

# Detecting multiple DLAs per spectrum in SDSS DR12 with Gaussian processes

Ming-Feng Ho,<sup>1</sup>★ Simeon Bird<sup>1</sup> and Roman Garnett<sup>2</sup> 

<sup>1</sup> Department of Physics & Astronomy, University of California, Riverside, 900 University Avenue, Riverside, CA 92521, USA

<sup>2</sup> Department of Computer Science and Engineering, Washington University in St. Louis, One Brookings Drive, St. Louis, MO 63130, USA

Accepted 2020 June 16. Received 2020 May 27; in original form 2020 April 1

## ABSTRACT

We present a revised version of our automated technique using Gaussian processes (GPs) to detect damped Lyman  $\alpha$  absorbers (DLAs) along quasar (QSO) sightlines. The main improvement is to allow our GP pipeline to detect multiple DLAs along a single sightline. Our DLA detections are regularized by an improved model for the absorption from the Lyman  $\alpha$  forest that improves performance at high redshift. We also introduce a model for unresolved sub-DLAs that reduces misclassifications of absorbers without detectable damping wings. We compare our results to those of two different large-scale DLA catalogues and provide a catalogue of the processed results of our GP pipeline using 158 825 Lyman  $\alpha$  spectra from SDSS data release 12. We present updated estimates for the statistical properties of DLAs, including the column density distribution function, line density (dN/dX), and neutral hydrogen density ( $\Omega_{\text{DLA}}$ ).

**Key words:** methods: statistical – intergalactic medium – quasars: absorption lines – galaxies: statistics.

## 1 INTRODUCTION

Damped Ly  $\alpha$  absorbers (DLAs) are absorption line systems with high neutral hydrogen column densities ( $N_{\text{H I}} > 10^{20.3} \text{ cm}^{-2}$ ) discovered in sightlines of quasar (QSO) spectroscopic observations (Wolfe et al. 1986). The gas that gives rise to DLAs is dense enough to be self-shielded from the ultraviolet background (UVB; Cen 2012) yet diffuse enough to have a low star formation rate (Fumagalli et al. 2014). DLAs dominate the neutral-gas content of the Universe after reionization (Gardner et al. 1997; Noterdaeme et al. 2012; Zafar et al. 2013; Crighton et al. 2015). Simulations tell us DLAs are connected with galaxies over a wide range of halo masses (Prochaska & Wolfe 1997; Haehnelt, Steinmetz & Rauch 1998; Pontzen et al. 2008), and at  $z \geq 2$  are formed from the accretion of neutral hydrogen gas on to dark matter haloes (Bird et al. 2014, 2015). The abundance of DLAs at different epochs of the Universe ( $2 < z < 5$ ) thus becomes a powerful probe to understand the formation history of galaxies (Gardner et al. 1997; Wolfe, Gawiser & Prochaska 2005).

Finding DLAs historically involves a combination of template fitting and visual inspection of spectra by the eyes of trained astronomers (Prochaska, Herbert-Fort & Wolfe 2005; Slosar et al. 2011). Recent spectroscopic surveys such as the Sloan Digital Sky Survey (SDSS) (York et al. 2000) have taken large amount of QSO spectra (Pâris et al. 2012, 2014) ( $\sim 500\,000$  in SDSS-IV, Pâris et al. 2018). Future surveys such as the Dark Energy Spectroscopic Instrument (DESI)<sup>1</sup> will acquire more than 1 million QSOs, making visual inspection of the spectra impractical. Moreover, the low signal-to-noise ratios of SDSS data make the task of detecting DLAs even harder, and induces noise related detection systematics. Since the

release of the SDSS DR14 QSO catalogue (Pâris et al. 2018), visual inspection is no longer performed on all QSO targets. A fully automated and statistically consistent method thus needs to be presented for current and future surveys.

We provide a catalogue of DLAs using SDSS DR12 with 158 825 QSO sightlines. We demonstrate that our pipeline is capable of detecting an arbitrary number of DLAs within each spectroscopic observation, which makes it suitable for future surveys. Furthermore, since our pipeline resides within the framework of Bayesian probability, we have the ability to make probabilistic statements about those observations with low signal-to-noise ratios. This property allows us to make probabilistic estimations of DLA population statistics, even with low-quality noisy data (Bird, Garnett & Ho 2017).

Other available searches of DLAs in SDSS include: a visual inspection survey (Slosar et al. 2011), visually guided Voigt-profile fitting (Prochaska et al. 2005; Prochaska & Wolfe 2009); and three automated methods: a template-fitting method (Noterdaeme et al. 2012), an unpublished machine-learning approach using Fisher discriminant analysis (Carithers 2012), and a deep-learning approach using a convolutional neural network (Parks et al. 2018). Although these methods have had some success in creating large DLA catalogues, they suffer from hard-to-control systematics due to reliance either on templates or black box training.

We present a revised version of our previous automated method based on a Bayesian model-selection framework (Garnett et al. 2017). In our previous model (Garnett et al. 2017), we built a likelihood function for the QSO spectrum, including the continuum and the non-DLA absorption, using Gaussian processes (GPs; Rasmussen & Williams 2005). The SDSS DR9 concordance catalogue was applied to learn the covariance of the GP model. In this paper, we use the effective optical depth of the Lyman-series forest to allow the mean model of the likelihood function to be

\* E-mail: mho026@ucr.edu

<sup>1</sup> <http://desi.lbl.gov>

adjustable to the mean flux of the QSO spectrum, which reduces the probability of falsely fitting high column density absorbers at high redshifts. We also improve our knowledge of low column density absorbers and build an alternative model for sub-DLAs, which are the H I absorbers with  $19.5 < \log_{10} N_{\text{H I}} < 20$ . These modifications allow us to extend our previous pipeline to detect an arbitrary number of DLAs within each QSO sightline without overfitting.

Alongside the revised DLA detection pipeline, we present the new estimates of DLA statistical properties at  $z > 2$ . Since the neutral hydrogen gas in DLAs will eventually accrete on to galactic haloes and fuel the star formation, these population statistics can give an independent constraint on the theory of galaxy formation. Our pipeline relies on a well-defined Bayesian framework and contains a full posterior density on the column density and redshift for a given DLA. We thus can properly propagate the uncertainty in the properties of each DLA spectrum to population statistics of the whole sample. Additionally, we are also able to account for low signal-to-noise ratio samples in our population statistics since the uncertainty will be reflected in the posterior probability. We thus substantially increase the sample size in our measurements by including these noisy observations.

## 2 NOTATION

We will briefly recap the notation we defined in Garnett et al. (2017). Imagine we are observing a QSO with a known redshift  $z_{\text{QSO}}$ . The underlying true emission function  $f(\lambda_{\text{rest}})$  ( $f: \mathcal{X} \rightarrow \mathbb{R}$ ) of the QSO is a mapping relation from rest-frame wavelength to flux. We will always assume the  $z_{\text{QSO}}$  is known and rescale the observed-frame wavelength  $\lambda_{\text{obs}}$  to the rest-frame wavelength with  $\lambda_{\text{rest}} (= \lambda_{\text{obs}}/(1 + z_{\text{QSO}}))$ . We will use  $\lambda$  to replace  $\lambda_{\text{rest}}$  in the rest of the text because we only work on  $\lambda_{\text{rest}}$ .

The QSO spectrum observed is not the intrinsic emission function  $f(\lambda)$ . Both the instrumental noise and absorption due to the intervening intergalactic medium along the line of sight will affect the observed flux. We thus denote the observed flux as a function  $y(\lambda)$ .

For a real spectroscopic observation, we measure the function  $y(\lambda)$  on a discrete set of samples  $\lambda$ . We thus denote the observed flux as a vector  $\mathbf{y}$ , which is defined as  $y_i = y(\lambda_i)$  with  $i$  representing  $i$ th pixel. For a given QSO observation, we use  $\mathcal{D}$  to represent a set of discrete observations  $(\lambda, \mathbf{y})$ .

We exclude missing values of the spectroscopic observations in our calculations. These missing values are due to pixel masking in the spectroscopic observations (e.g. bad columns in the CCD detectors). We will use NaN ('not a number') to represent those missing values in the text, and we will always ignore NaNs in the calculations.

## 3 BAYESIAN MODEL SELECTION

The classification approach used in our pipeline depends on Bayesian model selection. Bayesian model selection allows us to compute the probability that a spectroscopic sightline  $\mathcal{D}$  contains an arbitrary number of DLAs through evaluating the probabilities of a set of models  $\{\mathcal{M}_i\}$ , where  $i$  is a positive integer. This set of  $\mathcal{M}_i$  contains all potential models we want to classify: a model with no DLA and models having between one DLA and  $k$  DLAs.

For each  $\mathcal{M}_i$ , we want to compute the probability that best explains the data  $\mathcal{D}$  given a model  $\mathcal{M}$ . To do this, we have to marginalize the model parameters  $\theta$  and evaluate the model evidence,

$$p(\mathcal{D} | \mathcal{M}) = \int p(\mathcal{D} | \mathcal{M}, \theta) p(\theta | \mathcal{M}) d\theta. \quad (1)$$

Given a set of model pieces of evidence  $p(\mathcal{D} | \mathcal{M}_i)$  and model priors  $\text{Pr}(\mathcal{M}_i)$ , we are able to evaluate the posterior of a model given data based on Bayes's rule,

$$\text{Pr}(\mathcal{M} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathcal{M}) \text{Pr}(\mathcal{M})}{\sum_i p(\mathcal{D} | \mathcal{M}_i) \text{Pr}(\mathcal{M}_i)}. \quad (2)$$

We will select the model from  $\{\mathcal{M}_i\}$  with the highest posterior. Readers may think of this method as an application of Bayesian hypothesis testing. Instead of only getting the likelihoods conditioned on models, we get posterior probabilities for each model given data.

Let  $k$  be the maximum number of DLAs we will want to detect in a QSO spectrum. For our multi-DLA model selection, we will develop  $k + 2$  models, which include a null model for no DLA detection ( $\mathcal{M}_{\text{-DLA}}$ ), models for detecting exactly  $k$  DLAs ( $\mathcal{M}_{\text{DLA}(k)}$ ), and a model with sub-DLAs ( $\mathcal{M}_{\text{sub}}$ ). With a given spectroscopic sightline  $\mathcal{D}$ , we will compute the posterior probability of having exactly  $k$  DLAs in data  $\mathcal{D}$ ,  $\text{Pr}(\mathcal{M}_{\text{DLA}(k)} | \mathcal{D})$ .

## 4 GAUSSIAN PROCESSES

In this section, we will briefly recap how we use GPs to describe the QSO emission function  $f(\lambda)$ , following Garnett et al. (2017). The QSO emission function is a complicated function without a simple form derived from physically motivated parameters. We thus use a non-parametric framework, GPs, for modelling this physically unknown function  $f(\lambda)$ . A detailed introduction to GPs may be found in Rasmussen & Williams (2005).

### 4.1 Definition and prior distribution

We wish to use a GP to model the QSO emission function  $f(\lambda)$ . We can treat a GP as an extension of the joint Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  to infinite continuous domains. The difference is that a GP is a distribution over functions, not just a distribution over a finite number of random variables (although since we are dealing with pixelized variables here the distinction is less important).

A GP is completely specified by its first two central moments, a mean function  $\mu(\lambda)$  and a covariance function  $K(\lambda, \lambda')$ :

$$\begin{aligned} \mu(\lambda) &= \mathbb{E}[f(\lambda) | \lambda], \\ K(\lambda, \lambda') &= \mathbb{E}[(f(\lambda) - \mu(\lambda))(f(\lambda') - \mu(\lambda')) | \lambda, \lambda'] \\ &= \text{cov}[f(\lambda), f(\lambda') | \lambda, \lambda']. \end{aligned} \quad (3)$$

The mean vector describes the expected behaviour of the function, and the covariance function specifies the covariance between pairs of random variables. We thus will write the GP as

$$f(\lambda) \sim \mathcal{GP}(\mu(\lambda), K(\lambda, \lambda')). \quad (4)$$

We can write the prior probability distribution of a GP as

$$p(f) = \mathcal{GP}(f; \mu, K). \quad (5)$$

Real spectroscopic observations measure a discrete set of inputs  $\lambda$  and the corresponding  $f(\lambda)$ , so we get a multivariate Gaussian distribution

$$p(f) = \mathcal{N}(f(\lambda); \mu(\lambda), K(\lambda, \lambda')). \quad (6)$$

Assuming the dimension of  $\lambda$  and  $f$  is  $d$ , the form of the multivariate Gaussian distribution is written as

$$\begin{aligned} \mathcal{N}(f; \mu, \mathbf{K}) \\ = \frac{1}{\sqrt{(2\pi)^d \det \mathbf{K}}} \exp \left( -\frac{1}{2} (f - \mu)^\top \mathbf{K}^{-1} (f - \mu) \right). \end{aligned} \quad (7)$$

## 4.2 Observation model

We now have a GP model for a discrete set of wavelengths  $\lambda$  and true emission fluxes  $f$ . To build the likelihood function for observational data  $\mathcal{D} = (\lambda, y)$ , we have to incorporate the observational noise. Here, we assume the observational noise is modelled by an independent Gaussian variable for each wavelength pixel, allowing the noise realization to differ between pixels but neglecting inter-pixel correlations.

The noise variance for a given  $\lambda_i$  is written as  $v_i = \sigma(\lambda_i)^2$ .  $\sigma(\lambda_i)$  is the measurement error from a single observation on a given wavelength point  $\lambda$ . With the above assumptions, we can write down the mechanism of generating observations as

$$p(y | \lambda, f, v) = \mathcal{N}(y; f, V), \quad (8)$$

where  $V = \text{diag } v$ , which means we put the vector  $v$  on the diagonal terms of the diagonal square matrix  $V$ .

Given an observational model  $p(y | \lambda, f, v)$  and a GP emission model  $p(f | \lambda)$ , the prior distribution for observations  $y$  is obtained by marginalizing the latent function  $f$ :

$$\begin{aligned} p(y | \lambda, v) &= \int p(y | \lambda, f, v) p(f | \lambda) df \\ &= \int \mathcal{N}(y; f, V) \mathcal{N}(f; \mu, K) df \\ &= \mathcal{N}(y; \mu, K + V), \end{aligned} \quad (9)$$

where the Gaussians are closed under the convolution. Our observation model thus becomes a multivariate normal distribution described by a mean model  $\mu(\lambda)$ , covariance structure  $K(\lambda, \lambda')$ , and the instrumental noise  $V$ . The instrumental noise is derived from SDSS pipeline noise, so it is different from QSO-to-QSO; however, since  $K$  encodes the covariance structure of QSO emissions,  $K$  should be the same for all QSOs.

As explained in Garnett et al. (2017), there is no obvious choice for a prior covariance function  $K$  for modelling the QSO emission function. Most off-the-shelf covariance functions assume some sort of translation invariance, but this is not suitable for spectroscopic observations.<sup>2</sup> However, we understand the QSO emission function will be independent of the presence of a low-redshift DLA. We also assume that QSO emission functions are roughly redshift independent in the wavelength range of interest (Lyman limit to Lyman  $\alpha$ ), as accretion physics should not strongly vary with cosmological evolution. We thus build our own custom  $\mu$  and  $K$  for the GP prior to model the QSO spectra.

## 5 LEARNING A GP PRIOR FROM QSO SPECTRA

In this section, we will recap the prior modelling choices we made in Garnett et al. (2017) and the modifications we made to reliably detect multiple DLAs in one spectrum. We first build a GP model for QSO emission in the absence of DLAs, the null model  $\mathcal{M}_{\text{-DLA}}$ . Our model with DLAs ( $\mathcal{M}_{\text{DLA}}$ ) extends this null model. With the model priors and model evidence of all models we are considering, we compute the model posterior with Bayesian model selection.

The GP prior is completely described by the first two moments, the mean and covariance functions, which we derive from data. We must consider the mean flux of QSO emission, the absorption effect

due to the Lyman  $\alpha$  forest, and the covariance structure within the Lyman series.

### 5.1 Data

Our training set to learn our GP null model comprises the spectra observed by SDSS BOSS DR9 and labelled as containing (or not) a DLA by Lee et al. (2013).<sup>3</sup> The DR9 data set includes 54 468 QSO spectra with  $z_{\text{QSO}} > 2.15$ . We removed the following QSOs from the training set:

- (i)  $z_{\text{QSO}} < 2.15$ : QSOs with redshifts lower than 2.15 have no Lyman  $\alpha$  in the SDSS band.
- (ii) BAL: QSOs with broad absorption lines as flagged by the SDSS pipeline.
- (iii) Spectra with less than 200 detected pixels.
- (iv) ZWARNING: spectra whose analysis had warnings as flagged by the SDSS redshift estimation. Extremely noisy spectra (the TOO\_MANY\_OUTLIERS flag) were kept.

