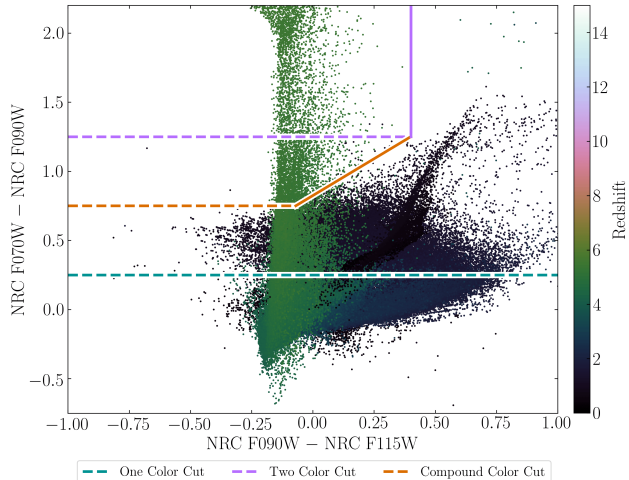


Figure 8. NIRCam color-color plot with the three selection criteria that we employ and compare in this paper. In lavender, we show the two-color selection method used in Section 3.1. In teal and orange, we show the One Color and Compound Color Cut selection methods, respectively.

as is often used to select high redshift galaxies. In Figure 8, we show the F070W dropout selection color space marked to show the two-color selection we have used up to this point (lavender), single-color selection (teal), and the compound color selection (orange). Because dust obscuration in a galaxy results in redder colors, for each dropout selection criteria, the slope of our angled cut corresponds to the reddening vector for the Calzetti et al. (2000) dust prescription for that filter combination⁴. To simulate the different selection methods, we repeated our previous analysis using these alternate selection methods on the DEEP survey, with a detection SNR of 3.0, but we fix the solid lines shown in Figure 8 and explore how changing the color indicated by the dashed lines impacts the recovery of high-redshift galaxies.

In Figure 9, we plot the TPD, completeness, and accuracy for F070W, F090W, and F115W dropouts, comparing the three color cut methods as shown in Figure 8. Not surprisingly, the One Color Cut method leads to a larger TPD for all three dropout selection techniques, as fewer objects are excluded. The completeness for the three selection methods is very similar to the TPD, in that the One Color Cut method results in larger completeness at all color cuts. The accuracy values of the One Color Cut and Two Color Cut methods are very similar for all three dropout selection techniques, likely due to the $\text{SNR} > 3$ detection threshold. At such a low SNR value, the large noise scatter on the mock galaxy colors leads to similar accuracy levels with or without the UV continuum cut. If we use a detection $\text{SNR} > 10$, the accuracy for the Two

Color Cut method is larger at all color limits than that for the One Color Cut Method, as shown for F070W dropouts in the bottom-left panel of Figure 9. The Compound Color Cut method results in the highest accuracy levels at bluer color cuts, but the third angled cut removes a significant fraction of high-redshift mock galaxies for all three dropout methods. We also find that the using a Compound Color Cut does remove relatively brighter ($m_{\text{AB},\text{F115W}} < 27$) interloper galaxies but in addition a number of faint ($m_{\text{AB},\text{F115W}} \sim 30$) high-redshift galaxies also are culled. The impact of the angled color cut on accuracy is lessened at higher redshift, where there are less dusty and quiescent mock galaxies in the JAGUAR catalog, and less reason for using a third color cut. The key result from this analysis is that for F070W and F090W dropouts, it is possible to get a significant number of candidates with an accuracy level greater than 70% by employing a Compound Color Cut.

We have shown results using the Two Color Cut method with a fixed UV continuum cut of < 0.4 . To explore how changing this second color cut affects the resulting TPD, completeness, and accuracy, we looked at selecting high-redshift dropout candidates by fixing the Lyman break cut and varying the UV continuum color cut (In Figure 8, this would amount to fixing the dashed lavender line and changing the solid lavender line). For this analysis, we set $\text{F070W} - \text{F090W} > 1.0$, $\text{F090W} - \text{F115W} > 1.0$, and $\text{F115W} - \text{F150W} > 1.0$, and looked at mock galaxies at the DEEP survey depth. We show how TPD, completeness, and accuracy vary with the second color limit and the detection SNR in Figure 10.

In these plots, we show that while TPD and completeness increases as the color cut becomes more inclusive, the accuracy falls, especially for F070W dropouts, due to the larger number of low-redshift interlopers. Because of these results, we have adopted a uniform color cut in our Two Color Method selection of $\text{F090W} - \text{F115W} < 0.4$ for F070W dropouts, $\text{F115W} - \text{F150W} < 0.4$ for F090W dropouts, and $\text{F150W} - \text{F200W} < 0.4$ for F115W dropouts. These color cuts correspond to UV slope $\beta < -0.52$ for F070W dropouts, $\beta < -0.61$ for F090W dropouts, and $\beta < -0.70$ for F115W dropouts.

3.3. The Impact of Ly α Emission on Color Selection

The presence of the Ly α emission line can contribute flux to the filters used in selecting high-redshift galaxies, potentially impacting the numbers of galaxies that are recovered by a given cut. In the JAGUAR mock catalog, the median Ly α rest-frame Equivalent Width (EW) is 74 Å for mock galaxies at $z > 4.1$, the redshift where Ly α enters the NIRCam F070W filter. At $z = 7$, an emission line with this EW would result in F090W magnitude difference of $\Delta m_{\text{AB}} = 0.29$.

Ly α is a resonant line, and its emission is highly dependent on the geometry of the gas in the galaxy as well as the surrounding IGM (Neufeld 1991; Giavalisco et al. 1996; Kunth et al. 1998; Frye et al. 2002; Shapley et al. 2003), although this resonance is not modeled for the JAGUAR mock catalog galaxies. At $z > 6$, the IGM has been observed to be increasingly neutral which significantly reduces the frac-

⁴ For F070W dropouts we use a slope of 1.07 and an intercept of 0.82, for F090W dropouts we use a slope of 1.03 and an intercept of 0.84, and for F115W dropouts we use a slope of 1.07 and an intercept of 0.82.

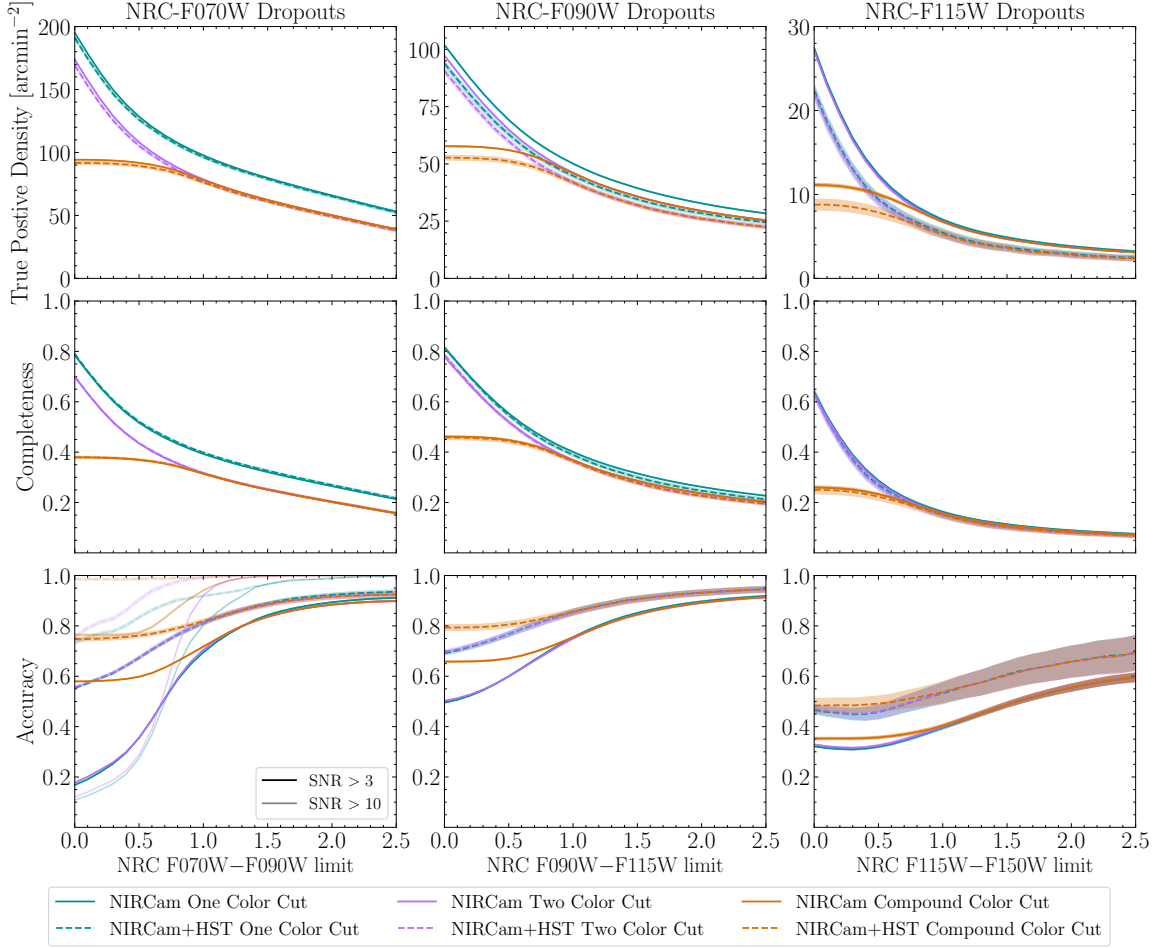


Figure 9. TPD (top row), completeness (middle row), and accuracy (bottom row) as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right) with the One Color Cut (teal), Two Color Cut (lavender), and Compound Color Cut (orange) methods for NIRCcam (solid) and *HST*+NIRCcam (dashed) surveys. Above a given Lyman break color cut, the Two Color and Compound Color cut results become identical. At bluer cuts, the Compound Color selection has a higher accuracy but a lower TPD and completeness. In general, the One Color selection has a lower accuracy, except at the reddest cuts. This is clearly seen at higher detection SNR, as illustrated by comparing the accuracy at SNR > 3 and SNR > 10 for F070W dropouts in the bottom-left panel.

tion of galaxies with observed Ly α in emission (Stark et al. 2010; Pentericci et al. 2011, 2014; Caruana et al. 2012, 2014; Schenker et al. 2012, 2014; Treu et al. 2013; Tilvi et al. 2014). To explore how Ly α emission affects our ability to recover high-redshift galaxies with NIRCcam, we used a version of the JAGUAR mock catalog that was created without modeling Ly α but is otherwise identical. We repeated our color-cut analysis at the DEEP survey depth, with a two-color selection and a detection SNR of 3. We plot these results in Figure 11.