### 5.2 Modelling decisions

Consider a set of QSO observations  $\mathcal{D} = (\lambda, y)$ ; we always shift the observer's frame  $\lambda_{\text{obs}}$  to rest-frame  $\lambda$  so that we can set the emissions of Lyman series from different spectra to the same rest wavelengths. The assumption here is that the  $z_{\text{QSO}}$ s of QSOs are known for all the observed spectra, which is not precisely true for the spectroscopic data we have here. Accurately estimating the redshift of QSOs is beyond the scope of this paper, and is tackled elsewhere (Fauber et al. 2020).

The observed magnitude of a QSO varies considerably, based on its luminosity distance and the properties of the black hole. For the observation  $y$  to be described by a GP, it is necessary to normalize all flux measurements by dividing by the median flux observed between 1310 and 1325 Å, a wavelength region that is unaffected by the Lyman  $\alpha$  forest.

We model the same wavelength range as in Garnett et al. (2017):

$$\lambda \in [911.75\text{\AA}, 1215.75\text{\AA}], \quad (10)$$

going from the QSO rest-frame Lyman limit to the QSO rest-frame Lyman  $\alpha$ . The spacing between pixels is  $\Delta\lambda = 0.25$  Å. Note that we prefer not to include the region past the Lyman limit. This is partly due to the relatively small amount of data in that region and partly because the non-Gaussian Lyman break associated with Lyman limit systems can confuse the model. In particular, it occasionally tries to model a Lyman break with a wide DLA profile with a high column density. We shall see this is especially a problem if the QSO redshift is slightly inaccurate. The code considers the prior probability of a Lyman break at a higher redshift than the putative QSO rest frame to be zero and thus is especially prone to finding other explanations for the large absorption trough.

To model the relationship between flux measurements and the true QSO emission spectrum, we have to add terms corresponding to instrumental noise and weak Lyman  $\alpha$  absorption to the intrinsic correlations within the emission spectrum. Instrumental noise was already added in equation (9) as a matrix  $V$ .

The remaining part of the modelling is to define the GP covariance structure for QSOs across different redshifts. In Garnett et al. (2017), Lyman  $\alpha$  absorbers were modelled by a single additive noise term,  $\Omega$ ,

<sup>2</sup>Detailed explanations are in Garnett et al. (2017), section 4.2.1.

<sup>3</sup>However, we use the DR12 pipeline throughout.

accounting for the effect of the forest as extra noise in the emission spectrum. This is not completely physical: it assumes that the Lyman  $\alpha$  forest is just as likely to cause emission as absorption.

Here, we rectify this by not only including the Lyman  $\alpha$  perturbation term in our GP as  $\Omega$ , but introducing a redshift-dependent mean flux ( $\mu(z)$ ) with a dependence on the absorber redshift ( $z(\lambda_{\text{obs}})$ ). We model the overall mean model with a redshift-dependent absorption function and a mean emission vector:  $\mu(z) = a(z) \circ \mu$ . The notation  $\circ$  refers to Hadamard product, which is the element-wise product between two vectors or matrices. The covariance matrix is decomposed into  $\mathbf{A}_F(\mathbf{K} + \Omega)\mathbf{A}_F$ , where  $\text{diag}(\mathbf{A}_F) = a(z)$  and  $\mathbf{A}_F$  is a diagonal matrix.<sup>4</sup> The  $\mathbf{K}$  matrix describes the covariance between different emission lines in the QSO spectrum, which we will learn from data. The  $\mathbf{A}_F$  matrix is applied to  $\mathbf{K}$  because we assume that  $\mathbf{K}$  is learned before the absorption noise  $a(z)$  is applied. See Section 5.4 for how we learn the covariance.

Combining all modelling decisions, the model prior for an observed QSO emission is

$$p(\mathbf{y} | \lambda, \mathbf{v}, z_{\text{QSO}}, \mathcal{M}_{\text{-DLA}}) = \mathcal{N}(\mathbf{y}; \mu(z), \mathbf{A}_F(\mathbf{K} + \Omega)\mathbf{A}_F + \mathbf{V}). \quad (11)$$

The mean emission flux is now redshift- and wavelength-dependent, so the optimization steps will differ slightly from Garnett et al. (2017). We will address the modifications in the following subsections.

### 5.3 Redshift-dependent mean flux vector

In this paper, instead of using a single mean vector  $\mu$  to describe all spectra, we adjust the mean model of the GP to fit the mean flux of each QSO spectrum. For modelling the effect of forest absorption on the flux, we adopt an empirical power law with effective optical depth  $\tau_0(1+z)^\beta$  for Ly  $\alpha$  forest (Kim et al. 2007):

$$a(z) = \exp(-\tau_0(1+z)^\beta), \quad (12)$$

where the absorber redshift  $z$  is related to the observer's wavelength  $\lambda_{\text{obs}}$  as

$$\begin{aligned} 1+z &= \frac{\lambda_{\text{obs}}}{\lambda_{\text{Ly}\alpha}} \\ &= \frac{\lambda_{\text{obs}}}{1215.7 \text{ \AA}} \\ &= (1+z_{\text{QSO}}) \frac{\lambda}{1215.7 \text{ \AA}}, \end{aligned} \quad (13)$$

so the absorber redshift  $z(\lambda_{\text{obs}}) = z(\lambda, z_{\text{QSO}})$  is a function of the QSO redshift and the wavelength.

In Garnett et al. (2017), we assumed the absorption from the forest would only play a role in the additive noise term ( $\omega$ ) in our likelihood model  $p(\mathbf{y} | \lambda, \mathbf{v}, \omega, z_{\text{QSO}}, \mathcal{M}_{\text{-DLA}})$  with the form:

$$\omega'(\lambda, \lambda_{\text{obs}}) = \omega(\lambda) s(z(\lambda_{\text{obs}}))^2; \quad (14)$$

$$s(z) = 1 - \exp(-\tau_0(1+z)^\beta) + c_0, \quad (15)$$

where  $z$  is the absorber redshift. The  $\omega(\lambda)$  term represents the global absorption noise, and the  $s(z)$  corresponds to the absorption effect contributed by the Lyman  $\alpha$  absorbers along the line of sight as a function of the absorber redshift  $z$ .

Thus in our earlier model, the Lyman  $\alpha$  forest introduces additional fluctuations in the observed spectrum  $\mathbf{y}$ . This assumption worked

well for low-redshift spectra, because mean absorption due to the Lyman  $\alpha$  forest at low redshifts is relatively small. At high redshifts, however, the suppression of the mean flux induced by many Lyman  $\alpha$  absorbers is substantial, see Fig. 1. In our earlier model, essentially all high-redshift QSO spectra were substantially more absorbed than the mean emission model  $\mu$  due to absorption from the Lyman  $\alpha$  forest. To explain this absorption, our model would fit multiple DLAs with large column densities.

We have improved the modelling of the Lyman  $\alpha$  forest by allowing the mean GP model  $\mu$  to be redshift dependent, having a mean optical depth following the measurement of Kim et al. (2007):

$$\begin{aligned} \tau_{\text{eff}}(z) &= \tau_0(1+z)^\gamma \\ &= 0.0023 \times \exp(1+z)^{3.65}. \end{aligned} \quad (16)$$

There are other measurements of  $\tau_{\text{eff}}$  at higher precision than Kim et al. (2007; e.g. Becker et al. 2013). However, they are derived from SDSS data while Kim et al. (2007) were derived from high-resolution spectra. We therefore choose to use Kim et al. (2007) to preserve the likelihood principle that priors should not depend on the data set in question.

We include the effect of the whole Lyman series with a similar model, but however accounting for the different atomic coefficients of the higher order Lyman lines:

$$\begin{aligned} \tau_{\text{eff}, \text{HI}}(z(\lambda_{\text{obs}}); \gamma, \tau_0) \\ = \sum_{i=2}^N \tau_0 \frac{\lambda_{1i} f_{1i}}{\lambda_{12} f_{12}} (1+z_{1i}(\lambda_{\text{obs}}))^\gamma \times I_{(z_{1i}(\min(\lambda_{\text{obs}}), z_{\text{QSO}}))}(z). \end{aligned} \quad (17)$$

Here,  $f_{1i}$  represents the oscillator strength and  $\lambda_{1i}$  corresponds to the transition wavelength from the  $n=1$  to  $n=i$  atomic energy level. We model the Lyman series up to  $N=32$ , with  $i=2$  being Ly  $\alpha$  and  $i=3$  Ly  $\beta$ . The absorption redshift  $z_{1i}$  for the  $n=1$  to  $n=i$  transition is defined by

$$1+z_{1i} = \frac{\lambda_{\text{obs}}}{\lambda_{1i}}. \quad (18)$$

The optical depth at the line centre is estimated by

$$\tau_0 = \sqrt{\pi} \frac{e^2}{m_e c} \frac{N_\ell f_{\ell u} \lambda_{\ell u}}{b}, \quad (19)$$

where  $\ell$  indicates the lower energy level and  $u$  is the upper energy level. For Lyman  $\alpha$ , we have  $\lambda_{\ell u} = 1215.7 \text{ \AA}$  and  $f_{\ell u} = 0.4164$ ; for Lyman  $\beta$ , we have  $\lambda_{\ell u} = 1025.7$  and  $f_{\ell u} = 0.07912$ . Given equation (19), we have the effective optical depth for the Lyman  $\beta$  forest:

$$\tau_\beta = \frac{f_{31} \lambda_{31}}{f_{21} \lambda_{21}} \tau_0 = \frac{0.07912 \times 1025.7}{0.4164 \times 1215.7} \times 0.0023 = 0.0004. \quad (20)$$

The mean prior of the GP model for each spectrum is rewritten as

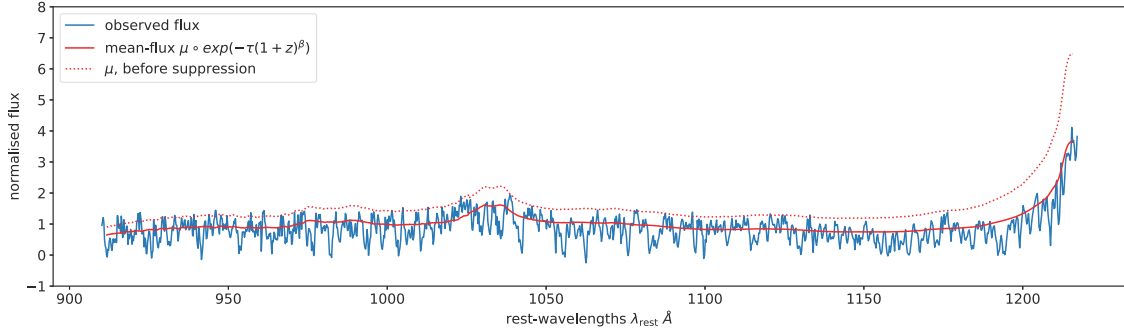
$$\mu(z) = \mu \circ \exp(-\tau_{\text{eff}, \text{HI}}(z; \gamma = 3.65, \tau_0 = 0.0023)). \quad (21)$$

We will simply write  $\tau_{\text{eff}, \text{HI}}(z) = \tau_{\text{eff}, \text{HI}}(z; \gamma = 3.65, \tau_0 = 0.0023)$  in the following text for simplicity. The new  $\mu$  is estimated via

$$\mu = \frac{1}{N_{\text{-NaN}}} \sum_{y_{ij} \neq \text{NaN}} y_{ij} \cdot \exp(+\tau_{\text{eff}, \text{HI}}(z_{ij})). \quad (22)$$

Equation (22) rescales the mean observed fluxes back to the expected continuum before the suppression due to Lyman series absorption, hopefully recovering approximately the true QSO emission function  $f$ . Fig. 1 shows the retrained mean QSO emission model for an example QSO. The mean model,  $\mu$ , is much closer to the peak emission flux above the absorbed forest.

<sup>4</sup> $\mathbf{A}_F^T = \mathbf{A}_F$  because it is diagonal.



**Figure 1.** The effect of the shift to the GP mean vector from the Lyman  $\alpha$  forest effective optical depth model ( $\mu \circ \exp(-\tau_0(1+z)^\beta)$ ). The dotted red curve shows the mean emission model before application of the forest suppression. The solid red curve is the mean model including the forest suppression.

For model consistency, we account for the mean suppression from weak absorbers in our redshift-dependent noise model  $\omega$  with:

$$\omega'(\lambda, \lambda_{\text{obs}}) = \omega(\lambda) s_F(z(\lambda_{\text{obs}}))^2; \quad (23)$$

$$\text{where } s_F(z(\lambda_{\text{obs}})) = 1 - \exp(-\tau_{\text{eff,H I}}(z(\lambda_{\text{obs}}); \beta, \tau_0)) + c_0. \quad (24)$$

$\tau_0$ ,  $\beta$ , and  $c_0$  are parameters that are learned from the data. Fig. 2 shows the mean model and absorption noise variance we use, compared to the model from Garnett et al. (2017).

Note that the mean flux model introduces degeneracies between the parameters of equation (24). For example,  $c_0$  may be compensated by the overall amplitude of pixel-wise noise vector  $\omega$ . For this reason, we should not ascribe strict physical interpretations to the optimal values of equation (24). The optimized  $\omega'$  is simply an empirical relation modelling the pixel-wise and redshift-dependent noise in the null model given SDSS data.

After introducing the effective optical depth into our GP mean model, we decrease the number of large DLAs we detect at high redshifts and thus measure lower  $\Omega_{\text{DLA}}$  at high redshifts (see Section 10.3 for more details). This is because, for high-redshift QSOs, the mean optical depth may be close to unity. To explain this unexpected absorption, the previous code will fit multiple high column density absorbers to the raw emission model, artificially increasing the number of DLAs detected. With the mean suppressed, there is substantially less raw absorption to explain, and so this tendency is avoided.

#### 5.4 Learning the flux covariance

$\mathbf{K}$  and  $\Omega$  (equation 11) are optimized to maximize the likelihood of generating the data,  $\mathcal{D}$ . The mean flux model is not optimized, but follows the effective optical depth reported in Kim et al. (2007). Thus, we remove the effect of forest absorption before we train the covariance function and train on  $\mathcal{D}' = \{\lambda, \mathbf{y} \circ \exp(+\tau_{\text{eff,H I}}(z)) - \mu(z)\}$  to find the optimal parameters for  $\mathbf{K}$  and  $\Omega$ .

We assume the same likelihood as Garnett et al. (2017) for generating the whole training data set ( $\mathbf{Y}$ ):

$$p(\mathbf{Y} | \lambda, \mathbf{V}, \mathbf{M}, \omega, z_{\text{QSO}}, \mathcal{M}_{\text{-DLA}}) = \prod_{i=1}^{N_{\text{spec}}} \mathcal{N}(\mathbf{y}_i; \mu, \mathbf{K} + \Omega + \mathbf{V}_i), \quad (25)$$

where  $\mathbf{Y}$  means the matrix containing all the observed flux in the training data, and the product on the right-hand side says we are combining all likelihoods from each single spectrum. The noise

matrix  $\Omega = \text{diag } \omega'$  is the diagonal matrix that represents the Lyman  $\alpha$  forest absorption from equation (24).

$\mathbf{M}$  is a low-rank decomposition of the covariance matrix  $\mathbf{K}$  we want to learn:

$$\mathbf{K} = \mathbf{M}\mathbf{M}^T, \quad (26)$$

where  $\mathbf{M}$  is an  $(N_{\text{pixels}} \times k)$  matrix. Without this low-rank decomposition, we would need to learn  $N_{\text{pixel}}^2 = 1217 \times 1217$  free parameters. With equation (26), we can limit the number of free parameters to be  $N_{\text{pixels}} \times k$ , where  $k \ll N_{\text{pixels}}$ ; also, it guarantees the covariance matrix  $\mathbf{K}$  to be positive semidefinite. Each column of the  $\mathbf{M}$  can be treated as an eigenspectrum of the training data, where we set the number of eigenspectra to be  $k = 20$ . We will optimize the  $\mathbf{M}$  matrix and the absorption noise in equation (24) simultaneously.

A modification performed in this work is to, instead of directly training on the observed flux, optimize the covariance matrix and noise model on the flux with Lyman  $\alpha$  forest absorption removed (de-forest flux):

$$\begin{aligned} \mathbf{y} &:= \mathbf{y} \circ \exp(+\tau_{\text{eff,H I}}(z)); \\ Y_{ij} &:= Y_{ij} \exp(+\tau_{\text{eff,H I}}(z))_{ij}. \end{aligned} \quad (27)$$

We may write this change into the likelihood:

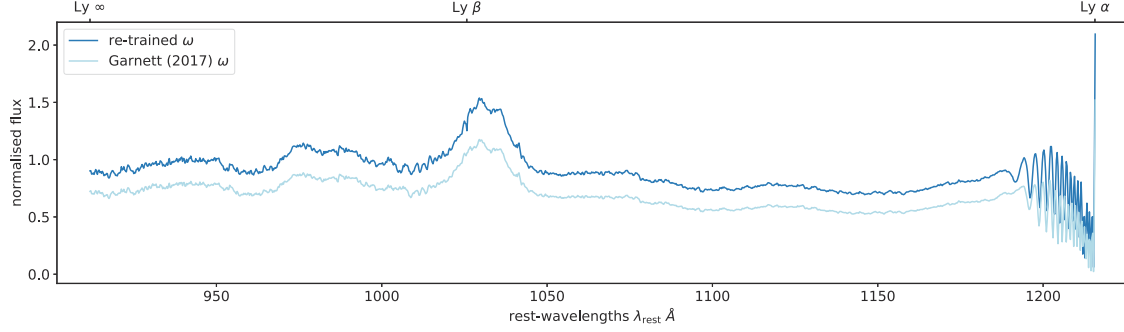
$$p(\mathbf{Y} \circ \exp(+\tau_{\text{eff,H I}}(z)) | \lambda, \mathbf{V}, \mathbf{M}, \omega, z_{\text{QSO}}, \mathcal{M}_{\text{-DLA}}) = \prod_{i=1}^{N_{\text{spec}}} \mathcal{N}(\mathbf{y}_i \circ \exp(+\tau_{\text{eff,H I}}(z_i)); \mu, \mathbf{K} + \Omega + \mathbf{V}_i), \quad (28)$$

where  $\mu$  is the mean model from equation (22). The rest of our optimization procedure follows the unconstrained optimization of Garnett et al. (2017).

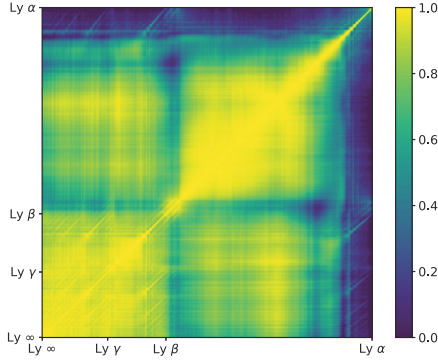
We use de-forest fluxes for training as we want our covariance matrix to learn the covariance in the true emission function. The emission function (like our kernel  $\mathbf{K}$ ) is independent of QSO emission redshift, whereas the absorption noise is not. We only implement the mean forest absorption of Kim et al. (2007), so we need an extra term to compensate for the variance of the forest around this mean. We thus still train the redshift- and wavelength-dependent absorption noise from data. The optimal values we learned for equation (24) are:

$$c_0 = 0.3050; \tau_0 = 1.6400 \times 10^{-4}; \beta = 5.2714. \quad (29)$$

As we might expect, the optimal  $\tau_0$  value is smaller than the  $\tau_0 = 0.01178$  learned in Garnett et al. (2017), which implies the effect of the forest is almost removed by applying the Lyman-series forest to the mean model. The final covariance matrix is shown in Fig. 3.



**Figure 2.** The difference between the original pixel-wise noise variance  $\omega$  (Garnett et al. 2017) and the retrained  $\omega$  from equation (28). The retrained  $\omega$  decreases because the fit no longer needs to account for the mean forest absorption.



**Figure 3.** The trained covariance matrix  $\mathbf{M}$ , which is almost the same as the covariance from Garnett et al. (2017). Note that we normalize the diagonal elements to be unity, so this is more like a correlation matrix than a covariance matrix. The values in the matrix are ranging from 0 to 1, representing the correlation between  $\lambda$  and  $\lambda'$  in the QSO emission.

### 5.5 Model evidence

Consider a given QSO observation  $\mathcal{D} = (\lambda, y)$  with known observational noise  $v(\lambda)$  and known QSO redshift  $z_{\text{QSO}}$ . The model evidence for  $\mathcal{M}_{\text{-DLA}}$  can be estimated using

$$p(\mathcal{D} | \mathcal{M}_{\text{-DLA}}, v, z_{\text{QSO}}) \propto p(y | \lambda, v, z_{\text{QSO}}, \mathcal{M}_{\text{-DLA}}), \quad (30)$$

which is equivalent to evaluating a multivariate Gaussian

$$p(y | \lambda, v, z_{\text{QSO}}, \mathcal{M}_{\text{-DLA}}) = \mathcal{N}(y; \mu \circ \exp(-\tau_{\text{eff,H I}}), \mathbf{A}_F(\mathbf{K} + \mathbf{\Omega})\mathbf{A}_F + \mathbf{V}). \quad (31)$$

Here,  $\exp(-\tau_{\text{eff,H I}}) = \text{diag } \mathbf{A}_F$  describes the absorption due to the forest and modifies the mean vector  $\mu$ , the covariance matrix  $\mathbf{K}$ , and the noise matrix  $\mathbf{\Omega}$  to account for the Lyman  $\alpha$  forest effective optical depth.

## 6 A GP MODEL FOR QSO SIGHTLINES WITH MULTIPLE DLAs

In Section 5, we learned a GP prior for QSO spectroscopic measurements without any DLAs for our null model  $\mathcal{M}_{\text{-DLA}}$ . Here, we extend the null model  $\mathcal{M}_{\text{-DLA}}$  to a model with  $k$  intervening DLAs,  $\mathcal{M}_{\text{DLA}(k)}$ .

Our complete DLA model,  $\mathcal{M}_{\text{DLA}}$ , will be the union of the models with  $i$  DLAs:  $\mathcal{M}_{\text{DLA}} = \{\mathcal{M}_{\text{DLA}(i)}\}_{i=1}^k$ . We consider only until  $k = 4$ , as DLAs are rare events and our sample only contains one spectrum with four DLAs.

### 6.1 Absorption function

Before we model a QSO spectrum with intervening DLAs, we need to have an absorption profile model for a DLA. Damped Lyman alpha absorbers, or DLAs, are neutral hydrogen (H I) absorption systems with saturated lines and damping wings in the spectroscopic measurements. Having saturated lines means the column density of the absorbers on the line of sight is high enough to absorb essentially all photons. The damping wings are due to natural broadening in the line.

The optical depth from each Lyman series transition is

$$\tau(\lambda; z_{\text{DLA}}, N_{\text{H I}}) = N_{\text{H I}} \frac{\pi e^2 f_{1u} \lambda_{1u}}{m_e c} \phi(v, b, \gamma), \quad (32)$$

where  $e$  is the elementary charge,  $\lambda_{1u}$  is the transition wavelength from the  $n = 1$  to  $n = u$  energy level ( $\lambda_{12} = 1215.6701 \text{ \AA}$  for Lyman  $\alpha$ ), and  $f_{1u}$  is the oscillator strength of the transition. The line profile  $\phi$  is a Voigt profile:

$$\phi(v, b, \gamma) = \int \frac{dv}{\sqrt{2\pi}\sigma_v} \exp(-v^2/2\sigma_v^2) \times \frac{4\gamma_{\ell u}}{16\pi^2[v - (1 - v/c)v_{\ell u}]^2 + \gamma_{\ell u}^2}, \quad (33)$$

which is a convolution between a Lorentzian line profile and a Gaussian line profile. The  $\sigma_v$  is the one-dimensional velocity dispersion,  $\gamma_{\ell u}$  is a parameter for Lorentzian profile,  $v$  is the frequency, and  $u$  represents the upper energy level and  $\ell$  represents the lower energy level.

Both profiles are parametrized by the relative velocity  $v$ , which means both profiles are distributions in the one-dimensional velocity space:

$$v = c \left( \frac{\lambda}{\lambda_{1u}} \frac{1}{(1 + z_{\text{DLA}})} - 1 \right). \quad (34)$$

The standard deviation of the Gaussian line profile is related to the broadening parameter  $b = \sqrt{2}\sigma_v$ , and if we assume the broadening is entirely due to thermal motion:

$$b = \sqrt{\frac{2kT}{m_p}}. \quad (35)$$

Introducing the damping constant  $\Gamma = 6.265 \times 10^8 \text{ s}^{-1}$  for Lyman  $\alpha$ , we have the parameter  $\gamma_{\ell u}$  to describe the width of the Lorentzian profile

$$\gamma_{\ell u} = \frac{\Gamma \lambda_{\ell u}}{4\pi}. \quad (36)$$

Our default DLA profile includes Ly  $\alpha$ , Ly  $\beta$ , and Ly  $\gamma$  absorptions. We fix the broadening parameter  $b$  by setting  $T = 10^4$  K, which increases the width of the DLA profile by  $13 \text{ km s}^{-1}$ , small compared to the effect of the Lorentzian wings. Thus, for a given QSO and a true emission function  $f(\lambda)$ , the function for the observed flux  $y(\lambda)$  is

$$y(\lambda) = f(\lambda) \exp(-\tau(\lambda; z_{\text{DLA}}, N_{\text{H I}})) \exp(-\tau_{\text{eff, H I}}(\lambda_{\text{obs}})) + \epsilon, \quad (37)$$

where  $\epsilon$  is additive Gaussian noise including measurement noise and absorption noise.

Suppose we have a DLA at redshift  $z_{\text{DLA}}$  with column density  $N_{\text{H I}}$ . We can model the spectrum with an intervening DLA by calculating the DLA absorption function:

$$\mathbf{a} = \exp(-\tau(\lambda; z_{\text{DLA}}, N_{\text{H I}})). \quad (38)$$

We apply the absorption function to the GP prior of  $\mathbf{y}$  with

$$\begin{aligned} p(\mathbf{y} | \boldsymbol{\lambda}, \mathbf{v}, z_{\text{QSO}}, z_{\text{DLA}}, N_{\text{H I}}, \mathcal{M}_{\text{DLA}}) \\ = \mathcal{N}(\mathbf{y}; \mathbf{a} \circ (\mathbf{a}_{\text{F}} \circ \boldsymbol{\mu}), \mathbf{A}(\mathbf{A}_{\text{F}}(\mathbf{K} + \boldsymbol{\Omega})\mathbf{A}_{\text{F}})\mathbf{A} + \mathbf{V}), \end{aligned} \quad (39)$$

where  $\mathbf{A} = \text{diag } \mathbf{a}$ .

For a model with  $k$  DLAs with  $k \in \mathbb{N}$ , we simply take the element-wise product of  $k$  absorption functions:

$$\begin{aligned} \mathbf{a}_{(k)} &= \prod_{i=1}^k \mathbf{a}(\lambda; z_{\text{DLA}_i}, N_{\text{H I}_i}); \\ \text{diag } \mathbf{A}_{(k)} &= \mathbf{a}_{(k)}. \end{aligned} \quad (40)$$

The prior for  $\mathcal{M}_{\text{DLA}(k)}$  would therefore be

$$\begin{aligned} p(\mathbf{y} | \boldsymbol{\lambda}, \mathbf{v}, z_{\text{QSO}}, \{z_{\text{DLA}_i}\}_{i=1}^k, \{N_{\text{H I}_i}\}_{i=1}^k, \mathcal{M}_{\text{DLA}(k)}) \\ = \mathcal{N}(\mathbf{y}; \mathbf{a}_{(k)} \circ (\mathbf{a}_{\text{F}} \circ \boldsymbol{\mu}), \mathbf{A}_{(k)}(\mathbf{A}_{\text{F}}(\mathbf{K} + \boldsymbol{\Omega})\mathbf{A}_{\text{F}})\mathbf{A}_{(k)} + \mathbf{V}). \end{aligned} \quad (41)$$

Here, we briefly review our notations in equation (41):  $\mathbf{a}_{(k)}$ , which is parametrized by  $(\{z_{\text{DLA}_i}\}_{i=1}^k, \{N_{\text{H I}_i}\}_{i=1}^k)$ , represents the absorption function with  $k$  DLAs in one spectrum. Note that each DLA is parametrized by a pair of  $(z_{\text{DLA}}, N_{\text{H I}})$ .  $\mathbf{a}_{\text{F}}$  corresponds to the absorption function from the Lyman series absorptions, which is derived from Kim et al. (2007) in the form of equation (21). The covariance matrix  $\mathbf{K}$  and the absorption model  $\boldsymbol{\Omega}$  are both learned from data, as described in Section 5.4.  $\mathbf{V}$  is the noise variance matrix given by the SDSS pipeline, so each sightline would have different  $\mathbf{V}$ .

## 6.2 Model evidence: DLA(1)

The model evidence of our DLA model is given by the integral:

$$\begin{aligned} p(\mathcal{D} | \mathcal{M}_{\text{DLA}(1)}, z_{\text{QSO}}) &\propto \int p(\mathbf{y} | \boldsymbol{\lambda}, \mathbf{v}, \theta, z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}) \\ &\times p(\theta | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}) d\theta, \end{aligned} \quad (42)$$

where we integrated out the parameters,  $\theta = (z_{\text{DLA}}, \log_{10} N_{\text{H I}})$ , with a given parameter prior  $p(\theta | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)})$ .

However, equation (42) is intractable, so we approximate it with a quasi-Monte Carlo method (QMC). QMC selects  $N = 10\,000$  samples with an approximately uniform spatial distribution from a Halton sequence to calculate the model likelihood, approximating the model evidence by the sample mean:

$$p(\mathcal{D} | \mathcal{M}_{\text{DLA}(1)}, z_{\text{QSO}}) \simeq \frac{1}{N} \sum_{i=1}^N p(\mathcal{D} | \theta_i, z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}). \quad (43)$$

## 6.3 Model evidence: Occam's razor effect for DLA(k)

For higher order DLA models, we have to integrate out not only the nuisance parameters of the first DLA model  $\mathcal{M}_{\text{DLA}(1)}$ ,  $(\theta_1)$  but also the parameters from  $\mathcal{M}_{\text{DLA}(2)}$  to  $\mathcal{M}_{\text{DLA}(k)}$ ,

$$\begin{aligned} p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, z_{\text{QSO}}) &\propto \int p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, \{\theta_i\}_{i=1}^k) \\ &\times p(\{\theta_i\}_{i=1}^k | \mathcal{M}_{\text{DLA}(k)}, \mathcal{D}, z_{\text{QSO}}) d\{\theta_i\}_{i=1}^k, \end{aligned} \quad (44)$$

which means we are marginalizing  $\{\theta_i\}_{i=1}^k$  in a parameter space with  $2 \times k$  dimensions. The parameter prior of multi-DLAs is a multiplication between a non-informative prior  $p(\theta_i | \mathcal{M}_{\text{DLA}(1)}, z_{\text{QSO}})$  and the posterior of the  $(k-1)$  multi-DLA model,

$$\begin{aligned} p(\{\theta_i\}_{i=1}^k | \mathcal{M}_{\text{DLA}(k)}, \mathcal{D}, z_{\text{QSO}}) \\ = p(\{\theta_i\}_{i=1}^{k-1} | \mathcal{M}_{\text{DLA}(k-1)}, \mathcal{D}, z_{\text{QSO}}) p(\theta_k | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}). \end{aligned} \quad (45)$$

We can approximate this integral using the same QMC method. For example, if we want to sample the model evidence for  $\mathcal{M}_{\text{DLA}(2)}$ , we would need  $N = 10\,000$  samples for each parameter dimension  $\{\theta_i\}_{i=1}^2$ , which results in sampling from two independent Halton sequences with  $10^8$  samples in total. If we want to sample up to  $\mathcal{M}_{\text{DLA}(k)}$  with  $N$  samples for each  $\{\theta_i\}$  from  $i = 1, \dots, k$ , we would need to have:

$$\begin{aligned} p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, z_{\text{QSO}}) &\simeq \frac{1}{N} \sum_{j^{(1)}=1}^N \frac{1}{N} \sum_{j^{(2)}=1}^N \frac{1}{N} \sum_{j^{(3)}=1}^N \dots \frac{1}{N} \sum_{j^{(k)}=1}^N \\ &p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, \{\theta_{1j^{(1)}}\}, \{\theta_{2j^{(2)}}\}, \{\theta_{3j^{(3)}}\}, \dots, \{\theta_{kj^{(k)}}\}, z_{\text{QSO}}), \end{aligned} \quad (46)$$

where  $\{j^{(1)}, j^{(2)}, j^{(3)}, \dots, j^{(k)}\}$  indicate the indices of QMC samples. The above equation (46) is thus in principle evaluated with  $N^k$  samples.

In practice, we only sample  $N = 10\,000$  points from  $p(\{\theta_i\}_{i=1}^k | \mathcal{M}_{\text{DLA}(k)}, \mathcal{D}, z_{\text{QSO}})$  instead of sampling  $N^k$  points, as a uniform sampling of the first DLA model may be reweighted to cover parameter space for the higher order models. A  $N^{k-1}$  factor of normalization is thus left behind in the summation,

$$\begin{aligned} p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, z_{\text{QSO}}) &\simeq \frac{1}{N^k} \sum_{j=1}^N p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, \{\theta_{ij}\}_{i=1}^k, z_{\text{QSO}}) \\ &\simeq \frac{1}{N^{k-1}} \left( \frac{1}{N} \sum_{j=1}^N p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, \{\theta_{ij}\}_{i=1}^k, z_{\text{QSO}}) \right) \\ &\simeq \frac{1}{N^{k-1}} \text{mean}_j (p(\mathcal{D} | \mathcal{M}_{\text{DLA}(k)}, \{\theta_{ij}\}_{i=1}^k, z_{\text{QSO}})). \end{aligned} \quad (47)$$

The additional  $\frac{1}{N^{k-1}}$  factor penalizes models with more parameters than needed, and can be viewed as an implementation of Occam's razor. This Occam's razor effect is caused by the fact that all probability distributions have to be normalized to unity. A model with more parameters, which means having a wider distribution in the likelihood space, results in a bigger normalization factor.