The presence of Ly α emitted by a galaxy has a subtle effect on dropout selection. We can illustrate this by looking at the TPD and completeness for the F070W dropouts. At blue selection colors, these values are higher for the sample without Ly α emission, and then at redder selection colors they are higher for the sample that includes Ly α emission. For F070W dropouts, we select objects at $z > 4.11$, which includes objects where Ly α is entering the F070W band, enhancing the flux, and making the F070W-F090W

color bluer than it would otherwise be without Ly α emission. At the same time, for objects at a redshift where Ly α sits in the F090W filter, this contributes to the flux in this band, causing these mock galaxies to be *redder* in the Lyman break color cut, and *bluer* in the UV continuum cut. Mock galaxies with Ly α emission are then both bluer and redder in the Lyman break cut depending on their redshifts, which impacts their selection as seen in the top and middle left panels of Figure 11. This effect is also observed for the TPD in the F090W and F115W dropout panels, but at less significance and redder selection colors. We find that the accuracy for dropout samples without Ly α emission is higher than samples with the emission line, with the highest significance for F115W dropouts. Similar results were seen for observations of galaxies in the *Hubble* Ultra Deep Field (HUDF) with the Multi Unit Spectroscopic Explorer (MUSE) in Inami et al. (2017), where these authors present HST color cuts to select for Ly α emitters at $2.9 < z < 6.7$.

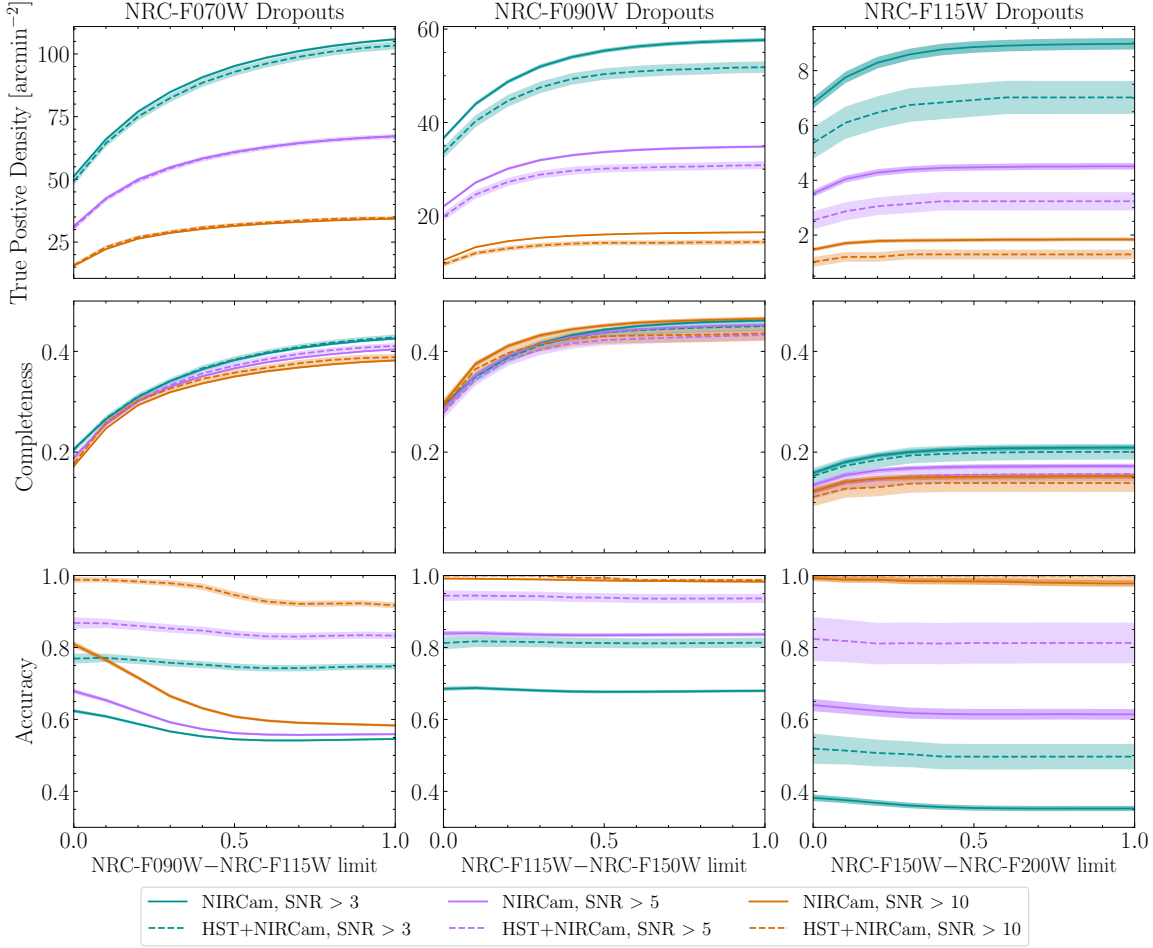


Figure 10. TPD (top row), completeness (middle row), and accuracy (bottom row) as a function of second color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right) as a function of detection SNR; SNR > 3 (teal), SNR > 5 (lavender), and SNR > 10 (orange), for NIRCcam (solid) and *HST*+NIRCcam (dashed) surveys. In each panel, we fix the first color cut to F070W - F090W > 1.0 (left panel), F090W - F115W > 1.0 (middle panel), and F115W - F150W > 1.0 (right panel).

3.4. Alternate Color Selection Criteria

Thus far, we have only explored NIRCcam color selections using three adjacent photometric bands (along with non-detections in photometric bands shortward of the Lyman break). Lyman break selection, however, uses a pair of observed colors: one that spans the Lyman break at a particular redshift and one that covers the relatively featureless UV stellar continuum from massive stars. In this section, we examine the TPD, completeness, and accuracy for alternate UV continuum color cuts which utilize two unique photometric bands (“Four-band color selection”) and a scenario where the UV continuum cut attempts to span the entire rest-UV portion of the galaxy SED (“long UV baseline”).

In the three-band selection methods we have outlined thus far, mock galaxies can artificially be driven into or out of the selection boxes because of noise in the common photometric band. To help explore this effect, we also explored selecting Lyman break galaxies using photometry with four distinct NIRCcam bands. While a four-band color selection criterion would require additional deep observations, it has the added

benefit that noise in a single photometric band cannot affect both colors being used to select the galaxy.

For the four-band analysis, we updated our selection criteria and re-ran the selection tests as was done in previous sections. For F070W dropouts, we compared F070W - F090W and F115W - F150W colors. For F090W dropouts, we compared F090W - F115W and F150W - F200W colors. Finally, for F115W dropouts, we compared F115W - F150W and F200W - F277W colors. In all cases, we used the DEEP survey depth, with a 3σ detection, and explore the two-color cut selection, varying the Lyman break cut (we fixed the UV continuum cut in each case using a similar test to what was done in Section 3.1 for the three-band selection). In Figure 12, we plot the TPD and accuracy for the four-band selection criteria compared to the three-band selection criteria.

For both F070W and (with less significance) F090W dropouts, the four-band color selection results in a larger TPD and completeness and a higher accuracy at all color cuts we explored. For three-band selection, dropout galaxies that are at redshifts where the Lyman break has entered one of

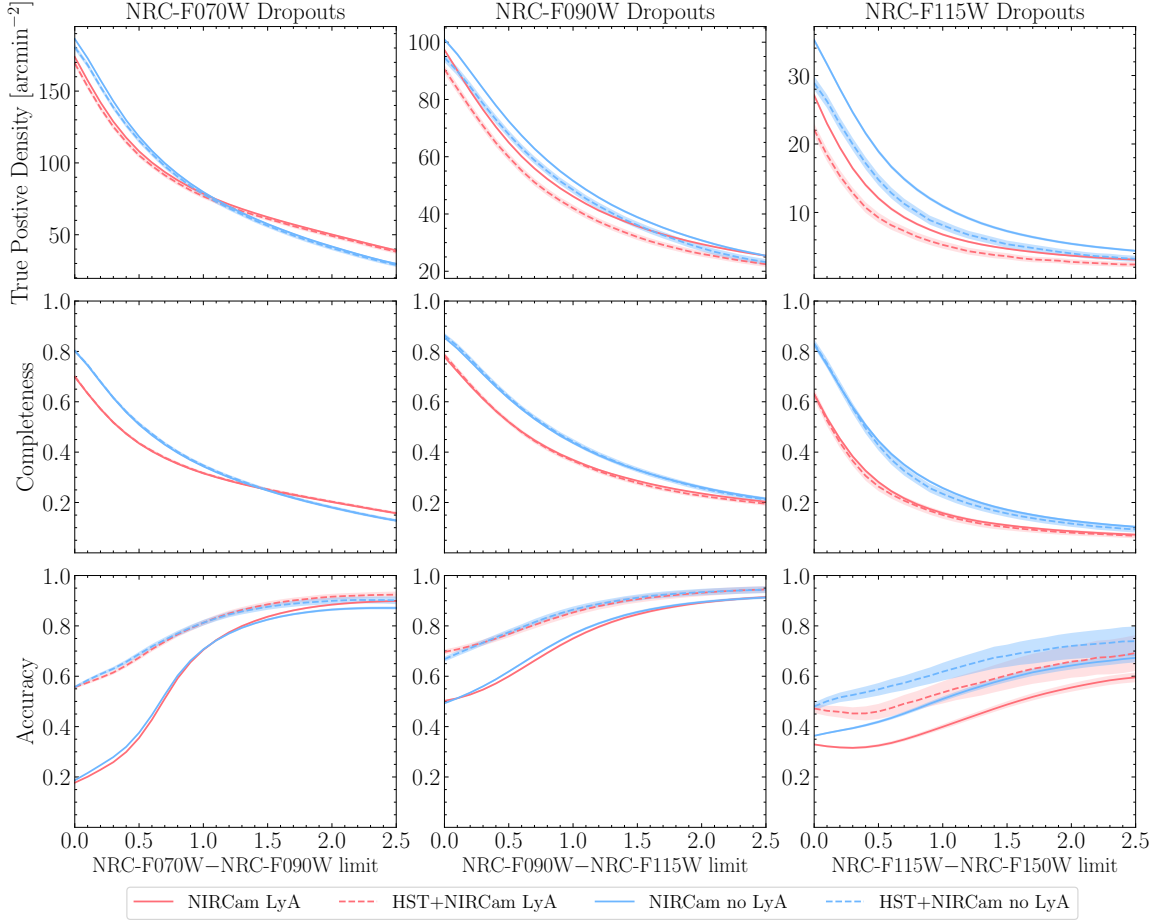


Figure 11. TPD (top row), completeness (middle row), and accuracy (bottom row) as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right), with (red) and without (blue) Ly α emission for NIRCcam (solid) and *HST*+NIRCcam (dashed) surveys. We additionally require UV continuum color cut of F090W - F115W < 0.4 (left), F115W - F150W < 0.4 (middle), and F150W - F200W < 0.4 (right). Ly α emission results in selection with a lower TPD and completeness at redder color cuts, and a higher TPD and completeness at bluer color cuts. The accuracy is similar between the two catalogs, except for F115W dropouts.

the bands used in the UV color cut will be rejected for being too red, which results in fewer total high-redshift galaxies selected. With four-band color selection, the UV color cut samples a longer wavelength region of the SED, and this effect is not observed, leading to a higher TPD. For F115W dropouts, this effect is less significant (due to the declining number of very high-redshift galaxies), and the addition of the F277W SNR requirement results in lower TPD values, although this is also at an increased accuracy. As a result, our results demonstrate that four-band selection is recommended for F070W and F090W dropouts, or at very large areas (where the impact to the recovered TPD is minimal) for F115W dropouts.

We also examined a selection criteria where the UV continuum cut spans a longer wavelength range across the rest-frame ultraviolet. For F070W dropouts, in this alternate selection criteria, the UV continuum cut is F090W - F200W, for F090W dropouts, the alternate UV continuum cut is F115W - F277W, and for F115W dropouts, the alternate UV continuum cut is F150W - F335M. As with the four-band color

selection, we updated our selection criteria and re-ran the selection tests for the DEEP survey depth and a detection SNR > 3.0, and explore the two-color cut selection, varying the Lyman break cut (we fixed the UV continuum cut in each case). In Figure 13, we plot the updated color-color diagrams for F070W (left column), F090W (middle column), and F115W (right column) dropout selection with a longer UV baseline. Here, we compare the TPD (second row), completeness (third row), and accuracy (fourth row) between cuts estimated with this longer UV baseline color criteria to those made with the original criteria from Section 3.1. For F070W dropouts, the use of a longer UV baseline results in a higher TPD and completeness but at similar accuracy values. This is likely because of the increased sensitivity for the F200W filter compared to the F115W filter, which leads to more objects satisfying the detection threshold. For F090W dropouts, the TPD and completeness is lower when using a longer UV baseline at similar (but slightly lower) accuracy values. For F115W dropouts, while the TPD decreases by a factor of two, the completeness only decreases slightly, but the accuracy

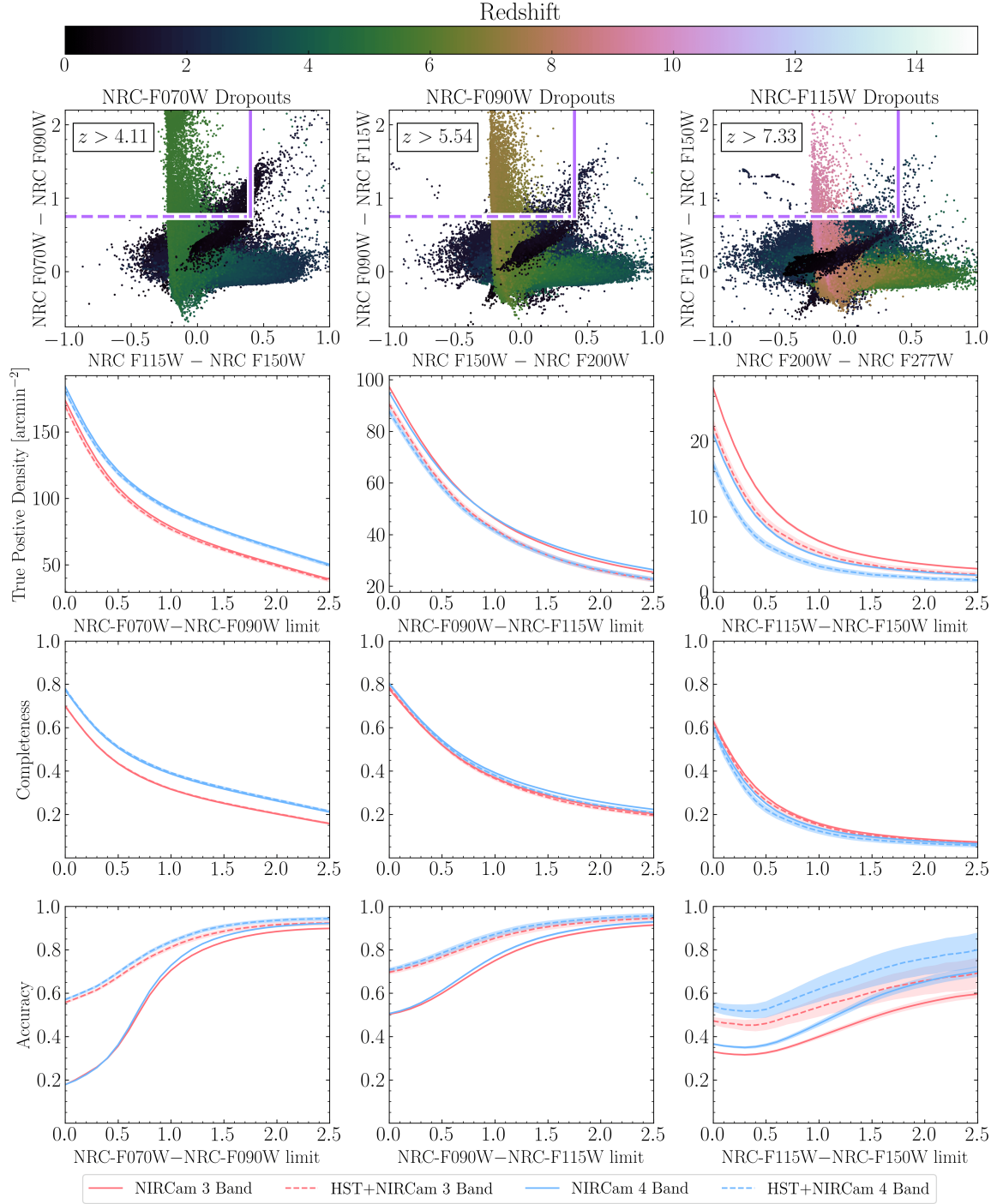


Figure 12. Color-color plots (first row) and the resulting TPD (second row), completeness (third row) and accuracy (fourth row) plots as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right) with three-band (red) and four-band (blue) selection for NIRCcam (solid) and *HST*+NIRCcam (dashed) surveys. For comparison, we require a UV continuum color cut of < 0.4 . TPD and completeness are higher with four-band selection than with three-band selection for F070W dropouts, similar between the two methods for F090W dropouts, and lower with four-band selection for F115W dropouts. With the exception of F115W dropouts, the accuracy of the four-band color selection is relatively consistent with that of the three-band selection.

increases when utilizing a longer UV baseline. These differences reflect how a longer UV baseline results in a larger

spread in color values, as seen in the first row of Figure 13, causing more objects to scatter outside of the selection boxes.

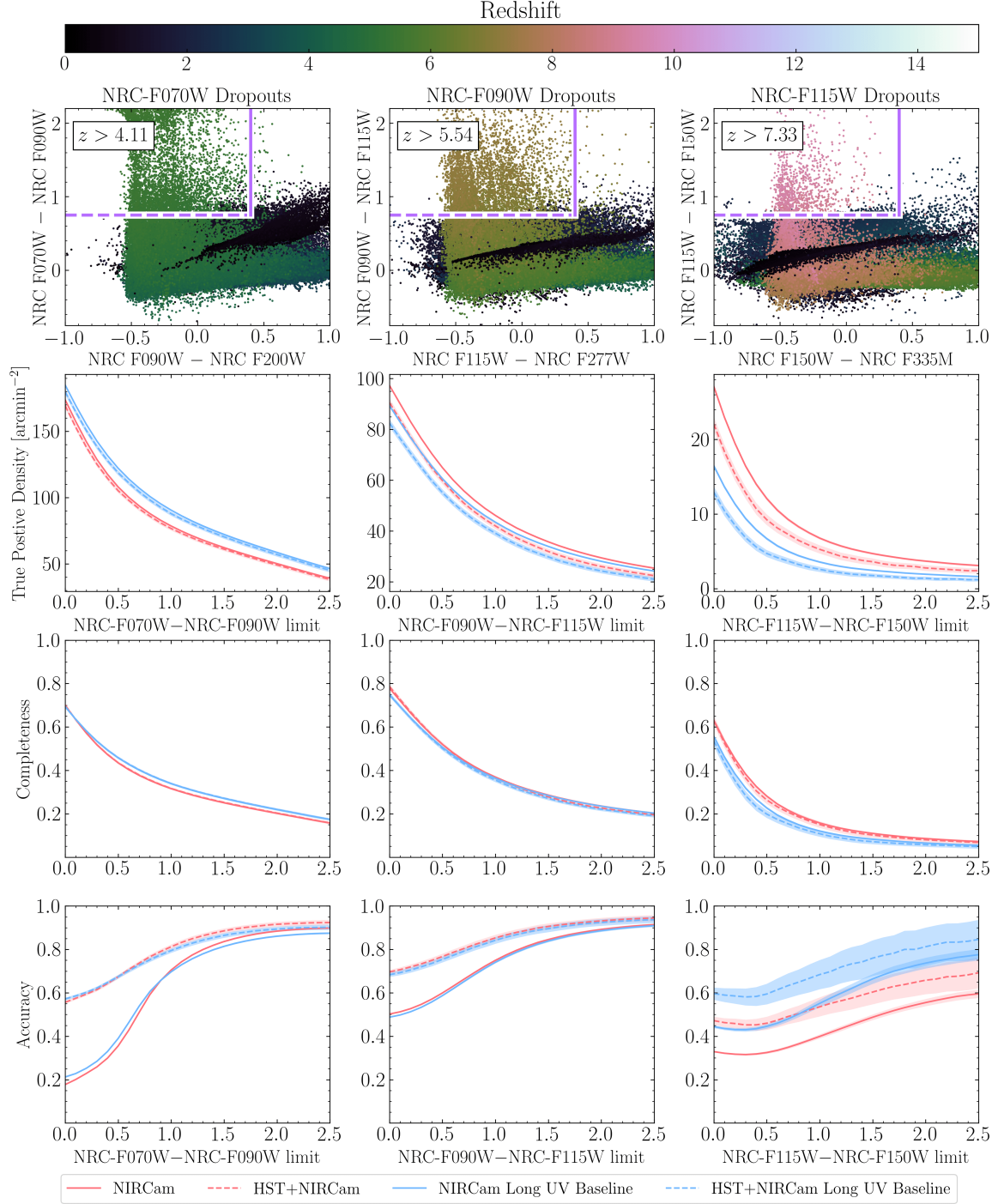


Figure 13. Color-color plots (first row) and the resulting TPD (second row), completeness (third row), and accuracy (fourth row) plots as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right), with the three-band selection from Section 3.1 (red) and with a three-band selection that utilizes a longer UV baseline (blue) for NIRCам (solid) and *HST*+NIRCам (dashed) surveys. For comparison, we require a UV continuum color cut of < 0.4 .