The motivation for us to draw  $N$  samples from the multi-DLA likelihood function  $p(\{\theta_i\}_{i=1}^k | \mathcal{M}_{\text{DLA}(k)}, \mathcal{D}, z_{\text{QSO}})$  is that we believe the prior density we took from the posterior density of  $\mathcal{M}_{\text{DLA}(k-1)}$  is representative enough even without  $N^k$  samples. For example, if we have two peaks in our likelihood density  $p(\mathcal{D} | \mathcal{M}_{\text{DLA}(1)}, \theta_1, z_{\text{QSO}})$ , we expect the sampling for  $\theta_2$  in  $p(\mathcal{D} | \mathcal{M}_{\text{DLA}(2)}, \{\theta_1, \theta_2\}, z_{\text{QSO}})$

would concentrate on sampling the density of the first highest peak in  $p(\mathcal{D} | \mathcal{M}_{\text{DLA}(1)}, \theta_1, z_{\text{QSO}})$  density. Similarly, while we are sampling for  $\mathcal{M}_{\text{DLA}(3)}$ , we expect  $\theta_3$  and  $\theta_2$  would cover the first- and the second-highest peaks.

To avoid multi-DLAs overlapping with each other, we inject a dependence between any pair of  $z_{\text{DLA}}$  parameters. Specifically, if any pair of  $z_{\text{DLA}}$ s have a relative velocity smaller than  $3000 \text{ km s}^{-1}$ , then we set the likelihood of this sample to NaN.

#### 6.4 Additional penalty for DLAs and sub-DLAs

In Section 6.3, we apply a penalty, Occam's razor, to regularize DLA models using more parameters than needed. This effect is due to the normalization (to unity) of the evidence.

In a similar fashion, and for a similar reason, we apply an additional regularization factor between the non-DLA and DLA models (including sub-DLAs). This additional factor ensures that when both models are a poor fit to a particular observational spectrum, the code prefers the non-DLA model, rather than preferring the model with more parameters and thus greater fitting freedom. We directly inject this Occam's razor factor in the model selection:

$$\Pr(\mathcal{M}_{\text{DLA}} | \mathcal{D}) = \frac{\Pr(\mathcal{M}_{\text{DLA}}) p(\mathcal{D} | \mathcal{M}_{\text{DLA}})^{\frac{1}{N}}}{\left( \frac{\Pr(\mathcal{M}_{\text{DLA}}) p(\mathcal{D} | \mathcal{M}_{\text{DLA}})}{\Pr(\mathcal{M}_{\text{sub}}) p(\mathcal{D} | \mathcal{M}_{\text{sub}})} \right)^{\frac{1}{N}} + \Pr(\mathcal{M}_{\text{-DLA}} | \mathcal{D})}, \quad (48)$$

where  $N = 10^4$  is the number of samples we used to approximate the parametrized likelihood functions. We evaluated the impact of this regularization factor on the area under the curve (AUC) in the receiver-operating characteristics (ROC) plot.<sup>5</sup> For  $N = 10^4$ , the AUC changed from 0.949 to 0.960. We considered other penalty values and found that the AUC increased up to  $N = 10^4$  and then plateaued.

In addition, we found by examining specific examples that this penalty regularized a relatively common incorrect DLA detection: finding objects in short, very noisy low redshift ( $z \sim 2.2$ ) spectra. In these spectra our earlier model would prefer the DLA model purely because of its large parameter freedom. In particular a high column density DLA, large enough that the damping wings exceed the width of the spectrum, would be preferred. Such a fit exploits a degeneracy in the model between the mean observed flux and the DLA column density when the spectrum is shorter than the putative DLA. The Occam's razor penalty avoids these spurious fits by penalizing the extra parametric freedom in the DLA model.

#### 6.5 Parameter prior

Here, we briefly recap the priors on model parameters chosen in Garnett et al. (2017). Suppose we want to make an inference for the column density and redshift of an absorber  $\theta = (N_{\text{H I}}, z_{\text{DLA}})$  from a given spectroscopic observation, the joint density for the parameter prior would be

$$p(\theta | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}) = p(N_{\text{H I}}, z_{\text{DLA}} | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}). \quad (49)$$

Suppose the absorber redshift and the column density are conditionally independent and the column density is independent of the QSO redshift  $z_{\text{QSO}}$ :

$$p(\theta | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}) = p(z_{\text{DLA}} | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}) \times p(N_{\text{H I}} | \mathcal{M}_{\text{DLA}(1)}) \quad (50)$$

We set a bounded uniform prior density for the absorber redshift  $z_{\text{DLA}}$ :

$$p(z_{\text{DLA}} | z_{\text{QSO}}, \mathcal{M}_{\text{DLA}(1)}) = \mathcal{U}[z_{\text{min}}, z_{\text{max}}], \quad (51)$$

where we define the finite prior range to be

$$z_{\text{min}} = \max \left\{ \frac{\lambda_{\text{Ly}\infty}}{\lambda_{\text{Ly}\alpha}} (1 + z_{\text{QSO}}) - 1 + 3000 \text{ km s}^{-1}/c, \frac{\min \lambda_{\text{obs}}}{\lambda_{\text{Ly}\alpha}} - 1 \right\} \quad (52)$$

$$z_{\text{max}} = z_{\text{QSO}} - 3000 \text{ km s}^{-1}/c; \quad (53)$$

which means we have a prior belief that the centre of the absorber is within the observed wavelengths. The range of observed wavelengths is either from  $\text{Ly } \infty$  to  $\text{Ly } \alpha$  of the QSO rest frame ( $\lambda_{\text{rest}} \in [911.75 \text{ \AA}, 1216.75 \text{ \AA}]$ ) or from the minimum observed wavelength to  $\text{Ly } \alpha$ . We also apply a conservative cut-off of  $3000 \text{ km s}^{-1}$  near to  $\text{Ly } \infty$  and  $\text{Ly } \alpha$ . The  $-3000 \text{ km s}^{-1}$  cut-off for  $z_{\text{max}}$  helps to avoid proximity ionization effects due to the QSO radiation field. Furthermore, the  $+3000 \text{ km s}^{-1}$  cutoff for  $z_{\text{min}}$  avoids a potentially incorrect measurement for  $z_{\text{QSO}}$ . An underestimated  $z_{\text{QSO}}$  can produce a Lyman-limit trough within the region of the QSO expected to contain only Lyman-series absorption, and the code can incorrectly interpret this as a DLA.

For the column density prior, we follow Garnett et al. (2017). We first estimate the density of DLAs column density  $p(N_{\text{H I}} | \mathcal{M}_{\text{DLA}})$  using the BOSS DR9 Lyman  $\alpha$  forest sample. We choose to put our prior on the base-10 logarithm of the column density  $\log_{10} N_{\text{H I}}$  due to the large dynamic range of DLA column densities in SDSS DR9 samples.

We thus estimate the density of logarithm column densities  $p(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{DLA}(1)})$  using univariate Gaussian kernels on the reported  $\log_{10} N_{\text{H I}}$  values in DR9 samples. Column densities from DLAs in DR9 with  $N_{\text{DLA}} = 5854$  are used to non-parametrically estimate the logarithm  $N_{\text{H I}}$  prior density, with

$$p_{\text{KDE}}(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{DLA}(1)}) = \frac{1}{N_{\text{DLA}}} \sum_{i=1}^{N_{\text{DLA}}} \mathcal{N}(\log_{10} N_{\text{H I}}; l_i, \sigma^2), \quad (54)$$

where  $l_i$  is the logarithm column density  $\log_{10} N_{\text{H I}}$  of the  $i$ th sample. The bandwidth  $\sigma^2$  is selected to be the optimal value for a normal distribution, which is the default setting for MATLAB.

We further simplify the non-parametric estimate into a parametric form with

$$p_{\text{KDE}}(\log_{10} N_{\text{H I}} = N | \mathcal{M}_{\text{DLA}(1)}) \simeq q(\log_{10} N_{\text{H I}} = N) \propto \exp(aN^2 + bN + c); \quad (55)$$

where the parameters  $(a, b, c)$  for the quadratic function are fitted via standard least-squared fitting to the non-parametric estimate of density  $p_{\text{KDE}}(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{DLA}(1)})$  with the range  $\log_{10} N_{\text{H I}} \in [20, 22]$ . The optimal values for the quadratic terms were

$$a = -1.2695; b = 50.863; c = -509.33. \quad (56)$$

Note that we have the same values as in Garnett et al. (2017).

Finally, we choose to be conservative about the data-driven column density prior. We thus take a mixture of a non-informative lognormal prior with the data-driven prior to make a non-restrictive prior on a large dynamical range:

$$p(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{DLA}(1)}) = \alpha q(\log_{10} N_{\text{H I}} = N) + (1 - \alpha) \mathcal{U}[20, 23]. \quad (57)$$

Here we choose the mixture coefficient  $\alpha = 0.97$ , which favours the data-driven prior. We still include a small component of a non-

<sup>5</sup>See Section 10.1 for how we compute our ROC plot.

informative prior so that we are able to detect DLAs with a larger column density than in the training set, if any are present in the larger DR12 sample. Note that  $\alpha = 0.97$  is 7 per cent higher than the coefficient chosen in Garnett et al. (2017), which was  $\alpha = 0.90$ . Our previous prior slightly overestimated the number of very large DLAs.

## 6.6 Sub-DLA parameter prior

As reported in Bird et al. (2017), the column density distribution function (CDDF) exhibited an edge feature: an overdetection of DLAs at low column densities ( $\sim 10^{20} \text{ cm}^{-2}$ ). This did not affect the statistical properties of DLAs as we restrict column density to  $N_{\text{H I}} \geq 10^{20.3} \text{ cm}^{-2}$  for both line densities ( $dN/dX$ ) and total column densities ( $\Omega_{\text{DLA}}$ ). However, to make our method more robust, here we describe a complementary method to avoid overestimating the number of low column density absorbers.

The excess of DLAs at  $\sim 10^{20} \text{ cm}^{-2}$  is due to our model excluding lower column density absorbers such as sub-DLAs. Since we limited our column density prior of DLAs to be larger than  $10^{20} \text{ cm}^{-2}$ , the code cannot correctly classify a sub-DLA. Instead it correctly notes that a sub-DLA spectrum is more likely to be a DLA with a minimal column density than an unabsorbed spectrum.

To resolve our ignorance, we introduce an alternative model  $\mathcal{M}_{\text{sub}}$  to account the model posterior of those low column density absorbers in our Bayesian model selection. The likelihood function we used for sub-DLAs is identical to the one we built for DLA model  $\mathcal{M}_{\text{DLA}(1)}$  in equation (39) but has a different parameter prior on the column densities  $p(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{sub}})$ . We restricted our prior belief of sub-DLAs to be within the range  $\log_{10} N_{\text{H I}} \in [19.5, 20]$ , and, as we do not have a catalogue of sub-DLAs for learning the prior density, we put a uniform prior on  $\log_{10} N_{\text{H I}}$ :

$$p(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{sub}}) = \mathcal{U}[19.5, 20]. \quad (58)$$

We place a lower cut-off at  $\log_{10} N_{\text{H I}} = 19.5$  because the relatively noisy SDSS data offer limited evidence for absorbers with column densities lower than this limit.

## 7 MODEL PRIORS

Bayesian model selection allows us to combine prior information with evidence from the data-driven model to obtain a posterior belief about the detection of DLAs  $p(\mathcal{M}_{\text{DLA}} | \mathcal{D})$  using Bayes' rule. For a given spectroscopic observation  $\mathcal{D}$ , we already have the ability to compute the model evidence for a DLA ( $p(\mathcal{D} | \mathcal{M}_{\text{DLA}})$ ) and no DLA ( $p(\mathcal{D} | \mathcal{M}_{\text{-DLA}})$ ). However, to compute the model posteriors, we need to specify our prior beliefs in these models. Here, we approximate our prior belief  $\Pr(\mathcal{M}_{\text{DLA}})$  using the SDSS DR9 DLA catalogue.

Consider a QSO observation  $\mathcal{D} = (\lambda, y)$  at  $z_{\text{QSO}}$ . We want to find our prior belief that  $\mathcal{D}$  contains a DLA. We count the fraction of QSO sightlines in the training set containing DLAs with redshift less than  $z_{\text{QSO}} + z'$ , where  $z' = 30\,000 \text{ km s}^{-1}/c$  is a small constant. If  $N$  is the number of QSO sightlines with redshift less than  $z_{\text{QSO}} + z'$ , and  $M$  is the number of sightlines in this set containing DLAs in the QSO rest-frame wavelengths range we search, then our empirical prior for  $\mathcal{M}_{\text{DLA}}$  is

$$\Pr(\mathcal{M}_{\text{DLA}} | z_{\text{QSO}}) = \frac{M}{N}. \quad (59)$$

We can break down our DLA prior  $\Pr(\mathcal{M}_{\text{DLA}} | z_{\text{QSO}})$  for multiple DLAs in a QSO sightline  $\Pr(\mathcal{M}_{\text{DLA}(k)} | z_{\text{QSO}})$  via

$$\Pr(\mathcal{M}_{\text{DLA}(k)} | z_{\text{QSO}}) \simeq \left(\frac{M}{N}\right)^k - \left(\frac{M}{N}\right)^{k+1}. \quad (60)$$

For example,  $\frac{M}{N}$  represents our prior belief of having at least one DLA in the sightline, and  $(\frac{M}{N})^2$  represents having at least two DLAs.  $\frac{M}{N} - (\frac{M}{N})^2$  is thus our prior belief of having exactly one DLA at the sightline.

## 7.1 Sub-DLA model prior

The CDDF of Bird et al. (2017) exhibited an edge effect at  $\log_{10} N_{\text{H I}} \sim 20$  due to a lack of sampling at lower column densities. We thus construct an alternative model for lower column density absorbers (sub-DLAs, DLAs' lower column density cousins) to regularize DLA detections. We use the same GP likelihood function as the DLA model  $\mathcal{M}_{\text{DLA}}$  to compute our sub-DLA model evidence  $p(\mathcal{D} | \mathcal{M}_{\text{sub}})$  but with a different column density prior  $p(\log_{10} N_{\text{H I}} | \mathcal{M}_{\text{sub}})$ .

There is no sub-DLA catalogue available for us to estimate the empirical prior directly. We, therefore, approximate our sub-DLA model prior by rescaling our DLA model prior:

$$\Pr(\mathcal{M}_{\text{sub}} | z_{\text{QSO}}) \propto \Pr(\mathcal{M}_{\text{DLA}} | z_{\text{QSO}}), \quad (61)$$

and we require our prior beliefs to sum to unity:

$$\Pr(\mathcal{M}_{\text{-DLA}} | z_{\text{QSO}}) + \Pr(\mathcal{M}_{\text{sub}} | z_{\text{QSO}}) + \Pr(\mathcal{M}_{\text{DLA}} | z_{\text{QSO}}) = 1. \quad (62)$$

The scaling factor between the DLA prior and sub-DLA prior should depend on our prior probability density of the column density of the absorbers. Here, we assume the density of sub-DLA  $\log_{10} N_{\text{H I}}$  is a uniform density with a finite range of  $\log_{10} N_{\text{H I}} \in [19.5, 20]$ . We believe there are more sub-DLAs than DLAs as high column density systems are generally rarer. We thus assume the probability of finding sub-DLAs at a given  $\log_{10} N_{\text{H I}}$  is the same as the probability of finding DLAs at the most probable  $\log_{10} N_{\text{H I}}$ , which is

$$\begin{aligned} p(\log_{10} N_{\text{H I}} = N | \{\mathcal{M}_{\text{DLA}}, \mathcal{M}_{\text{sub}}\}) \\ = \alpha q(N | \mathcal{M}_{\text{DLA}}) \mathbb{I}_{(20, 23)}(N) \\ + \alpha \max(q(N | \mathcal{M}_{\text{DLA}}) \mathbb{I}_{(19.5, 20)}(N) \\ + (1 - \alpha) \mathcal{U}[19.5, 23]. \end{aligned} \quad (63)$$

Since  $q(N | \mathcal{M}_{\text{DLA}})$  has a simple quadratic functional form, we can solve the maximum value analytically, which is  $\max(q(N | \mathcal{M}_{\text{DLA}})) \simeq q(N = 20.03 | \mathcal{M}_{\text{DLA}})$ .

We thus can use our prior knowledge about the logarithm of column densities for different absorbers to rescale model priors:

$$\Pr(\mathcal{M}_{\text{sub}} | z_{\text{QSO}}) = \frac{Z_{\text{sub}}}{Z_{\text{DLA}}} \Pr(\mathcal{M}_{\text{DLA}} | z_{\text{QSO}}), \quad (64)$$

where the scaling factor is

$$\frac{Z_{\text{sub}}}{Z_{\text{DLA}}} = \frac{\int_{19.5}^{20} p(N | \{\mathcal{M}_{\text{DLA}}, \mathcal{M}_{\text{sub}}\}) dN}{\int_{20}^{23} p(N | \{\mathcal{M}_{\text{DLA}}, \mathcal{M}_{\text{sub}}\}) dN}, \quad (65)$$

which is the odds of finding absorbers in the range of  $\log_{10} N_{\text{H I}} \in [19.5, 20]$  compared to finding absorbers in  $\log_{10} N_{\text{H I}} \in [20, 23]$ . Note that we will treat the model posteriors of the sub-DLA model as part of the non-detections of DLAs in the following analysis sections.