3.5. NIRCам Long Wavelength Rejection Colors

One of the key features of the NIRCам instrument is a dichroic beam splitter which allows observations in two filters simultaneously, one at short wavelengths (0.6 - 2.3 μm), and one at longer wavelengths (2.4 - 5.0 μm). Future deep

extragalactic surveys with NIRCам will need to utilize the dichroic to increase the efficiency of the observational strategy. In this section, we will discuss the usage of NIRCам observations made at longer wavelengths to help select high-

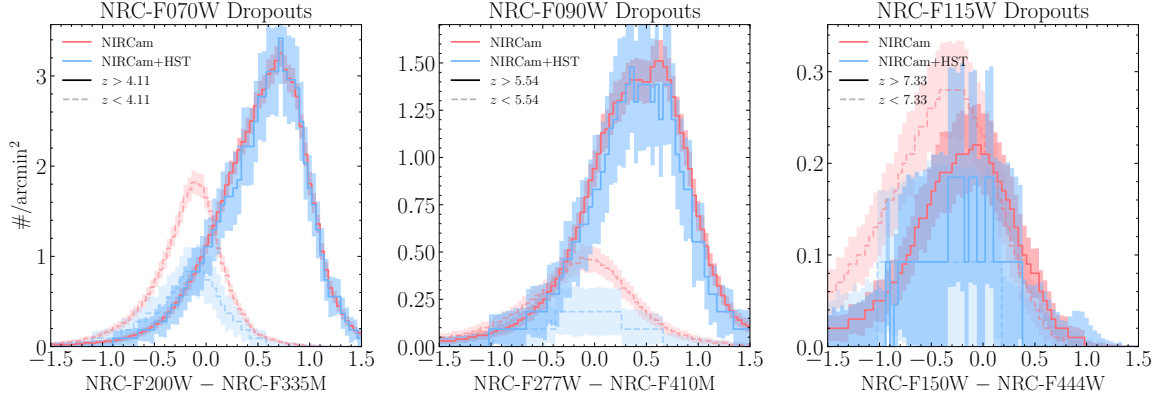


Figure 14. Color histograms of mock galaxies at the DEEP survey depth with F070W - F090W > 1.0 (left), F090W - F115W > 1.0 (middle), and F115W - F150W > 1.0 (right), as well as a second color cut as described in Figure 5, with a detection SNR > 3 for all three plots. In the left panel, we plot F200W - F335M color, in the middle panel we plot F277W - F410M color, and in the right panel we plot F150W - F444W color. In each panel, true high-redshift objects are plotted with a solid line and lower-redshift interlopers are plotted with a dashed line. We show results from NIRCcam-only photometry only with a red line, and HST+NIRCcam with a blue line. For F070W and F090W dropouts, the true high-redshift objects are found at redder short-to-long-wavelength colors than the interlopers. For F115W dropouts, because NIRCcam photometry does not cover the optically red portion of the SED in interloper galaxies, the difference between interloper and true high-redshift galaxy colors is less pronounced.

redshift galaxies (following work done in previous sections) as well as reject low-redshift interlopers.

For Lyman break selection, in this current work we only explore photometric bands that cover the rest-frame UV, as at longer wavelengths the addition of flux due to strong emission lines and the 4000Å +Balmer break results in redder UV-to-optical colors. In Section 3.4, we demonstrate the TPD, completeness, and accuracy for a UV continuum color cut with a longer wavelength range. For F070W dropouts, the UV extends only to the F200W filter, and so longer wavelength data will only probe the rest-frame optical and near-IR. However, for F090W and F115W dropouts, observations can be made with the F277W and F335M filters respectively, which can be done simultaneously alongside shorter wavelength observations.

In addition, longer wavelength data can be used to reject low-redshift interloper galaxies by virtue of the overall color differences between short and longer wavelength observations between these two samples. In true high-redshift galaxies, NIRCcam long-wavelength filters cover the rest-frame optical (see Figure 1), which may have boosted flux due to the 4000Å +Balmer break and optical line emission. In interlopers, however, NIRCcam long-wavelength data samples the continuum drop-off in the near-IR (in the absence of significant very hot dust emission). By comparing a short to a long wavelength filter, interloper mock galaxies in JAGUAR are observed to be systematically bluer than the true high-redshift mock galaxies.

In Figure 14, we plot color distributions for mock galaxies observed as part of our DEEP survey, with a detection SNR > 3.0. In each panel, we plot the distribution of dropout galaxies above our redshift cuts with solid lines, and the interloper galaxies in dashed lines, and we plot in red and blue the distributions with NIRCcam data alone and NIRC-

Cam+HST data respectively. In the left panel, we plot the F200W - F335M color distribution, and the distribution of high-redshift dropout mock galaxies is significantly redder than the interloper distribution. For F070W dropout galaxies at $z > 4.11$, the NIRCcam F335M filter covers the [OIII]λ5007 emission line, leading to the red color. The [OIII]λ5007 emission line has been inferred to be strong (with EW values of 500) at these redshifts from Spitzer IRAC observations (Labbé et al. 2013; Smit et al. 2015; De Barros et al. 2019). In the middle panel, we plot the F277W - F410M color for F090W dropouts and interlopers, and we see a similar behavior, as the F410M filter covers [OIII]λ5007 for true high-redshift dropout galaxies. The difference is not as great as what is observed for the F070W dropouts, as the F410M filter no longer samples the near-IR wavelength range for the interloper galaxies. We find that a color cut at F200W - F335M > 0.0 or F277W - F410M > 0.0 aids in rejecting interloper galaxies.

The situation for F115W dropouts is more complex because NIRCcam short-to-long-wavelength colors are very similar in both true high-redshift galaxies and interlopers. In the right panel of Figure 14, we plot the F150W - F444W (the longest wavelength wide-band NIRCcam filter) colors for F115W dropout galaxies at $z > 7.33$ and interloper galaxies. We see that there is a tendency for interlopers to be at slightly bluer colors than true high-redshift mock galaxies, although with less significance owing to the small numbers of these galaxies in a given sample. We explored other color combinations besides F150W - F444W, but each had similar or worse results for rejecting interlopers.

To explore the use of these cuts for F070W, F090W, and F115W dropout selection, we calculated the TPD, completeness, and accuracy values for mock galaxies at the DEEP survey depth with SNR > 3.0, but required F200W - F335M

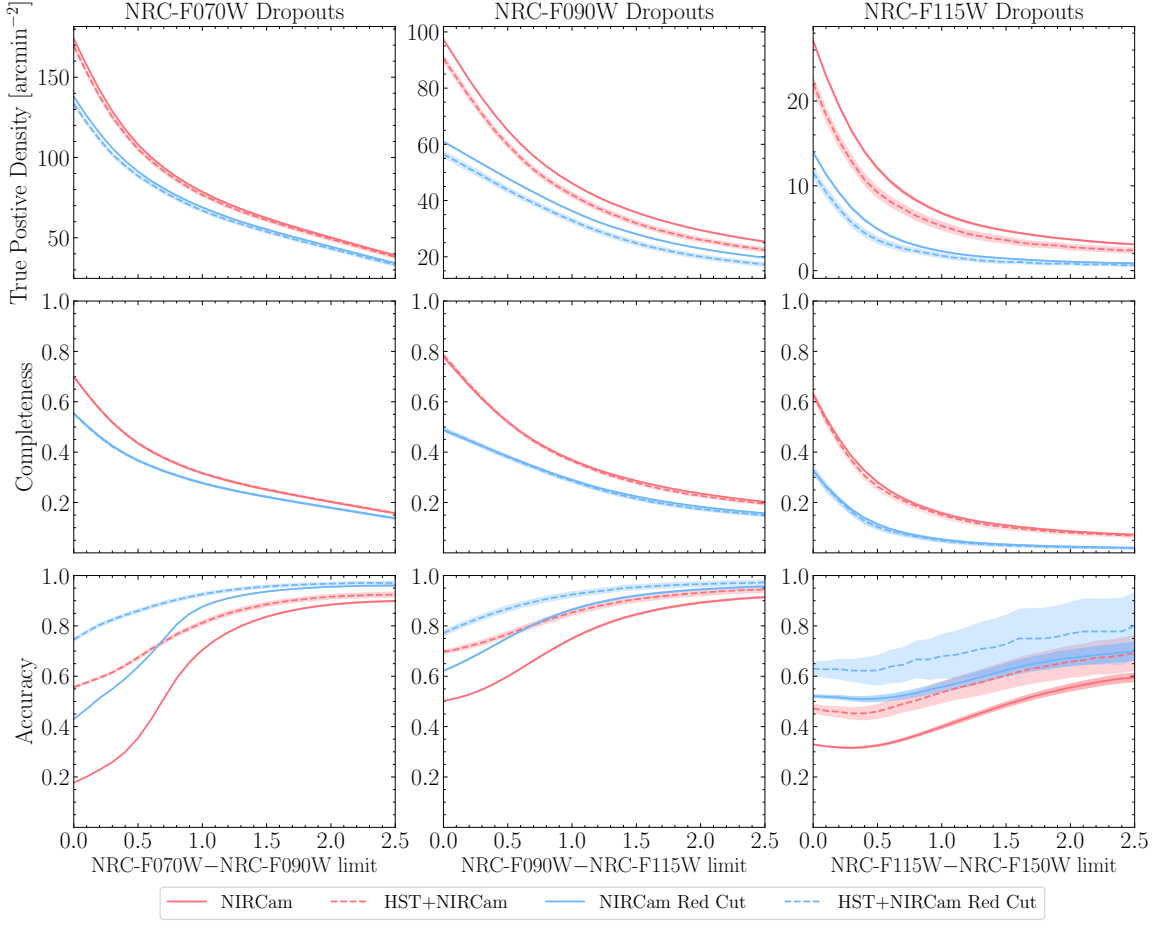


Figure 15. TPD (top row), completeness (middle row), and accuracy (bottom row) as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right), with (blue) and without (red) an additional short-to-long wavelength color cut for NIRCam (solid) and *HST*+NIRCam (dashed) surveys. For F070W dropouts, we require $F200W - F335M > 0.0$, for F090W dropouts, we require $F277W - F410M > 0.0$, and for F115W dropouts, we require $F150W - F444W > 0.0$. We additionally require UV continuum color cut of $F090W - F115W < 0.4$ (left), $F115W - F150W < 0.4$ (middle), and $F150W - F200W < 0.4$ (right). While the addition of these short-to-long-wavelength color cuts results in a decrease in the TPD and completeness, the accuracy increases significantly for all three selection criteria.