## 8 CATALOGUE

The original parameter prior in Garnett et al. (2017) is uniformly distributed in  $z_{\text{DLA}}$  between the Lyman limit ( $\lambda_{\text{rest}} = 911.76 \text{ \AA}$ ) and the Ly  $\alpha$  emission of the QSO. In Bird et al. (2017), we chose the minimum value of  $z_{\text{DLA}}$  to be at the Ly  $\beta$  emission line of the QSO rest frame (instead of the Lyman limit) to avoid the region containing unmodelled Ly  $\beta$  forest. The primary reason for this was that the original absorption noise model did not include Ly  $\beta$  absorption. With the updated model from equation (24) we are able to model this absorption. Hence, for our new public catalogue, we sample  $z_{\text{DLA}}$  to be from Ly  $\infty$  to Ly  $\alpha$  in the QSO rest frame and for the convenience of future investigators our public catalogue contains DLAs throughout the whole available spectrum, including Ly  $\beta$  to Ly  $\infty$ . There is still some contamination in the blue end of high-redshift spectra from the Ly  $\beta$  forest and occasional Lyman breaks from a misestimated QSO redshift. In practice we shall see that the contamination is not severe except for  $z_{\text{DLA}} > 3.75$ . However, in the interest of obtaining as reliable DLA statistics as possible, when computing population statistics we consider only 3000  $\text{\AA}$  redward of Ly  $\beta$  to 3000  $\text{\AA}$  blueward of Ly  $\alpha$  in the QSO rest frame.

In this paper, we computed the posterior probability of  $\mathcal{M}_{\text{-DLA}}$  to  $\mathcal{M}_{\text{DLA}(k)}$  models. For each spectrum, the catalogue includes:

- (i) The range of redshift DLA searched  $[z_{\text{min}}, z_{\text{max}}]$ ,
- (ii) The log model priors from  $\log \Pr(\mathcal{M}_{\text{-DLA}} | z_{\text{QSO}})$ ,  $\log \Pr(\mathcal{M}_{\text{sub}} | z_{\text{QSO}})$ , to  $\log \Pr(\{\mathcal{M}_{\text{DLA}(i)}\}_{i=1}^k | z_{\text{QSO}})$ ,
- (iii) The log model evidence  $\log p(\mathbf{y} | \boldsymbol{\lambda}, \mathbf{v}, z_{\text{QSO}}, \mathcal{M})$ , for each model we considered,
- (iv) The model posterior  $\Pr(\mathcal{M} | \mathcal{D}, z_{\text{QSO}})$ , for each model we considered,
- (v) The probability of having DLAs  $\Pr(\{\mathcal{M}_{\text{DLA}}\} | \mathcal{D}, z_{\text{QSO}})$ ,
- (vi) The probability of having zero DLAs  $\Pr(\mathcal{M}_{\text{-DLA}} | \mathcal{D}, z_{\text{QSO}})$ ,
- (vii) The sample log likelihoods  $\log p(\mathbf{y} | \boldsymbol{\lambda}, \mathbf{v}, z_{\text{QSO}}, \{z_{\text{DLA}(i)}\}_{i=1}^k, \{\log_{10} N_{\text{H I}(i)}\}_{i=1}^k, \mathcal{M}_{\text{DLA}(k)})$  for all DLA models we considered, and
- (viii) The maximum a posteriori (MAP) values of all DLA models we considered.

The full catalogue will be available alongside the paper: [http://tiny.cc/multidla\\_catalog\\_gp\\_dr12q](http://tiny.cc/multidla_catalog_gp_dr12q). The code to reproduce the entire catalogue will be posted in [https://github.com/rmgarnett/gp\\_dla\\_detection/tree/master/multi\\_dlas](https://github.com/rmgarnett/gp_dla_detection/tree/master/multi_dlas).

### 8.1 Running time

We ran our multi-DLA code on UCR's High-Performance Computing Center (HPCC) and Amazon Elastic Compute Cloud (EC2). The computation of model posteriors of  $\mathcal{M}_{\text{-DLA}}$ ,  $\mathcal{M}_{\text{sub}}$ ,  $\{\mathcal{M}_{\text{DLA}(i)}\}_{i=1}^4$  takes 7–11 s per spectrum on a 32-core node in HPCC and 3–5 seconds on a 48-core machine in EC2. For each spectrum, we have to compute  $10\,000 \times 5 + 1$  log likelihoods in the form of equation (11). If we scale the sample size from  $N = 10\,000$  to 100 000, it costs 38–52 s on a 32-core node in HPCC.

## 9 EXAMPLE SPECTRA

Here, we show a few examples of the fitted GP priors, both to compare our method to others and to aid the reader in understanding concretely how our method works.

We show an exmple where our older model detected only one DLA, as shown in Fig. 4, while our new code detects three DLAs in this single spectrum as shown in Fig. 5. Because the mean QSO model

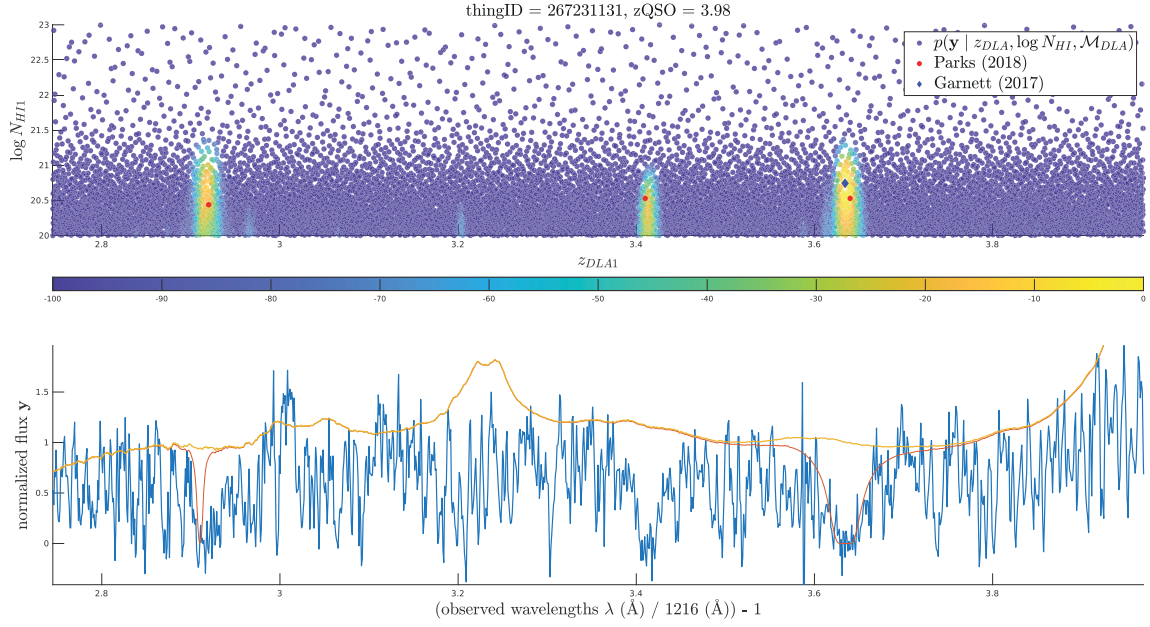
includes a redshift-dependent term corresponding to intervening absorbers, our new mean model can now fit the mean observed QSO spectrum better. Although we show the sample likelihoods in the  $\mathcal{M}_{\text{DLA}(1)}$  parameter space, our current code finds these three DLAs in the six-dimensional parameter space  $(z_{\text{DLA}(i)}, \log_{10} N_{\text{H I}(i)})_{i=1}^3$ .

In Fig. 6, we show a representative sample of a very common case in our  $\mathcal{M}_{\text{DLA}(1)}$  model. The red curve represents our GP prior on the given spectrum, and the orange curve is the curve with fitted DLAs provided by the CNN model presented in Parks et al. (2018).<sup>6</sup> We found Parks et al. (2018) underestimated the column densities of the underlying DLAs in the spectra due to not modelling Lyman  $\beta$  and Lyman  $\gamma$  absorption in DLAs, while the predictions of  $N_{\text{H I}}$  in our model are more robust since the predicted  $N_{\text{H I}}$  is constrained by  $\alpha$ ,  $\beta$ , and  $\gamma$  absorption. In the spectrum, Lyman  $\beta$  absorption is clearly visible (although noisy). In Fig. 6, Parks et al. (2018) have actually mistaken the Ly  $\gamma$  absorption line of the DLA for another, weaker, DLA. This demonstrates again the necessity of including other Lyman-series members in the modelling steps. Since Parks et al. (2018) broke down each spectrum into pieces during the training and testing phases, it is impossible for the CNN to use knowledge about other Lyman series lines associated with the DLAs. Another example, from a spectrum where we detect two DLAs and the CNN detects four (although at low significance) is shown in Fig. 7. Here, the CNN has mistaken both the Ly  $\beta$  and Ly  $\gamma$  absorption associated with the large DLA at  $z \sim 3$  (near the QSO rest frame) for separate DLAs at  $z = 2.4$  and  $z = 2.22$ , respectively. The large DLA at  $z \sim 3$  has been split into two of reduced column density and reduced confidence. The CNN has also missed the second genuine DLA at a rest-frame wavelength of 1025  $\text{\AA}$ , presumably due to the proximity of an emission line. Our code, able to model the higher order Lyman lines, has used the information contained within them to correctly classify this spectrum as containing two DLAs.

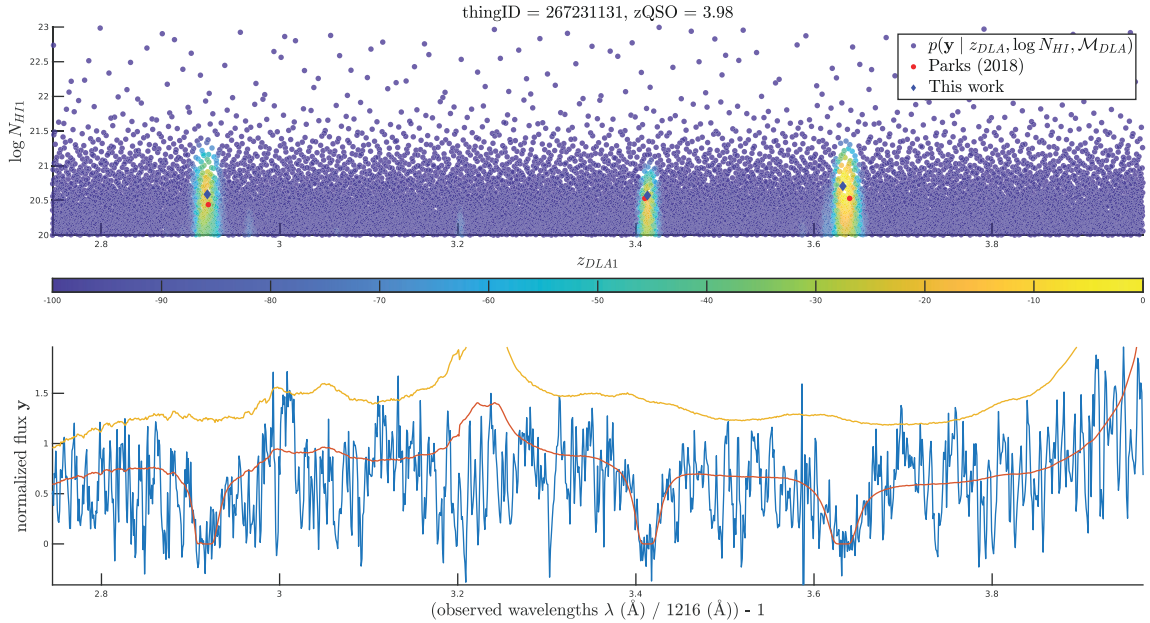
Fig. 9 shows an example that was problematic in both the models of Garnett et al. (2017) and Parks et al. (2018). This is an extremely noisy spectrum, where the length of the spectrum is not long enough for us to contain higher order Lyman-series absorption or even to see the full length of the putative Lyman  $\alpha$  absorption. By eye, distinguishing a DLA from the noise is challenging. If we examine the sample likelihoods from our model (shown in Fig. 10), we see that the DLA posterior probability is spread over the whole of parameter space; in other words, all models are a poor fit for this noise-dominated spectrum. The model selection is thus really comparing the likelihood function on the basis of how much parametric freedom it has. After implementing the additional Occam's razor factor between the null model and parametrized models (DLAs and sub-DLAs) described in Section 6.4, we found that the large DLA fitted to the noisy short spectrum by Garnett et al. (2017) was no longer preferred. This indicates that our Occam's razor penalty is effective. As shown in Fig. 16,  $\Omega_{\text{DLA}}$  at low redshifts is lower than the measurements in Bird et al. (2017), indicating that this class of error is common enough to have a measurable effect on the column density function. We checked that the addition of the Occam's razor penalty,  $\Omega_{\text{DLA}}$  is insensitive to the noise threshold used when selecting the spectra for our sample.

There are still some very high redshift QSOs ( $z_{\text{QSO}} \gtrsim 5$ ) where our code clearly detects too many DLAs in a single spectrum, even at low redshift. We exclude these spectra from our population statistics. At

<sup>6</sup>We used the version of Parks et al. (2018)'s catalogue listed in the published paper and found on Google Drive at <https://tinyurl.com/cnn-dlas>.



**Figure 4.** An example of finding DLAs using Garnett et al. (2017)’s model. Here, we use the single-DLA per spectrum version of Garnett’s model. Upper: sample likelihoods  $p(y | \theta, \mathcal{M}_{DLA})$  in the parameter space  $\theta = (z_{DLA}, \log_{10} N_{HI})$ . Red dots show the DLAs predicted by Parks et al. (2018), and the blue squares show the MAP prediction of the Garnett et al. (2017). Bottom: the observed spectrum (blue), the null model GP prior (orange), and the DLA model GP prior (red). So that the upper and bottom panels have the same x-axis, we rescale the observed wavelength to absorber redshift.



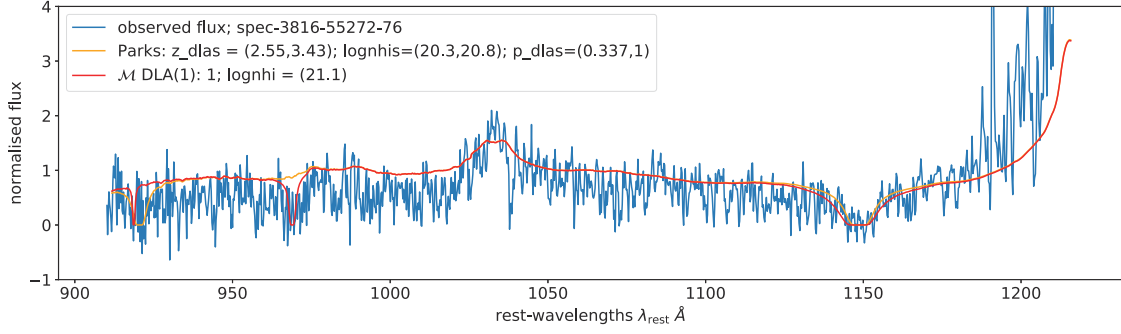
**Figure 5.** The same spectrum as Fig. 4, but using the multi-DLA model reported in this paper. Upper: sample likelihoods  $p(y | \theta, \mathcal{M}_{DLA})$  in the parameter space of the  $\mathcal{M}_{DLA(1)}$ , with  $\theta = (z_{DLA}, \log_{10} N_{HI})$ . Bottom: the observed spectrum (blue), the null model GP prior before the suppression of effective optical depth (orange), and the multi-DLA GP prior (red). The orange curve is slightly higher than the one in Fig. 4 because we try to model the mean spectrum before the forest. However, the DLA QSO model (red curve) matches the level of the observed mean flux better than Fig. 4 due to the inclusion of a term for the effective optical depth of the Lyman  $\alpha$  forest.

high redshift, the Lyman  $\alpha$  forest absorption is so strong as to render the observed flux close to zero. We thus cannot easily distinguish between the null model and the DLA models. It is also possible that at high redshifts, the mean flux of the forest is substantially different from the Kim et al. (2007) model we assume, and that this biases the fit. Finally, there are few such spectra, and so we cannot rule

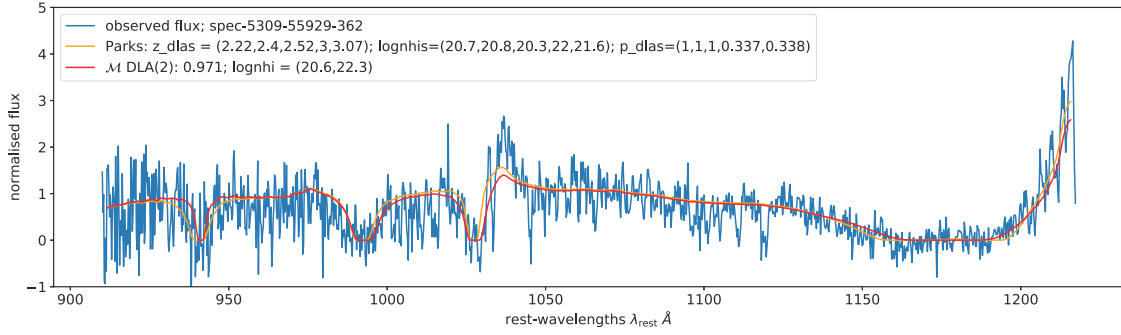
out the possibility that covariance of their emission spectra differs quantitatively from lower redshift QSOs.

## 10 ANALYSIS OF THE RESULTS

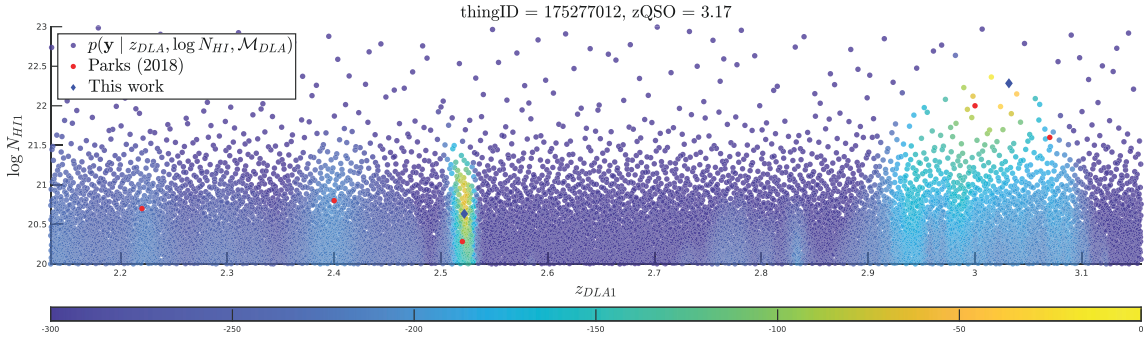
In this section, we present results from our classification pipeline, and we also present the statistical properties (CDDF, line densities



**Figure 6.** Blue: The normalized observed flux. The spectral ID represents `spec-plate-mjd-fiber_id`. Yellow: Parks' predictions on top of our null model. Our model predicts only one DLA while the CNN model in Parks et al. (2018) predicts two DLAs. One of the DLAs predicted by Parks et al. (2018) is coincident with the  $\text{Ly}\gamma$  absorption from our predicted DLA.  $z_{\text{dla}}$  corresponds to the DLA redshifts reported in Parks' catalogue, and  $\log n_{\text{hi}}$  corresponds to the column density estimations of Parks' catalogue.  $p_{\text{dla}}$  is the `dla_confidence` reported in Parks. Red: Our current model with the highest model posterior and the MAPs of column densities. In this spectrum, we show that it is crucial to include  $\text{Ly}\beta$  and  $\text{Ly}\gamma$  absorption from the DLA in the DLA profile. It not only helps to localize the DLA, but it also predicts  $N_{\text{H I}}$  more accurately using information from the  $\text{Ly}\beta$  region. The blue line shows the observed flux, the red curve is our multi-DLA GP prior, and the orange curve shows the predicted DLAs from Parks et al. (2018) subtracted from our mean model.



**Figure 7.** A spectrum in which we detect two DLAs. Blue: Normalized flux. Red: GP mean model with two intervening DLAs. Yellow: The predictions from Parks' catalogue. Pink: The MAP prediction of Garnett et al. (2017) on top of the GP mean model without mean flux suppression. The model posterior from Garnett et al. (2017) is listed in the legend (1) with the MAP value of  $\log_{10} N_{\text{H I}}$ . The column density estimate for the DLA near  $\lambda_{\text{rest}} = 1025 \text{ \AA}$  has large uncertainty (see Fig. 8). It is thus possible that this DLA could be a sub-DLA, as preferred by Parks et al. (2018).



**Figure 8.** The log sample likelihoods for the DLA model of the spectrum shown in Fig. 7, normalized to range from  $-\infty$  to 0. The DLA at  $z_{\text{DLA}} \sim 2.52$  could be a sub-DLA (as preferred by Parks et al. 2018), as the  $\log_{10} N_{\text{H I}}$  estimate is uncertain. However, we found that the two-DLA model posterior  $\log p(\mathcal{M}_{\text{DLA}(2)} | \mathbf{y}, \lambda, \mathbf{v}, z_{\text{QSO}}) = -638$  is still higher than the model posterior from combining one-DLA and one-sub-DLA, which is  $\log p(\mathcal{M}_{\text{DLA}(1)} + \mathcal{M}_{\text{sub}} | \mathbf{y}, \lambda, \mathbf{v}, z_{\text{QSO}}) = -691.47$ .

$dN/dX$ , and total column densities  $\Omega_{\text{DLA}}$  of the DLAs detected in our catalogue.

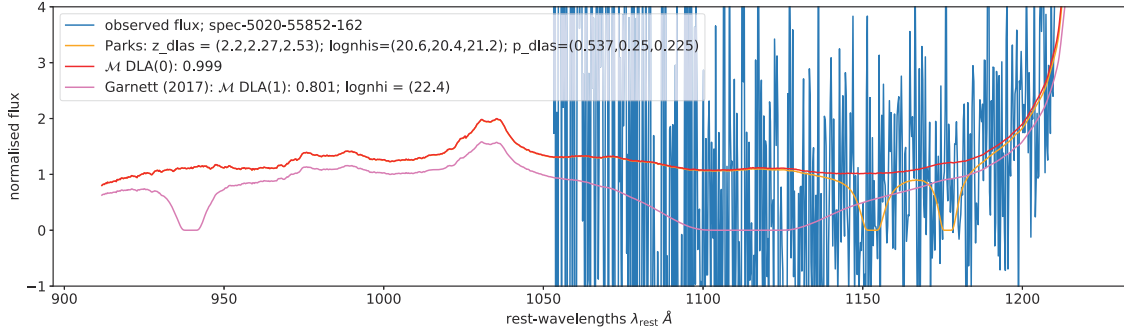
### 10.1 ROC analysis

To evaluate how well our multi-DLA classification reproduces earlier results, we rank our DLA detections using the log posterior odds

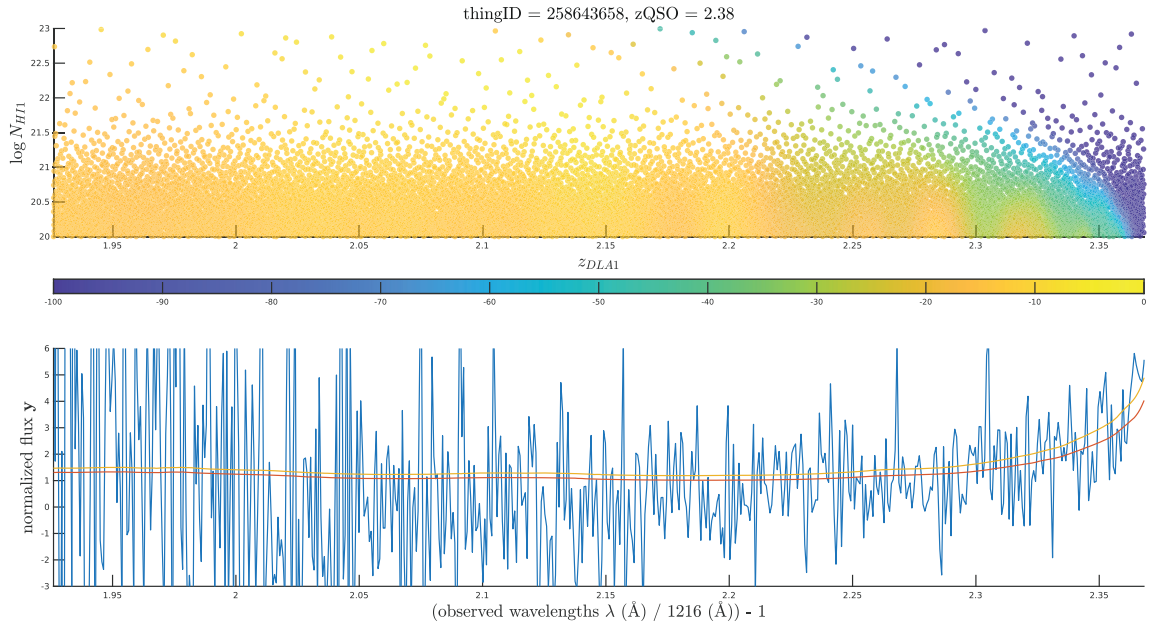
between the DLA model (summing up all possible DLA models  $\{\mathcal{M}_{\text{DLA}(i)}\}_{i=1}^k$ ) and the null model:

$$\begin{aligned} \log(\text{odds}) &= \log \Pr(\{\mathcal{M}_{\text{DLA}}\} | \mathcal{D}, z_{\text{QSO}}) - \log \Pr(\mathcal{M}_{-\text{DLA}} | \mathcal{D}, z_{\text{QSO}}), \end{aligned} \quad (66)$$

where the ranking is over all sightlines. From the top of the ranked list based on the log posterior odds, we calculate the true positive



**Figure 9.** A noisy spectrum at  $z_{\text{QSO}} = 2.378$  fitted with a large DLA by Garnett et al. (2017). Red: The model presented in this paper predicts no DLA detection in this spectrum. Pink: The MAP prediction of Garnett et al. (2017) on top the GP mean model without the mean-flux suppression. Gold: The prediction of Parks et al. (2018) subtracted from our mean model. Note that Parks et al. (2018) also indicate a detection of a DLA at  $z_{\text{DLA}} = 2.53$ , but outside the range of this spectrum.



**Figure 10.** Top: The sample likelihoods of the spectrum shown in Fig. 9. The colour bar indicates the normalized log likelihoods ranging from  $-\infty$  to 0. Bottom: The orange curve indicates the GP mean model before mean-flux suppression, the red curve represents the mean model after suppression, and the blue line is the normalized flux of this spectrum. The x-axis of this spectrum is rescaled to be the same as the  $z_{\text{DLA}}$  presented in the upper panel.

rate and false positive rate for each rank:

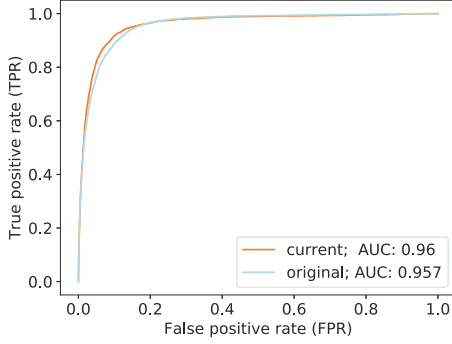
$$\begin{aligned} \text{TPR} &= \frac{\text{TP}}{\text{TP} + \text{FN}}; \\ \text{FPR} &= \frac{\text{FP}}{\text{FP} + \text{TN}}. \end{aligned} \quad (67)$$

The true positive rate is the fraction of sightlines where we detect DLAs (ordered by their rank) divided by the number of sightlines with DLAs detected by earlier catalogues. The false positive rate is the number of detections of DLAs divided by the number of sightlines where earlier catalogues did not detect DLAs. In Fig. 11, we show the TPR and FPR in a receiver-operating characteristics (ROC) plot to show how well our classification performs. We have compared to the concordance DLA catalogue (Lee et al. 2013) in the hope that it approximates ground truth, there being no completely reliable DLA catalogue.

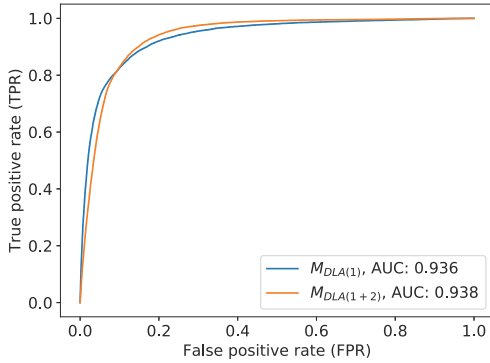
We also want to know how well our pipeline can identify the number of DLAs in each spectrum. The DR9 concordance catalogue

does not count multiple DLA spectra, and so we compare our multi-DLA detections to the catalogue published by Parks et al. (2018). Each DLA detected in Parks et al. (2018) comes with a measurement of their confidence of detection ( $\text{dla\_confidence}$  or  $p_{\text{DLA}}^{\text{Parks}}$ ) and an MAP redshift and column density estimate. We compare our multiple DLA catalogue to those spectra with  $p_{\text{DLA}}^{\text{Parks}} > 0.98$ . The resulting ROC plot is shown in Fig. 12. We count a maximum of two DLAs in each spectrum: three or more DLAs in a single sightline are extremely rare and do not provide a large enough sample for an ROC plot. Parks' catalogue is not a priori more reliable than ours, especially in spectra with multiple DLAs, but comparing the first two DLAs is a reasonable way to validate our method's ability to detect multiple DLAs.

These spectra are counted by breaking down each two-DLA sightline (either in Parks or our catalogue) into two single observations. For example, if there are two DLAs detected in Parks and one DLA detected in our pipeline for an observation  $\mathcal{D}$ , we will assign one



**Figure 11.** The ROC plot made by ranking the sightlines in BOSS DR9 samples using the log posterior odds of containing at least one DLA. Ground truths are from the DR9 concordance catalogue. The orange curve shows the ROC plot of our current multi-DLA model, and the blue curve is derived from Garnett et al. (2017). In this plot, we consider only the model containing at least one DLA  $p(\{\mathcal{M}_{\text{DLA}}\} | \mathcal{D})$ , rather than the multiple DLA models, as the concordance catalogue contains only one DLA per spectrum.



**Figure 12.** The ROC plot for sightlines with one and two DLA detections, by using the catalogue of Parks et al. (2018) (with `dla_confidence` > 0.98) as ground truth.

ground-truth detection to  $p(\mathcal{M}_{\text{DLA}(1)} | \mathcal{D})$  and assign one ground-truth detection to  $p(\mathcal{M}_{\text{-DLA}} | \mathcal{D})$ . On the other hand, if there is only one DLA detected in Parks and two DLAs detected in our pipeline, we will assign one ground-truth detection to  $p(\mathcal{M}_{\text{DLA}(2)} | \mathcal{D})$  and one ground-truth non-detection to  $p(\mathcal{M}_{\text{DLA}(2)} | \mathcal{D})$ .

In Fig. 13, we also analyse the MAP estimate of the parameters ( $z_{\text{DLA}}$ ,  $\log_{10} N_{\text{H I}}$ ) by comparing with the reported values in DR9 concordance DLA catalogue. The median difference between these two is  $-2.2 \times 10^{-4}$  ( $-66.6 \text{ km s}^{-1}$ ) and the interquartile range is  $2.2 \times 10^{-3}$  ( $662 \text{ km s}^{-1}$ ). For the log column density estimate, the median difference is 0.040, and the interquartile range is 0.26. The medians and interquartile ranges of the MAP estimate are very similar to the values reported in Garnett et al. (2017) with the median of  $z_{\text{DLA}}$  slightly smaller and the median of  $\log_{10} N_{\text{H I}}$  slightly larger. Note that the DR9 concordance catalogue is not the ground truth, so small variations in comparison to Garnett et al. (2017) can be considered to be negligible. As shown in Fig. 13, both histograms are roughly diagonal, although the scatter in column density MAP is large. Note that our DLA-detection procedure is designed to evaluate the model evidence across all of parameter space: a single sample MAP cannot convey the full posterior probability distribution. In Section 10.2, we thus describe a procedure to propagate the posterior density in the parameter space directly to column density statistics.

## 10.2 CDDF analysis

We follow Bird et al. (2017) in calculating the statistical properties of the modified DLA catalogue presented in this paper. We summarize the properties of DLAs using the averaged binned CDDF, the incident probability of DLAs ( $dN/dX$ ), and the averaged matter density as a function of redshift ( $\Omega_{\text{DLA}}(z)$ ).

To plot these summary statistics, we need to convert the probabilistic detections in the catalogue to the expected average number of DLAs and their corresponding variances. We first describe how we compute the expected number of DLAs in a given column density and redshift bin. Next, we show how we derive the CDDF,  $dN/dX$ , and  $\Omega_{\text{DLA}}(z)$  from the expected number of DLAs. A sample of  $n$  observed spectra contains a sequence of  $n$  model posteriors  $p_{\text{DLA}}^1, p_{\text{DLA}}^2, \dots, p_{\text{DLA}}^n$  defined by

$$p_{\text{DLA}}^i = p(\{\mathcal{M}_{\text{DLA}}\} | y_i, \lambda_i, v_i, z_{\text{QSO}i}), \quad (68)$$

where  $i = 1, 2, \dots, n$  is the index of the spectrum, and the DLA model here includes all computed DLA models  $\{\mathcal{M}_{\text{DLA}}\} = \{\mathcal{M}_{\text{DLA}(i)}\}_{i=1}^k$ , so that  $k = 4$  is the maximum possible number of DLAs in each spectrum in our model.