> 0.0 for F070W dropouts, $F277W - F410M > 0.0$ for F090W dropouts, and $F150W - F444W > 0.0$ for the F115W dropouts, which we plot compared to the TPD, completeness and accuracy made without the cuts in Figure 15. The bottom panels show the increase in accuracy that can be achieved through these short-to-long wavelength color cuts, although this is at the expense of the TPD and completeness plotted in the top and middle panels, which both drop by almost half for F115W dropouts.

This sets up a potential example observational strategy for F070W, F090W, and F115W dropouts. For F070W dropouts, short-wavelength observations would need to be made at F070W, F090W, and F115W (or either F150W or F200W) for the dropout selection, but these data could be supplemented by simultaneous observations with F335M, as well as longer-wavelength data (F356W and F410M), which is important for any potential SED fitting of these galaxies. For F090W dropouts, it is much more straightforward. The short wavelength data necessary would be at F090W and F115W,

which could be observed simultaneously with the F277W and F410M filters. Similarly, for F115W dropouts, the short wavelength data necessary would be at F115W and F150W, which could be observed simultaneously with the F335M (or F277W) and F444W filters.

3.6. Interlopers and χ^2_{opt}

In Bouwens et al. (2015), the authors explore the usage of a statistic they refer to as χ^2_{opt} , defined as

$$\chi^2_{\text{opt}} = \sum_i \text{SGN}(f_i)(f_i/\sigma_i)^2 \quad (1)$$

where for each undetected ($\text{SNR} < 2$) photometric band to the blue of the Lyman break, f_i is the flux in that band, σ_i is the uncertainty in that band, and $\text{SGN}(f_i)$ is 1 if $f_i > 0$ and -1 if $f_i < 0$. This statistic was designed to measure whether, for objects with a non-significant detection in the filters to the blue of the Lyman break, the flux is biased towards positive values. For actual high-redshift galaxies, the distribution

of χ_{opt}^2 should be centered at 0, while for lower-redshift interlopers, the distribution will be biased towards positive values.

We explore the efficacy of the Bouwens et al. (2015) χ_{opt}^2 statistic in discriminating low-redshift interlopers in *JWST* surveys using our three-band analysis at the DEEP survey depth. We calculate the χ_{opt}^2 for each object selected with the two color cut adopted throughout this work assuming a detection SNR > 3 and a blue non-detection SNR < 2 . We separate them based on their true redshifts and measure the distribution of the results, presented in Figure 16. We note that since there are no blue rejection filters for F070W dropouts, their χ_{opt}^2 cannot be calculated. With F090W dropouts, there are too few blue bands to find a clear delineation between high and low redshift mock galaxies. However, for the redder F115W dropout band, 10 - 30% of interloper galaxies could be reliably rejected without affecting the number of high redshift galaxies selected by adopting a $\chi_{\text{opt}}^2 \gtrsim 5$. This cutoff is relatively unaffected by an increase the detection SNR.

3.7. Brown Dwarf Interlopers

In addition to the mock galaxies in the JAGUAR catalog, we also explored how ultracool brown dwarf stars may be selected as dropout candidates following the work of Wilkins et al. (2014), Finkelstein et al. (2015), and Ryan & Reid (2016). Brown dwarfs have stellar spectra that become redder at cooler temperatures, with stronger molecular absorption features that can mimic the red dropout colors of high-redshift galaxies. While these studies conclude that ultracool dwarfs will be relatively rare ($\sim 1 \text{ arcmin}^{-2}$), extended deep *JWST* surveys will likely contain a number of dwarfs due to the 9.7 arcmin^2 FOV of NIRCам. To that end, we used a subsample of the published spectra for L and T dwarfs from the SpeX Prism Spectral Library⁵ and calculated the fluxes of these objects through the NIRCам wide filters: F070W, F090W, F115W, F150W, and F200W, as these spectra have wavelength coverage to $2.5 \mu\text{m}$. We supplemented these observational data with a set of L and T dwarf model spectra (which extend to $50 \mu\text{m}$) from Sonora18 (Marley et al. in prep) at a range of surface temperatures ($T = 200 - 2300 \text{ K}$) and a fixed surface gravity of $\log(g) = 5.0$. For these model spectra, we calculated the NIRCам fluxes through the NIRCам wide filters F070W, F090W, F115W, F150W, and F200W, F277W, F356W, and F410M and the NIRCам medium filters F335M and F410M. We simulate these real and model objects at a range of distances between 0.1 kpc and 40 kpc, and add noise at the DEEP *HST*+NIRCам survey depth. We note that brown dwarfs are unresolved in extragalactic surveys and a stellarity parameter has been used to remove these sources from deep *HST* catalogs (see Section 3.5.1. in Bouwens et al. 2015). We do not simulate this in our current work, and caution that while morphology can be used for rejecting stellar contaminants, compact high-redshift galaxies may also be similarly unresolved.

We plot the positions of the noisy brown dwarf candidates on the F070W, F090W, and F115W dropout color-color plots in Figure 17, where we impose the same red filter detection SNR (> 3) and blue filter non-detection SNR (< 2). In the figure, the colors of the points indicate the optical spectral type of the object given in the SpeX Library (as shown in the color bar on the right) or estimated from the temperature of the Sonora18 model spectrum, and the size of the point indicates the simulated distance of the brown dwarf. While brown dwarfs do not have colors similar to F070W dropouts, a large population of brown dwarfs would be selected as F090W dropouts, and a smaller number would be identified as F115W dropouts.

For the F090W dropouts there is a general trend between optical spectral type and redder F090W - F115W and F115W - F150W colors, and most of the sources selected are at larger distances ($> 10 \text{ kpc}$), echoing results from Ryan & Reid (2016) demonstrating that *JWST* will be able to detect brown dwarfs in the Milky Way halo. A color selection at $F115W - F150W < 0.3$ would select against many late L and T dwarfs. Late T-dwarfs have very red F090W - F115W colors and blue F115W - F150W colors, and would also be selected as F090W dropouts. To aid in differentiating true high-redshift galaxies from brown dwarfs, we used the *Spitzer* IRAC photometry for a sample of 86 late M, L, and T dwarfs provided by Patten et al. (2006). After converting the Channel 1 ($3.6 \mu\text{m}$) and Channel 2 ($4.5 \mu\text{m}$) fluxes to AB magnitudes, we find that M and L dwarfs have $[3.6] - [4.5] < -0.3$ (roughly analogous to NIRCам F356W - F444W < -0.3), which is significantly bluer than the bulk of true F090W dropout galaxies. T dwarfs in the Patten et al. sample, however, have red $[3.6] - [4.5]$ colors and are not as easily separated from F090W dropouts. To find methods for removing T dwarfs from F090W dropout samples, we looked at the long-wavelength NIRCам colors of these stars using the Sonora18 model spectra. In Figure 18, we plot the F090W - F115W color vs. F335M - F356W color for both JAGUAR mock galaxies and Sonora18 model brown dwarfs with F115W, F150W SNR > 3 and F435W, F606W, and F070W SNR < 2 (*HST*+NIRCам) with the model brown dwarf points colored as they are in Figure 17. From this Figure, we show that Late L and all T dwarfs can be reliably separated from true F090W dropouts by requiring a color cut at $F335M - F356W < 0.75$ (black vertical line), along with the F090W - F115W color cut (lavender horizontal dashed line), although this potentially removes a small population of high-redshift mock galaxies with strong optical line emission.

It is important to note that we do not simulate the on-sky density of objects in our sample in these plots. Ryan & Reid (2016) explored the actual number density of brown dwarfs in the thick and thin disk of the Milky Way, and concluded that only a few ultracool dwarfs would contaminate extragalactic surveys. Our results demonstrate the importance in the use of a Two Color or Compound Color method as well as observations at longer wavelengths to help mitigate this contamination.

⁵ Compiled by Adam Burgasser and found online at <http://pono.ucsd.edu/~adam/browndwarfs/spexprism/>

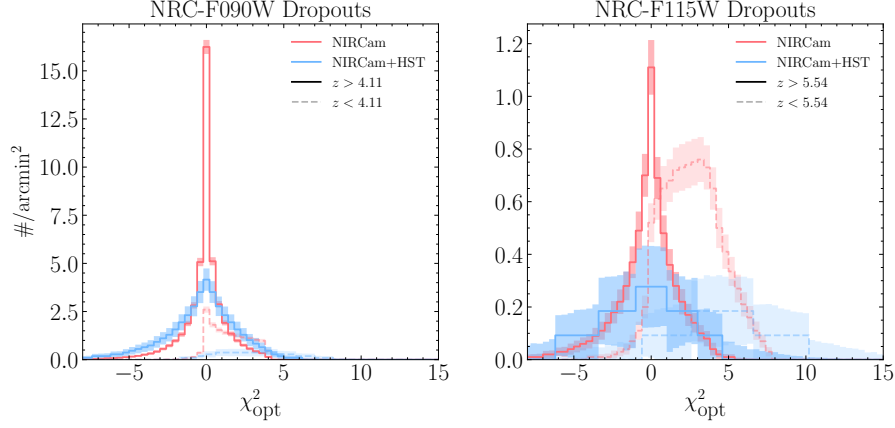


Figure 16. χ^2_{opt} distributions for F090W dropouts (left), F115W dropouts (middle), and F150W dropouts (right) for high-redshift objects (solid) and low redshift interlopers (dashed) in NIRC+HST (red) and *HST*+NIRC+HST (blue) surveys. The use of χ^2_{opt} is more effective for higher redshift dropout galaxies where there are more photometric bands for a given mock galaxy at wavelengths shorter than the Lyman break. For the reddest dropout bands, a $\chi^2_{\text{opt}} \gtrsim 5$ would accurately reject outliers without impacting the selection of true high redshift galaxies.