Suppose the region of interest is in a specific bin  $\Theta$ , an interval in the parameter space of column density or DLA redshift  $\Theta \in \{N_{\text{H I}}, z_{\text{DLA}}\}$ . To compute the posterior of having DLAs in each spectrum in a given bin  $\Theta$ ,  $p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | \Theta)$ , we integrate over the sample likelihoods in the bin and multiply the model posterior by the total  $p_{\text{DLA}}^i$  for spectrum  $i$ :

$$p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | \Theta) \propto p_{\text{DLA}}^i \times \int_{\Theta} p(y_i | \{\mathcal{M}_{\text{DLA}}\}, \lambda_i, v_i, z_{\text{QSO}i}, \theta) d\theta. \quad (69)$$

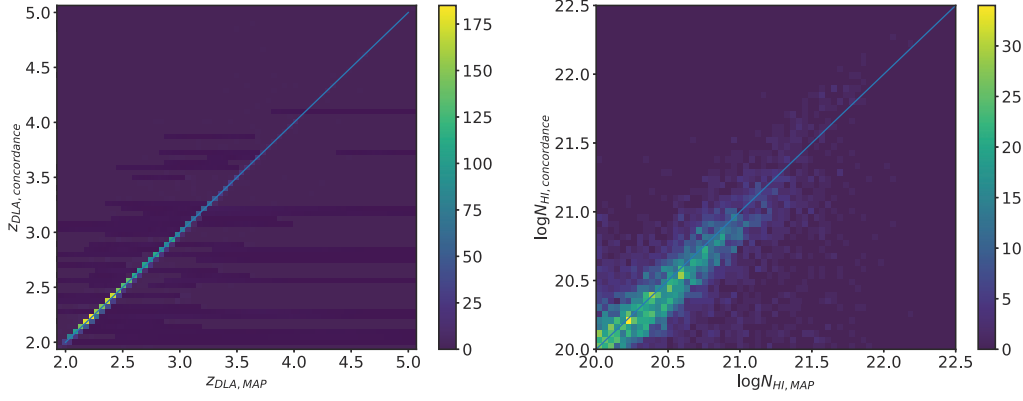
$\theta$  is either  $z_{\text{DLA}}$  or  $\log_{10} N_{\text{H I}}$  and  $\theta \in \Theta = (\Theta, \bar{\Theta})$ .

We calculate the posterior probability of having  $N$  DLAs by noting that the full likelihood follows the Poisson–Binomial distribution. Consider a sequence of trials with a probability of success equal to  $p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | \Theta) \in [0, 1]$ . The probability of having  $N$  DLAs out of a total of  $n$  trials is the sum of all possible  $N$  DLAs subsets in the whole sample:

$$\Pr(N) = \sum_{\text{DLA} \in F_N} \prod_{i \in \text{DLA}} p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | \Theta) \times \prod_{j \in \text{DLA}^c} (1 - p_{\text{DLA}}^j(\{\mathcal{M}_{\text{DLA}}\} | \Theta)), \quad (70)$$

where  $F_N$  corresponds to all subsets of  $N$  integers that can be selected from the sequence  $\{1, 2, \dots, n\}$ . The above expression means we select all possible  $N$  choices from the entire sample, calculate the probability of those  $N$  choices having DLAs and multiply that by the probability of the other  $n - N$  choices having no DLAs. If all  $p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | \Theta)$  are equal, the Poisson–Binomial distribution reduces to a Binomial distribution.

The above Poisson–Binomial distribution is not trivial to compute given our large sample size. The technical details of how to evaluate equation (70) efficiently are described in Bird et al. (2017). In short, we use Le Cam (1960)’s theorem to approximate those spectra with  $p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | \Theta) < p_{\text{switch}} = 0.25$  by an ordinary Poisson distribution, and evaluate the remaining samples with the discrete Fourier transform (Fernandez & Williams 2010). Our catalogue contains the posteriors of samples in a given spectrum. Combined with the above probabilistic description of the total number of DLAs in the entire sample, we are able to obtain not only the point estimation of  $\Pr(N)$  but also its probabilistic density interval.



**Figure 13.** The MAP estimates of the DLA parameters  $\theta = (z_{\text{DLA}}, \log_{10} N_{\text{HI}})$  for DLAs detected by our model in spectra observed by SDSS DR9, compared to the values reported in the concordance catalogue. The straight line indicates a perfect fit. Note that the concordance  $\log_{10} N_{\text{HI}}$  values are not ground truth, so the scatter in column density predictions was expected.

We thus compute the CDDF in a given bin  $\Theta = N_{\text{HI}} \in [N_{\text{HI}}, N_{\text{HI}} + \Delta N_{\text{HI}}]$  with

$$f(N) = \frac{F(N)}{\Delta N \Delta X(z)}, \quad (71)$$

where  $F(N) = \mathbb{E}(N \mid N_{\text{HI}} \in [N_{\text{HI}}, N_{\text{HI}} + \Delta N_{\text{HI}}])$  is the expected number of absorbers at a given sightline within a column density interval. Thus, the CDDF  $f(N)$  is the expected number of absorbers per unit column density per unit absorption distance, within a given column density bin.

The definition of absorption distance  $\Delta X(z)$  is

$$X(z) = \int_0^z (1+z')^2 \frac{H_0}{H(z')} dz', \quad (72)$$

which includes the contributions of the Hubble function  $H^2(z)/H_0^2 = \Omega_{\text{M}}(1+z)^3 + \Omega_{\Lambda}$ , with  $\Omega_{\text{M}}$  the matter density and  $\Omega_{\Lambda}$  the dark energy density.

The incident rate of DLAs  $dN/dX$  is defined as

$$\frac{dN}{dX} = \int_{10^{20.3}}^{\infty} f(N \mid N_{\text{HI}}, X \in [X, X + dX]) dN_{\text{HI}}, \quad (73)$$

which is the expected number of DLAs per unit absorption distance.

The total column density  $\Omega_{\text{DLA}}$  is defined as

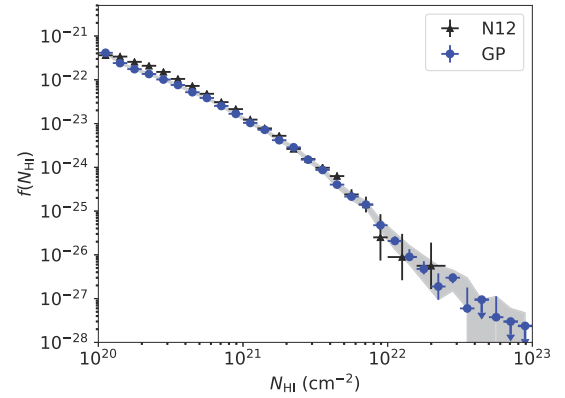
$$\Omega_{\text{DLA}} = \frac{m_{\text{P}} H_0}{c \rho_{\text{c}}} \int_{10^{20.3}}^{\infty} N_{\text{HI}} f(N \mid N_{\text{HI}}, X \in [X, X + dX]) dN_{\text{HI}}, \quad (74)$$

where  $\rho_{\text{c}}$  is the critical density at  $z = 0$  and  $m_{\text{P}}$  is the proton mass.

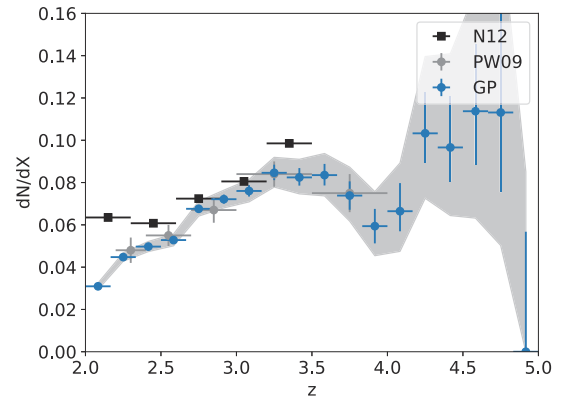
### 10.3 Statistical properties of DLAs

Based on the above calculations, we show our CDDF in Fig. 14,  $\frac{dN}{dX}$  in Fig. 15, and  $\Omega_{\text{DLA}}$  in Fig. 16.<sup>7</sup> Note that for determining the statistical properties of DLAs, we limit the samples of  $z_{\text{DLA}}$  to the range redward of the Lyman  $\beta$  in the QSO rest frame, as in Bird et al. (2017).

Fig. 14 shows the CDDF from our DR12 catalogue in comparison to the DR9 catalogue of Noterdaeme et al. (2012). Our CDDF analysis combines all spectral paths with QSO redshift smaller than 5,  $z_{\text{DLA}} < 5$ . The CDDF statistics are dominated by the low-redshift

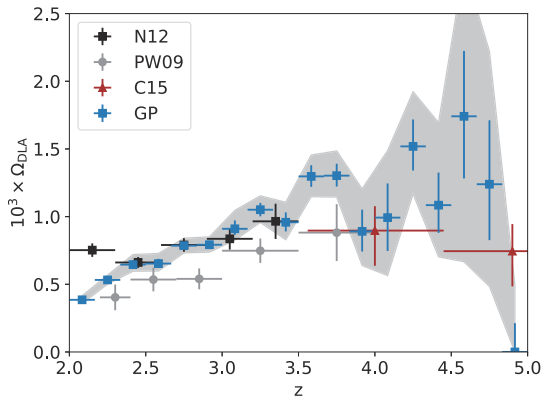


**Figure 14.** The CDDF based on the posterior densities for at least one DLA (blue, ‘GP’). The DLAs are derived from SDSS DR12 spectra using the method presented in this paper. We integrate all spectral lengths with  $z < 5$ . We also plot the CDDF of Noterdaeme et al. (2012) (N12; black) as a comparison. The error bars represent the 68 per cent confidence limits, while the grey filled band represents the 95 per cent confidence limits. Note that our CDDF completely overlaps with those of N12 for column densities in the range  $10^{21} \text{ cm}^{-2} < N_{\text{HI}} < 10^{22} \text{ cm}^{-2}$ .

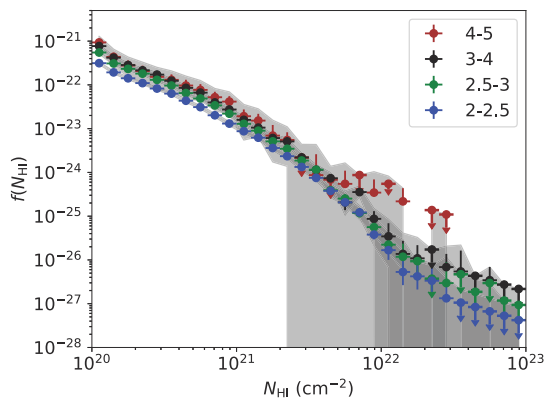


**Figure 15.** The line density of DLAs as a function of redshift from our DR12 multi-DLA catalogue (blue, ‘GP’). We also plot the results of Noterdaeme et al. (2012) (N12; black) and Prochaska & Wolfe (2009) (PW09; grey). Note that statistical error was not computed in Noterdaeme et al. (2012).

<sup>7</sup>The table files to reproduce Figs 14–16 will be posted in [http://tiny.cc/multiidla.catalog\\_gp.dr12q](http://tiny.cc/multiidla.catalog_gp.dr12q).



**Figure 16.** The total H I density in DLAs,  $\Omega_{\text{DLA}}$ , from our DR12 multi-DLA catalogue as a function of redshift (blue, ‘GP’), compared to the results of Noterdaeme et al. (2012) (N12; black), Prochaska & Wolfe (2009) (PW09; grey), and Crighton et al. (2015) (C15; red).



**Figure 17.** The redshift evolution (or non-evolution) of the CDDF. Labels show the absorber redshift ranges used to plot the CDDFs. In column density and redshift ranges with no detection at 68 per cent confidence, a down-pointing arrow is shown indicating the 68 per cent upper limit.

absorbers, as demonstrated in Fig. 17. The error bars represent the 68 per cent confidence interval, while the grey shaded area encloses the 95 per cent highest density region. The CDDF values in Fig. 14 are calculated from the posterior distribution directly. We note that there are only two DLAs with  $\text{MAP } \log_{10} N_{\text{HI}} > 22.5$  in our catalogue with high confidence ( $p_{\text{DLA}} > 0.99$ ). The non-zero values in the CDDF are due to uncertainty in  $\log_{10} N_{\text{HI}}$ , not positive detections.

Noterdaeme et al. (2012) contains multi-DLAs, but, as described in section 2.2 in their paper, they applied a stringent cut on their samples with  $\text{CNR} > 3$ , where CNR refers to the continuum-to-noise ratio. The CDDF of N12 in Fig. 14 is thus a subsample of their catalogue. We, on the other hand, use all data even those with low signal-to-noise ratios. Comparing to our previously published CDDF (Bird et al. 2017), the CDDF in this paper shows DLA detections at low  $N_{\text{HI}}$  are consistent with Noterdaeme et al. (2012). Introducing the sub-DLA as an alternative model successfully regularizes detections at  $\sim 10^{20} \text{ cm}^{-2}$ .<sup>8</sup>

Fig. 15 shows the line density of DLAs. Our results are again consistent with those of Prochaska & Wolfe (2009) and Noterdaeme et al. (2012) where they both agree. Our detections are between those

two catalogues at low-redshift bins and consistent with Prochaska & Wolfe (2009) in the highest redshift bin. Comparing to our previous  $dN/dX$  (Bird et al. 2017), we moderately regularize the detections of DLAs at high redshifts. This change shows that changing the mean model of the GP to include the mean flux absorption prevents the pipeline confusing the suppression due to the Lyman alpha forest with a DLA. While the change of posterior modes in  $dN/dX$  is large at high-redshift bins, we note that those changes are mostly within 95 per cent confidence interval of our previously published line densities. All analyses shown measure a peak in  $dN/dX$  at  $z \sim 3.5$ . This may be partially due to  $z_{\text{DLA}} = 3.5$  the SDSS colour selection algorithm systematic identified by Prochaska, Worsack & O’Meara (2009), which oversamples Lyman-limit systems (LLS), especially near the QSO, in the redshift range 3.0–3.6 (Worsack & Prochaska 2011; Fumagalli et al. 2013). Note, however, that in our analysis neighbouring redshift bins are highly correlated and so a statistical fluctuation is also a valid explanation. We have checked visually that our sub-DLA model successfully models spectra with an LLS in the proximate zone of the QSO emission peak.

Fig. 16 shows the total column density  $\Omega_{\text{DLA}}$  in DLAs in units of the cosmic density. Our results are mostly consistent with Noterdaeme et al. (2012) although we have slightly lower  $\Omega_{\text{DLA}}$  at  $z \sim 2$ . This is due to our Occam’s razor penalty, which suppresses DLAs in spectra which are not long enough to include the full width of the DLA. Since these are all low-redshift QSOs, this suppresses DLA detections at  $z < 2.3$ . As discussed in Noterdaeme et al. (2012), Sánchez-Ramírez et al. (2016), and Bird et al. (2017), the relatively low  $\Omega_{\text{DLA}}$  of Prochaska & Wolfe (2009) is due to the smaller sample size of the SDSS DR5 data set. We also compare our  $\Omega_{\text{DLA}}$  to that measured by Crighton et al. (2015) at high redshifts ( $z = 4$  and  $z = 5$ ). Crighton et al. (2015) used a small but higher signal-to-noise data set. Our results at  $z = 4$  and  $z = 5$  are consistent with those from Crighton et al. (2015). However, we note that the relatively small sample of Crighton et al. (2015) may bias it slightly low, as contributions from DLAs with  $N_{\text{HI}}$  higher than expected to be in the survey will not be included in their  $\Omega_{\text{DLA}}$  estimate. Our Bayesian analysis includes possible contributions of undetected DLAs with column density up to  $\log_{10} N_{\text{HI}} = 23$  in the error bars via the prior on the column density.

Compared to our previously published  $\Omega_{\text{DLA}}$  (Bird et al. 2017), we found a reduction in  $\Omega_{\text{DLA}}$  between  $z = 4$  and  $z = 5$ . This is due to the incorporation of a better mean flux vector model, which reduces the posterior density of high column density systems for high-redshift absorbers (although within the 95 per cent confidence bars of the earlier work). Our confidence intervals are also substantially smaller for  $z_{\text{DLA}} \gtrsim 3.7$  than in (Bird et al. 2017). This is due to our inclusion, for the first time, of information from the Lyman  $\beta$  absorption of the DLAs, which both constrains DLA properties and helps to distinguish DLAs from noise fluctuations.

We have tested the robustness of our method with respect to spectra with different SNRs and found that, as in Bird et al. (2017), the statistical properties predicted by our method are uncorrelated with the QSO SNR. Furthermore, the presence of a DLA is uncorrelated with the QSO redshift, fixing a statistical systematic in the earlier work.

As a cross-check of our wider catalogue, we also tested the CDDF, line densities, and total column densities of the DLAs in our catalogue with a full range of  $z_{\text{DLA}}$ , from Ly  $\infty$  to Ly  $\alpha$ . The CDDF was very similar to the CDDF excluding the Ly  $\beta$  region shown in Fig. 14, but with a moderate increase at high column density.  $dN/dX$  was almost

<sup>8</sup>Note again the artefact at  $\sim 10^{20} \text{ cm}^{-2}$  will not affect the analyses of  $dN/dX$  or  $\Omega_{\text{DLA}}$  as the definition of a DLA is absorbers with  $N_{\text{HI}} > 10^{20.3} \text{ cm}^{-2}$ .

identical to Fig. 15, indicating that the detection of DLAs is robust even though we extend our sampling range to  $\text{Ly } \infty$ . However,  $\Omega_{\text{DLA}}$  increases for  $3.5 < z_{\text{DLA}} < 4.0$ . By visual inspection we found that this is due to the spectra where the QSO redshift from the SDSS pipeline in error and a Lyman break trough appears at the blue end of the spectrum in a region the code expects to contain only  $\text{Ly } \beta$  absorption. As our model does not account for redshift errors, it explains the absorption due to these troughs by DLAs.

#### 10.4 Comparison to Garnett’s catalogue

To understand the effect of the modifications we made to our model in this paper, we visually inspected a subset of spectra with high model posteriors of a DLA in Garnett et al. (2017) ( $p_{\text{DLA}}^{\text{Garnett}}$ ) but low model posteriors in our current model ( $p_{\text{DLA}}$ ). In particular, we chose spectra with ( $p_{\text{DLA}}^{\text{Garnett}} - p_{\text{DLA}} > 0.99$ ).

A large fraction of these spectra falls within the  $\text{Ly } \beta$  emission region. One plausible explanation is that the  $\text{Ly } \beta$  emission region has a higher noise variance, which makes it harder to distinguish the DLA and sub-DLA models. We also checked that we are not unfairly preferring the sub-DLA model during model selection. Our model selection uses the sub-DLA model only to regularize the DLA model and does not consider cases where DLAs and sub-DLAs occur in the same spectrum. Thus a spectrum with a clear detection of a sub-DLA could fail to detect a true DLA at a different redshift. In light of this, we also tested if combining multi-DLA models with a sub-DLA affects our results.

We modified the DLA model, assuming that the DLA and sub-DLA models are independent, to include the sub-DLA model prior. We then considered an iterative sampling procedure: First, we sampled the  $k$ -th DLA likelihood. Next we used the  $k$ -th DLA parameter posterior as a prior to sample  $\mathcal{M}_{\text{DLA}(k)}$  and combine  $\mathcal{M}_{\text{DLA}(k)}$  with the sub-DLA model via sampling a non-informative prior. The full procedure can be written as

$$p(\{\theta_i\}_{i=1}^k \mid \mathcal{M}'_{\text{DLA}(k)}, \mathcal{D}, z_{\text{QSO}}) = (1 + p(\theta_{\text{sub}} \mid \mathcal{M}_{\text{sub}}, z_{\text{QSO}})) \times p(\{\theta_i\}_{i=1}^k \mid \mathcal{M}_{\text{DLA}(k)}, \mathcal{D}, z_{\text{QSO}}), \quad (75)$$

For computational simplicity, we only consider the modified model until  $\mathcal{M}'_{\text{DLA}(3)}$ ; the probability of  $\mathcal{M}'_{\text{DLA}(4)}$  is expected to be insignificant comparing to the total DLA model posterior,  $p(\{\mathcal{M}_{\text{DLA}}\} \mid \mathcal{D}, z_{\text{QSO}})$ .

In practice, however, we found that this made a small difference to our results, only marginally modifying the ROC curve and CDDF. Moreover, the ability of the sub-DLA model to regularize low column density DLAs was reduced, so we have preserved our default model.

#### 10.5 Comparison to Parks catalogue

In this section, we compare our results with Parks et al. (2018). We first show the differences between our MAP predictions and Parks’ predictions for DLA redshift and column density. We required  $p_{\text{DLA}}^{\text{Parks}} > 0.98$ . We measured the difference in posterior parameters when both pipelines predicted one DLA. As shown in Fig. 18, both histograms are roughly symmetric. We measure small median offsets between two pipelines with

$$\begin{aligned} \text{median}(z_{\text{DLA}}^{\text{MAP}} - z_{\text{DLA}}^{\text{Parks}}) &= 0.00010; \\ \text{median}(\log N_{\text{H I}}^{\text{MAP}} - \log N_{\text{H I}}^{\text{Parks}}) &= 0.016. \end{aligned} \quad (76)$$

We also compared our absorber redshift measurements and column density measurements to Parks’ catalogue for those spectra that we

both agree contain two DLAs. The differences between these two have small median offsets of  $\Delta z_{\text{DLA}} = 0.000052$  and  $\Delta \log_{10} N_{\text{H I}} = 0.006$  (and dominated by low column density systems).

We show the disagreement between multi-DLA predictions for our catalogue and Parks’ catalogue in Table 1. Note that though the multi-DLA detections between our method and Parks do not completely agree, the level of disagreement is small: 6.1 per cent. Moreover, if Parks predicts one or two DLAs, our method generally detects one or two DLAs. There are, however, some spectra where we detected  $>2$  DLAs, but Parks detected none. To understand the statistical effect of this discrepancy, we compare our DLA properties to those reported by Parks et al. (2018). We plot the CDDF and  $dN/dX$  of that catalogue. We assume  $p_{\text{DLA}}^{\text{Parks}} > 0.9$  represents a DLA and use  $z_{\text{DLA}}$  and  $\log_{10} N_{\text{H I}}$  reported in their catalogue in JSON format.<sup>9</sup> To compute the sightline path searched over, we assume their CNN model was searching the range  $\text{Ly } \infty$  to  $\text{Ly } \alpha$  in the QSO rest frame. Note this differs slightly from Parks et al. (2018), section 3.2, where a sightline search radius ranging from 900 to 1346 Å in the QSO rest frame is given. However, we know the centres of DLAs should be at a redshift between  $\text{Ly } \infty$  and  $\text{Ly } \alpha$  in the rest frame and modify our search paths accordingly.

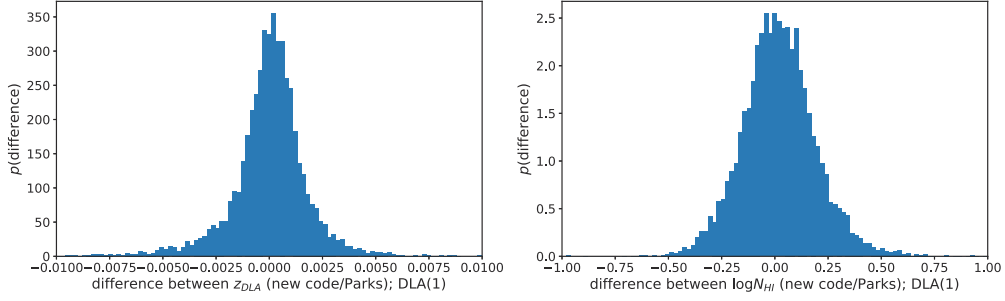
Fig. 19 shows that  $dN/dX$  is consistent with Noterdaeme et al. (2012) for  $z_{\text{DLA}} < 3.5$  (although lower than our measurement at higher redshift). The CNN is thus successfully detecting DLAs, especially the most common case of DLAs with a low column density. There are fewer DLAs detected at higher redshift, likely reflecting the increased difficulty for the CNN of distinguishing DLAs from the Lyman  $\alpha$  forest. This is discussed in Parks et al. (2018), who note that the CNN finds it difficult to detect a weak DLA in noisy spectra. However, as shown in Fig. 20, the CDDF measured by the CNN model is significantly discrepant with other surveys for large column densities. Note that the scale is logarithmic: the CNN is failing to detect  $>60$  per cent of DLAs with  $\log_{10} N_{\text{H I}} > 21$ . We noticed that large DLAs were often split into two objects with lower column density, which accounts for many of the discrepancies between our two data sets. We suspect this might be due to the limited size of the convolutional filters used by Parks et al. (2018). If the filter is not large enough to contain the full damping wings of a given DLA, the allowed column density would be artificially limited.

## 11 CONCLUSION

We have presented a revised pipeline for detecting DLAs in SDSS QSO spectra based on Garnett et al. (2017). We have extended the pipeline to reliably detect up to four DLAs per spectrum. We have performed modifications to our model for the Lyman  $\alpha$  forest to improve the reliability of DLA detections at high redshift and introduced a model for sub-DLAs to improve our measurement of low column density DLAs. Finally, we introduced a penalty on the DLA model based on Occam’s razor which meant that spectra for which both models are a poor fit generally prefer the no-DLA model.

Our results include a public DLA catalogue, with several examples shown above and further examples easily plotted using a PYTHON package. We have visually inspected several extreme cases to validate our results and compared extensively to several earlier DLA catalogues: the DR9 concordance catalogue (Lee et al. 2013) and a DR12 catalogue using a CNN (Parks et al. 2018). Our new pipeline had very good performance validated against both catalogues.

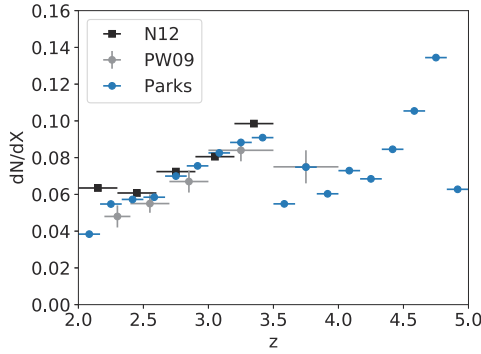
<sup>9</sup><https://tinyurl.com/cnn-dlas>



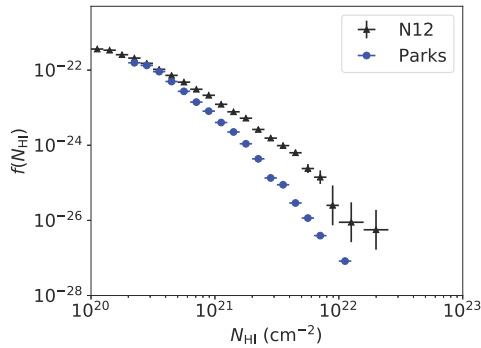
**Figure 18.** The difference between the MAP estimates of the DLA parameters  $\theta = (z_{\text{DLA}}, \log_{10} N_{\text{HI}})$ , against the predictions of Parks et al. (2018). We consider spectra that both catalogues agree contain one DLA.

**Table 1.** The confusion matrix for multi-DLAs detections between Garnett with multi-DLAs and Parks. Note we require both the model posteriors in Garnett and DLA confidence in Parks to be larger than 0.98. We also require  $\log_{10} N_{\text{HI}} > 20.3$ .

Parks Garnett with multi-DLAs	0 DLA	1 DLA	2 DLAs	3 DLAs	4 DLAs
0 DLA	138 726	6197	142	6	0
1 DLA	3050	8752	335	4	0
2 DLAs	293	570	566	28	0
3 DLAs	30	39	34	21	0
4 DLAs	5	9	6	1	0



**Figure 19.**  $dN/dX$  from Parks et al. (2018). The  $dN/dX$  agrees well with other surveys, but there is a moderate deficit of DLAs at high redshifts.



**Figure 20.** The CDDF from Parks et al. (2018), showing that the CNN algorithm substantially underestimates the number of DLAs in the high- $N_{\text{HI}}$  regime.

Based on the revised pipeline, we also presented a new measurement of the abundance of neutral hydrogen from  $z = 2$  to  $z = 5$  using similar calculations to Bird et al. (2017). The statistical properties

of DLAs were in good agreement with our previous results (Bird et al. 2017) and consistent with Noterdaeme et al. (2012), Prochaska & Wolfe (2009), and Crighton et al. (2015). The modifications made, including introducing a sub-DLA model, adjusting the mean flux, and penalizing complex models with Occam’s razor, remove overdetections of low column density absorbers and make more robust predictions for the properties of DLAs at  $z > 4$ . Similarly to previous work, we detect only a small increase in the CDDF for  $2 < z < 4$ , and a similarly moderate increase in the line density of DLAs and  $\Omega_{\text{DLA}}$  over this redshift range.

## ACKNOWLEDGEMENTS

The authors thank Bryan Scott and Reza Monadi for valuable comments. This work was partially supported by an [AWS Machi ne Learning Research Awards](#) allocation on EC2 and UCR’s HPCC. SB was supported by the National Science Foundation (NSF) under award number AST-1817256. RG was supported by the NSF under award number IIS-1845434.

## DATA AVAILABILITY

All the code to reproduce the data products is available in our GitHub repo: [https://github.com/rmgarnett/gp\\_dla\\_detection/tree/master/multi\\_dlas](https://github.com/rmgarnett/gp_dla_detection/tree/master/multi_dlas). The final data products are available in this Google Drive: [http://tiny.cc/multidla.catalog\\_gp\\_dr12q](http://tiny.cc/multidla.catalog_gp_dr12q), including a MAT (HDF5) catalogue and a JSON catalogue. README files are included in each folder to explain the content of the catalogues.

## REFERENCES

- Becker G. D., Hewett P. C., Worseck G., Prochaska J. X., 2013, *MNRAS*, 430, 2067  
 Bird S., Vogelsberger M., Haehnelt M., Sijacki D., Genel S., Torrey P., Springel V., Hernquist L., 2014, *MNRAS*, 445, 2313

- Bird S., Haehnelt M., Neeleman M., Genel S., Vogelsberger M., Hernquist L., 2015, *MNRAS*, 447, 1834
- Bird S., Garnett R., Ho S., 2017, *MNRAS*, 466, 2111
- Carithers W., 2012, Published internally to SDSS
- Cen R., 2012, *ApJ*, 748, 121
- Crighton N. H. M. et al., 2015, *MNRAS*, 452, 217
- Fauber L., Ho M.-F., Bird S., Shelton C. R., Garnett R., Korde I., 2020, *MNRAS*, preprint ([arXiv:2006.07343](https://arxiv.org/abs/2006.07343))
- Fernandez M., Williams S., 2010, *IEEE Trans. Aerosp. Electron. Syst.*, 46, 803
- Fumagalli M., O’Meara J. M., Prochaska J. X., Worseck G., 2013, *ApJ*, 775, 78
- Fumagalli M., O’Meara J. M., Prochaska J. X., Rafelski M., Kanekar N., 2014, *MNRAS*, 446, 3178
- Gardner J. P., Katz N., Weinberg D. H., Hernquist L., 1997, *ApJ*, 486, 42
- Garnett R., Ho S., Bird S., Schneider J., 2017, *MNRAS*, 472, 1850
- Haehnelt M. G., Steinmetz M., Rauch M., 1998, *ApJ*, 495, 647
- Kim T.-S., Bolton J. S., Viel M., Haehnelt M. G., Carswell R. F., 2007, *MNRAS*, 382, 1657
- Lee K.-G. et al., 2013, *AJ*, 145, 69
- Le Cam L., 1960, *Pac. J. Math.*, 10, 1181
- Noterdaeme p. et al., 2012, *A&A*, 547, L1
- Pâris I. et al., 2012, *A&A*, 548, A66
- Pâris I. et al., 2014, *A&A*, 563, A54
- Pâris I. et al., 2018, *A&A*, 613, A51
- Parks D., Prochaska J. X., Dong S., Cai Z., 2018, *MNRAS*, 476, 1151
- Pontzen A. et al., 2008, *MNRAS*, 390, 1349
- Prochaska J. X., Wolfe A. M., 1997, *ApJ*, 487, 73
- Prochaska J. X., Wolfe A. M., 2009, *ApJ*, 696, 1543
- Prochaska J. X., Herbert-Fort S., Wolfe A. M., 2005, *ApJ*, 635, 123
- Prochaska J. X., Worseck G., O’Meara J. M., 2009, *ApJ*, 705, L113
- Rasmussen C. E., Williams C. K. I., 2005, *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, Cambridge, Massachusetts
- Sánchez-Ramírez R. et al., 2016, *MNRAS*, 456, 4488
- Slosar A. et al., 2011, *J. Cosmol. Astropart. Phys.*, 2011, 001
- Wolfe A. M., Turnshek D. A., Smith H. E., Cohen R. D., 1986, *ApJS*, 61, 249
- Wolfe A. M., Gawiser E., Prochaska J. X., 2005, *ARA&A*, 43, 861
- Worseck G., Prochaska J. X., 2011, *ApJ*, 728, 23
- York D. G. et al., 2000, *AJ*, 120, 1579
- Zafar T., Péroux C., Popping A., Milliard B., Deharveng J.-M., Frank S., 2013, *A&A*, 556, A141

This paper has been typeset from a  $\text{\LaTeX}$  file prepared by the author.