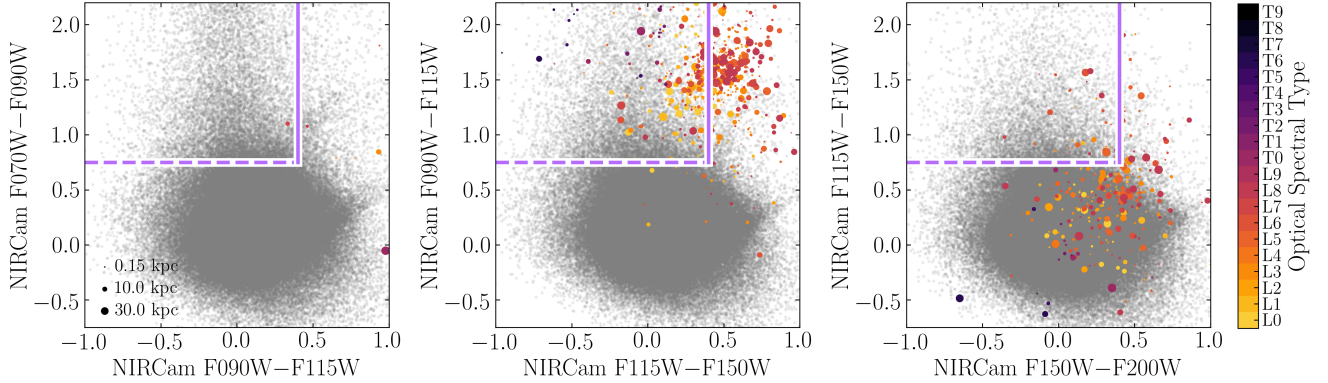


Figure 17. NIRC+HST color-color plots with the mock galaxies from Figure 3 plotted with grey points, overplotted with a selection of brown dwarfs, at different simulated distances and with at the DEEP *HST*+NIRC+HST survey depth. The points are colored by their optical spectral type, as shown on the color bar, and the sizes of the points indicate the distance from Earth, as shown in the bottom-left corner of the leftmost figure.

3.8. Dusty Star-Forming Galaxies

The JAGUAR catalog contains only a limited population of highly dust obscured star-forming mock galaxies as the mass and luminosity functions that were used to create JAGUAR are dependent on observations of the rest-frame optical and UV portions of a galaxy’s spectrum and are therefore missing extremely dusty galaxies. To further explore how dust affects dropout selection, we reproduced the entire catalog of JAGUAR star-forming mock galaxies, keeping the properties including mass, observed redshift, and star formation history the same, but assigning a random extinction value (parameterized by the color difference $E(B-V)$) to each object between $E(B-V) = 0-2$. To recreate these objects, we used the Flexible Stellar Population Synthesis code (FSPS, Conroy et al. 2009; Conroy & Gunn 2010), and used Padova isochrones along with the MILES spectral library (Sánchez-Blázquez et al. 2006). We chose to model the dust as a foreground screen using the Calzetti et al. (2000) prescrip-

tion at $z < 3$, and with SMC-bar-like dust (Gordon et al. 2003) at $z > 3$. While assigning a random $E(B-V)$ to all of the JAGUAR mock galaxies is nonphysical in light of observed trends between stellar mass and $E(B-V)$ out to $z \sim 6$ (Schaerer & de Barros 2010), as well as the complexity of actual dust geometry within galaxies, these extreme values for dust extinction will allow us to observe how obscuration affects the NIRC+HST colors of a diverse population of low-redshift interlopers.

In Figure 19, we plot the noise-free F070W, F090W, and F115W dropout space of the dusty star-forming mock galaxies, with the points colored by $E(B-V)$, as shown in the color bar on the right side of the figure. We indicate true high-redshift mock galaxies by blue circles. In each panel, dustier mock galaxies are found in a sequence that extends upwards and to the right. We also plot example two-color selection boxes, demonstrating the use of a second color cut to exclude lower redshift dusty interlopers. The lack of dusty

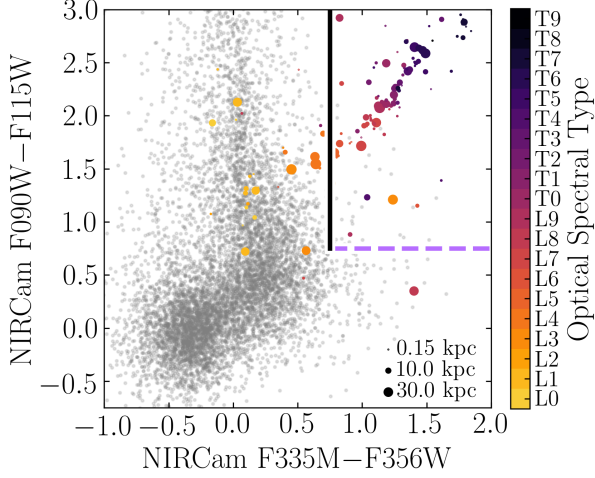


Figure 18. NIRCcam F090W - F115W color vs. F335M - F356W color plot. In grey, we plot JAGUAR mock galaxies in a DEEP survey with F115W, F150W SNR > 3 and F435W, F606W, and F070W SNR < 2 at $z > 5.54$, the redshift demarcation for F090W dropouts. The colored points are Sonora18 model brown dwarfs with optical spectral type as shown with the color bar on the right side, with the size of the markers symbolizing the distance from the Earth as given in the bottom right of the plot. While the central panel in Figure 17 shows that a quantity of L and T dwarfs contaminate F090W dropout selection, many of the late L and T dwarfs that satisfy the same F090W - F115W color cut (dashed lavender line) can be removed by also requiring F335M - F356W < 0.75 (black line).

mock galaxies that fall inside the selection boxes demonstrates the bias against selecting obscured galaxies using the Lyman dropout technique.

4. COMPARISON TO THE EMPIRICAL GALAXY GENERATOR

In anticipation of future deep extragalactic surveys, the ASTRODEEP collaboration developed the Empirical Galaxy Generator (EGG, Schreiber et al. 2017)⁶, which constructs mock catalogs including both photometry and morphologies. Similar to JAGUAR, EGG uses empirical prescriptions, starting with a derivation of the evolution of the stellar mass function from deep observations. In this section we compare the recovered TPD, completeness, and accuracy for the EGG catalog to what we found using the JAGUAR catalog. A few of the primary differences between JAGUAR and EGG that will influence the present analysis are the evolution of the stellar mass function, the treatment of galaxy morphologies and dust obscuration, and the inclusion of self-consistent nebular continuum and line emission.

The EGG team started with a framework for the evolution of the star-forming and quiescent galaxy mass function at

$z = 0.3 - 4.5$ based on observations from CANDELS (Grogin 2011; Koekemoer 2011), where they computed photometric redshifts using EAZY (Brammer et al. 2008) and galaxy stellar masses using FAST (Kriek et al. 2009). At $z = 4.5 - 7.5$, the authors rely on the stellar mass functions from Grazian et al. (2015). The resulting mass function evolution has a steeper low-mass slope than the prescription that underpins the JAGUAR catalog at $z > 1.5$ and the discrepancy is larger at higher redshifts. In addition, the evolution of the EGG mass function predicts fewer high-mass galaxies at $z > 4$ than JAGUAR. Both of these differences are likely a consequence of the necessary extrapolation that was done for each catalog due to lack of observational data.

The SEDs in EGG were generated by first assigning a $U - V$ and $V - J$ color to each mock galaxy based on the observed evolution of these colors for star-forming and quiescent galaxies. At this point, each mock galaxy was given an SED based on the average SED for observed CANDELS galaxies with those UVJ colors (from the FAST fits, using the Bruzual & Charlot (2003) stellar library). As the morphology of each EGG mock galaxy is defined to be a combination of a bulge and disk component, each component was assigned a separate SED. This process differs significantly from the SED creation in JAGUAR, which uses BEAGLE fits to 3D-HST objects to calculate the SEDs for each object. In the version of the EGG catalog generation tool we used in this analysis, v1.4.0 (*egg-gencat*), the authors included a simple prescription for emission lines, where the strength of each line is estimated using each mock galaxy’s SFR, metallicity, total infrared luminosity, and gas mass⁷, which we include to better compare to JAGUAR.

We used *egg-gencat* to create two catalogs, one with 100 square arcminutes and one with 10.8 square arcminutes, with a minimum stellar mass of $10^6 M_\odot$, at $z = 0.2 - 15$. We then constructed 500 noisy catalogs with each area in the exact manner as was done in Section 2.2 for the JAGUAR catalogs, although we modified this process to account for the combination of the disk and bulge components in each EGG mock galaxy. From these noisy catalogs, we measured the TPD, completeness, and accuracy as a function of color cuts, SNR, and survey depth following the analysis we performed for the noisy JAGUAR catalogs.

In Figure 20, we plot the TPD, completeness, and accuracy as a function of Lyman break color cut at the DEEP survey depth, with a detection SNR > 3 , and the Two Color Cut Scheme where the UV continuum cut is set at 0.4 for all three selection methods. We compare the EGG values to the JAGUAR values in each panel. For F070W dropouts, we find that EGG predicts almost double the TPD, and increased completeness, at all colors, but at significantly reduced accuracy, likely a result of the increased number of low-mass, faint galaxies in the EGG catalog. For F090W dropouts, the predicted TPD is more comparable between the EGG and the JAGUAR results, although the abundance of low-mass mock

⁶ <https://cschreib.github.io/egg/>

⁷ see: <https://github.com/cschreib/egg/blob/master/CHANGELOG>

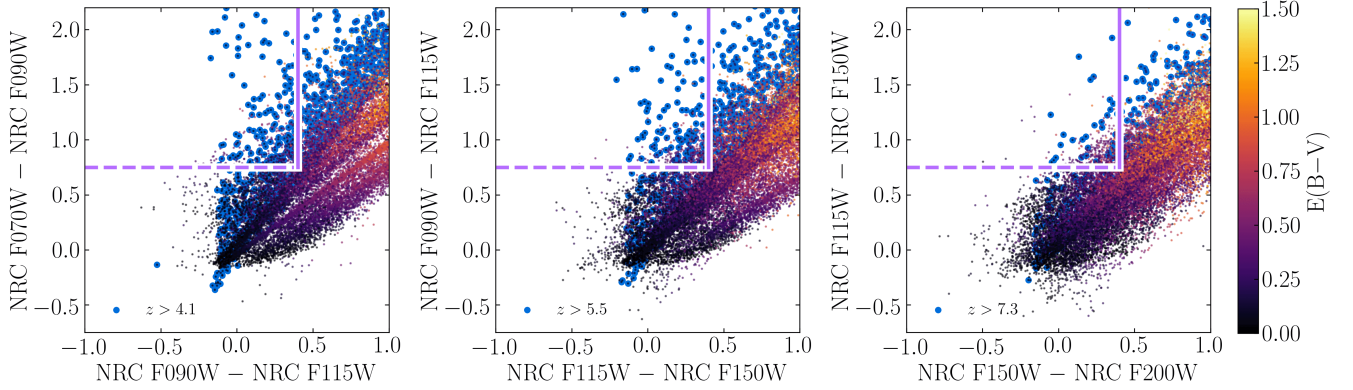


Figure 19. NIRCcam color-color plots as in Figure 2, but with mock galaxies with randomly assigned values for the extinction between $E(B-V)=0-2$. For mock galaxies at $z < 3$, we use the Calzetti et al. (2000) dust prescription, and at $z > 3$, we use the SMC-bar-like dust from Gordon et al. (2003). We represent actual high-redshift dropout mock galaxies in each panel with blue points. Dust causes the mock galaxy colors to be redder along each axis of the figure, reinforcing the need for a second color cut (such as the example vertical solid lavender line in each panel) to remove dusty low-redshift interlopers.

galaxies means that the predicted accuracy is significantly lower for EGG than for JAGUAR. For F115W dropouts, the EGG catalog results in higher completeness values but at very low accuracies.

These results highlight the difficulty in extrapolating the galaxy stellar mass function to low mass, especially at high redshift. Future deep *JWST* surveys will include a significant number of low-mass ($< 10^7 M_\odot$, see Figure 6) galaxies that are currently too faint to be observed with existing instruments, which will allow for a more robust measurement of the low-mass slope of the mass function.

5. DISCUSSION

JWST/NIRCam will be an exceptional instrument for selecting large samples of high-redshift galaxies. The NIRCam field of view spans 9.7 arcmin^2 , so with a single pointing as part of a DEEP survey, these results indicate that it may be possible to select ~ 900 $z = 5-7$ F070W dropouts at greater than 70% accuracy, ~ 500 $z = 6-9$ F090W dropouts at greater than 70% accuracy, and ~ 60 $z = 8-11$ F115W dropouts at greater than 50% accuracy. In Figure 6, we additionally show that these samples include mock galaxies at stellar masses of $10^6 - 10^7 M_\odot$, a stellar mass range that is difficult to access at these redshifts.

These results also demonstrate the importance of targeting fields with existing deep *HST* data. The addition of *HST* photometry has a limited effect in reducing the number of recovered high-redshift galaxies, but serves to significantly increase dropout selection accuracy, especially for blue color limits. While the majority of existing deep *HST* data are not at the XDF depths, the usage of shallower photometry should produce a result that is bracketed by the NIRCam only and *HST*+NIRCam results we have presented. Additional exploration of alternate survey designs at different *HST* and NIRCam depths can be accomplished using the NIRCPrepare-Mock software package described in Section 2.2.

A selection scheme should be chosen based on the sample purity and accuracy requirements for the science goal of a given survey. For F070W, F090W, and F115W dropouts, with two color selections, we find that a color cut of ~ 1 balances the trade-off between completeness (and TPD) and accuracy. Because of how the presence of significant quantities of dust affects the observed NIRCcam colors of low redshift interloper galaxies, we also recommend additional color cuts to focus on galaxies with blue UV slopes. For more inclusive (bluer) color cuts, the Compound Color Cut method provides the highest accuracy levels observed, cutting out a quantity of low redshift interloper galaxies. In addition, when existing deep optical data exist, we recommend the use of the χ^2_{opt} statistic to help increase survey accuracy, but with the caveat that there is significant overlap between the χ^2_{opt} distributions for true high-redshift galaxies and lower redshift interlopers.

6. CONCLUSIONS

We have demonstrated how *JWST*/NIRCam colors can be used to select samples of high-redshift galaxies through the Lyman dropout technique by simulating surveys with the JAGUAR catalog. We examined how the exact color cut affected the on-sky TPD of high-redshift sources, and the completeness and accuracy of the resulting sample. Our primary results are:

1. At increasingly redder color cuts, in all cases, the true positive density and completeness of the recovered sample decreases, as fewer objects are selected, but the overall accuracy increases, often plateauing at a maximum accuracy level for a given set of mock galaxies. This process is primarily driven by SNR, as the recovered mock galaxies in each selection method we explored are at higher redshifts than the interlopers, which are scattered into the selection box due to noise.
2. We explored three dropout selection methods which utilize either one single color cut, two color cuts, or a

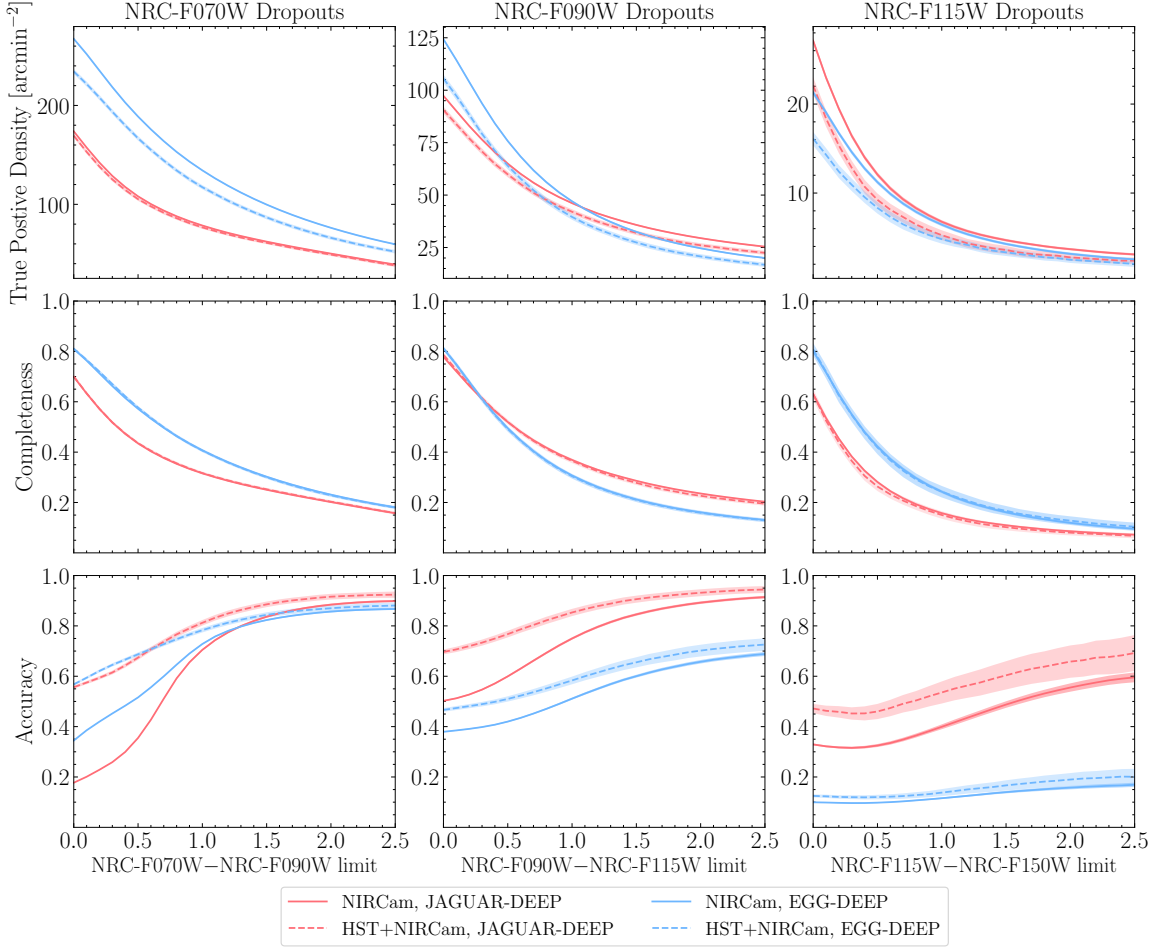


Figure 20. TPD (top row) and accuracy (bottom row) as a function of color cut for F070W dropouts (left), F090W dropouts (middle), and F115W dropouts (right) in the JAGUAR (red) and EGG (blue) catalogs for NIRCcam (solid) and *HST*+NIRCcam (dashed) surveys. In all three cases, we are using the DEEP survey depth, with a detection SNR > 3, and a two-color selection, with a UV continuum color cut of < 0.4. The evolution of the mass function in the EGG catalog predicts a much larger number of low-mass galaxies than the JAGUAR galaxies, and less high-redshift galaxies at all masses. This leads to a higher TPD for F070W and F090W dropouts, along with a much lower accuracy for F090W and F115W dropouts, where more faint low mass galaxies serve as high-redshift interlopers.

compound method with three color cuts. While the use of a second or third color cut significantly increases the accuracy of the recovered sample, it is at the expense of the TPD and completeness, especially for F070W dropouts. This is less true for F115W dropout selection due to the limited number of mock galaxies predicted to be observed per square arcminute in a given survey. Dust obscuration has the effect of moving mock galaxies to redder colors, reinforcing the usage of two or three color selection criteria for selecting against low-redshift dusty interlopers.

3. We find that the presence of Ly α emission has a complicated effect on the TPD, completeness, and accuracy of recovered samples as a function of color cut due to how the emission line contributes to flux in different bands at different redshifts. For F090W and F115W dropouts, we predict a higher TPD, complete-

ness, and accuracy for samples without Ly α emission at most moderate color selection cuts.

4. Filter selection for the UV continuum color cut plays a significant role in recovering galaxies. For F070W dropouts, it is recommended to use filters that probe longer wavelengths to increase the TPD, completeness, and accuracy of the recovered sample. For F090W and F115W dropouts, using redder filters in the UV continuum color cut leads to an increase in accuracy at a lower TPD and completeness. In addition, we recommend using the long wavelength NIRCcam filters to assist in rejecting interloper galaxies for all three dropout schemes.
5. The usage of the χ^2_{opt} statistic (Bouwens et al. 2015) for removing low-redshift interloper galaxies is only recommended for samples with deep observations in multiple filters at wavelengths shorter than the Lyman

break for the dropout sample, where assuming $\chi_{\text{opt}}^2 < 5$ would help to reject outliers.

6. While it will be possible for NIRCam to detect brown dwarf stars out to 10+ kpc, they are only a significant source of contamination in F090W dropout selection, which can be alleviated with color cuts using NIRCam data at 3 - 5 μm .

ACKNOWLEDGMENTS

ECL acknowledges support from the ERC Advanced Grant 695671 “QUENCH.” CCW acknowledges support from the National Science Foundation Astronomy and Astrophysics Fellowship grant AST-1701546.

REFERENCES

- Bouwens, R. J., Illingworth, G. D., Franx, M., & Ford, H. 2007, *ApJ*, 670, 928
- . 2008, *ApJ*, 686, 230
- Bouwens, R. J., Illingworth, G. D., Rosati, P., et al. 2003, *ApJ*, 595, 589
- Bouwens, R. J., Thompson, R. I., Illingworth, G. D., et al. 2004, *ApJL*, 616, L79
- Bouwens, R. J., Illingworth, G. D., Oesch, P. A., et al. 2015, *ApJ*, 803, 34
- Brammer, G. B., van Dokkum, P. G., & Coppi, P. 2008, *ApJ*, 686, 1503
- Bruzual, G., & Charlot, S. 2003, *MNRAS*, 344, 1000
- Bunker, A. J., Stanway, E. R., Ellis, R. S., & McMahon, R. G. 2004, *MNRAS*, 355, 374
- Bunker, A. J., Stanway, E. R., Ellis, R. S., McMahon, R. G., & McCarthy, P. J. 2003, *MNRAS*, 342, L47
- Bunker, A. J., Wilkins, S., Ellis, R. S., et al. 2010, *MNRAS*, 409, 855
- Calzetti, D., Armus, L., Bohlin, R. C., et al. 2000, *ApJ*, 533, 682
- Caruana, J., Bunker, A. J., Wilkins, S. M., et al. 2012, *MNRAS*, 427, 3055
- . 2014, *MNRAS*, 443, 2831
- Casey, C. M., Narayanan, D., & Cooray, A. 2014, *PhR*, 541, 45
- Cassata, P., Tasca, L. A. M., Le Fèvre, O., et al. 2015, *A&A*, 573, A24
- Charlot, S., & Fall, S. M. 2000, *ApJ*, 539, 718
- Chevallard, J., & Charlot, S. 2016, *MNRAS*, 462, 1415
- Conroy, C., & Gunn, J. E. 2010, *ApJ*, 712, 833
- Conroy, C., Gunn, J. E., & White, M. 2009, *ApJ*, 699, 486
- De Barros, S., Oesch, P. A., Labbé, I., et al. 2019, *MNRAS*, 489, 2355
- Dunlop, J. S. 2013, in *Astrophysics and Space Science Library*, Vol. 396, *The First Galaxies*, ed. T. Wiklind, B. Mobasher, & V. Bromm, 223
- Ellis, R. S., McLure, R. J., Dunlop, J. S., et al. 2013, *ApJL*, 763, L7
- Finkelstein, S. L., Ryan, Russell E., J., Papovich, C., et al. 2015, *ApJ*, 810, 71
- Frye, B., Broadhurst, T., & Benítez, N. 2002, *ApJ*, 568, 558
- Gardner, J. P., Mather, J. C., Clampin, M., et al. 2006, *SSRv*, 123, 485
- Giallongo, E., Grazian, A., Fiore, F., et al. 2015, *A&A*, 578, A83
- Giavalisco, M., Koratkar, A., & Calzetti, D. 1996, *ApJ*, 466, 831
- Gordon, K. D., Clayton, G. C., Misselt, K. A., Landolt, A. U., & Wolff, M. J. 2003, *ApJ*, 594, 279
- Grazian, A., Fontana, A., Santini, P., et al. 2015, *A&A*, 575, A96
- Grogin, N. A. e. a. 2011, *ApJS*, 197, 35
- Guhathakurta, P., Tyson, J. A., & Majewski, S. R. 1990, *ApJL*, 357, L9
- Guo, Y., McIntosh, D. H., Mo, H. J., et al. 2009, *MNRAS*, 398, 1129
- Illingworth, G. D., Magee, D., Oesch, P. A., et al. 2013, *ApJS*, 209, 6
- Inami, H., Bacon, R., Brinchmann, J., et al. 2017, *A&A*, 608, A2
- Infante, L., Zheng, W., Laporte, N., et al. 2015, *ApJ*, 815, 18
- Inoue, A. K., Shimizu, I., Iwata, I., & Tanaka, M. 2014, *MNRAS*, 442, 1805
- Ishigaki, M., Kawamata, R., Ouchi, M., et al. 2015, *ApJ*, 799, 12
- . 2018, *ApJ*, 854, 73
- Koekemoer, A. M. e. a. 2011, *ApJS*, 197, 36
- Kriek, M., van Dokkum, P. G., Labbé, I., et al. 2009, *ApJ*, 700, 221
- Kunth, D., Mas-Hesse, J. M., Terlevich, E., et al. 1998, *A&A*, 334, 11
- Labbé, I., Oesch, P. A., Bouwens, R. J., et al. 2013, *ApJL*, 777, L19
- Lorenzoni, S., Bunker, A. J., Wilkins, S. M., et al. 2013, *MNRAS*, 429, 150
- . 2011, *MNRAS*, 414, 1455
- Madau, P., & Haardt, F. 2015, *ApJL*, 813, L8
- McLeod, D. J., McLure, R. J., & Dunlop, J. S. 2016, *MNRAS*, 459, 3812
- McLure, R. J., Dunlop, J. S., Cirasuolo, M., et al. 2010, *MNRAS*, 403, 960
- McLure, R. J., Dunlop, J. S., Bowler, R. A. A., et al. 2013, *MNRAS*, 432, 2696
- Neufeld, D. A. 1991, *ApJL*, 370, L85
- Oesch, P. A., Bouwens, R. J., Illingworth, G. D., et al. 2015a, *ApJ*, 808, 104
- Oesch, P. A., Bouwens, R. J., Illingworth, G. D., Labbé, I., & Stefanon, M. 2018, *ApJ*, 855, 105
- Oesch, P. A., Bouwens, R. J., Illingworth, G. D., et al. 2013, *ApJ*, 773, 75
- . 2014, *ApJ*, 786, 108

- Oesch, P. A., van Dokkum, P. G., Illingworth, G. D., et al. 2015b, *ApJL*, 804, L30
- Oke, J. B., & Gunn, J. E. 1983, *ApJ*, 266, 713
- Ono, Y., Ouchi, M., Mobasher, B., et al. 2012, *ApJ*, 744, 83
- Onoue, M., Kashikawa, N., Willott, C. J., et al. 2017, *ApJL*, 847, L15
- Patten, B. M., Stauffer, J. R., Burrows, A., et al. 2006, *ApJ*, 651, 502
- Pentericci, L., Fontana, A., Vanzella, E., et al. 2011, *ApJ*, 743, 132
- Pentericci, L., Vanzella, E., Fontana, A., et al. 2014, *ApJ*, 793, 113
- Roberts-Borsani, G. W., Bouwens, R. J., Oesch, P. A., et al. 2016, *ApJ*, 823, 143
- Robertson, B. E., Ellis, R. S., Furlanetto, S. R., & Dunlop, J. S. 2015, *ApJL*, 802, L19
- Ryan, R. E., J., & Reid, I. N. 2016, *AJ*, 151, 92
- Salmon, B., Coe, D., Bradley, L., et al. 2018, *ApJL*, 864, L22
- Sánchez-Blázquez, P., Peletier, R. F., Jiménez-Vicente, J., et al. 2006, *MNRAS*, 371, 703
- Schaerer, D., & de Barros, S. 2010, *A&A*, 515, A73
- Schenker, M. A., Ellis, R. S., Konidaris, N. P., & Stark, D. P. 2014, *ApJ*, 795, 20
- Schenker, M. A., Stark, D. P., Ellis, R. S., et al. 2012, *ApJ*, 744, 179
- Schreiber, C., Elbaz, D., Pannella, M., et al. 2017, *A&A*, 602, A96
- Shapley, A. E., Steidel, C. C., Pettini, M., & Adelberger, K. L. 2003, *ApJ*, 588, 65
- Shen, S., Mo, H. J., White, S. D. M., et al. 2003, *MNRAS*, 343, 978
- Shibuya, T., Kashikawa, N., Ota, K., et al. 2012, *ApJ*, 752, 114
- Skelton, R. E., Whitaker, K. E., Momcheva, I. G., et al. 2014, *ApJS*, 214, 24
- Smit, R., Bouwens, R. J., Franx, M., et al. 2015, *ApJ*, 801, 122
- Song, M., Finkelstein, S. L., Livermore, R. C., et al. 2016, *ApJ*, 826, 113
- Spilker, J. S., Marrone, D. P., Aravena, M., et al. 2016, *ApJ*, 826, 112
- Stanway, E. R., Bunker, A. J., McMahon, R. G., et al. 2004, *ApJ*, 607, 704
- Stark, D. P. 2016, *ARA&A*, 54, 761
- Stark, D. P., Ellis, R. S., Chiu, K., Ouchi, M., & Bunker, A. 2010, *MNRAS*, 408, 1628
- Steidel, C. C., Adelberger, K. L., Giavalisco, M., Dickinson, M., & Pettini, M. 1999, *ApJ*, 519, 1
- Steidel, C. C., Adelberger, K. L., Shapley, A. E., et al. 2003, *ApJ*, 592, 728
- Steidel, C. C., Giavalisco, M., Pettini, M., Dickinson, M., & Adelberger, K. L. 1996, *ApJL*, 462, L17
- Tasca, L. A. M., Le Fèvre, O., Ribeiro, B., et al. 2017, *A&A*, 600, A110
- Tilvi, V., Papovich, C., Finkelstein, S. L., et al. 2014, *ApJ*, 794, 5
- Tomczak, A. R., Quadri, R. F., Tran, K.-V. H., et al. 2014, *ApJ*, 783, 85
- Treu, T., Schmidt, K. B., Trenti, M., Bradley, L. D., & Stiavelli, M. 2013, *ApJL*, 775, L29
- van der Wel, A., Franx, M., van Dokkum, P. G., et al. 2014, *ApJ*, 788, 28
- Vanzella, E., Giavalisco, M., Dickinson, M., et al. 2009, *ApJ*, 695, 1163
- Vanzella, E., Pentericci, L., Fontana, A., et al. 2011, *ApJL*, 730, L35
- Wang, T., Schreiber, C., Elbaz, D., et al. 2019, *Nature*, 572, 211
- Wilkins, S. M., Bunker, A. J., Stanway, E., Lorenzoni, S., & Caruana, J. 2011, *MNRAS*, 417, 717
- Wilkins, S. M., Stanway, E. R., & Bremer, M. N. 2014, *MNRAS*, 439, 1038
- Williams, C. C., Curtis-Lake, E., Hainline, K. N., et al. 2018, *ApJS*, 236, 33
- Williams, C. C., Labbe, I., Spilker, J., et al. 2019, *arXiv e-prints*, arXiv:1905.11996
- Zitrin, A., Zheng, W., Broadhurst, T., et al. 2014, *ApJL*, 793, L12