# An Unsupervised Learning Approach for Quasar Continuum Prediction

Zechang Sun [1]   Yuan-sen Ting [2][3]   Zheng Cai [1]

## Abstract

Modeling quasar spectra is a fundamental task in astrophysics as quasars are the tell-tale sign of cosmic evolution. We introduce a novel unsupervised learning algorithm, Quasar Factor Analysis (QFA), for recovering the intrinsic quasar continua from noisy quasar spectra. QFA assumes that the Ly$\alpha$ forest can be approximated as a Gaussian process, and the continuum can be well described as a latent factor model. We show that QFA can learn, through unsupervised learning and directly from the quasar spectra, the quasar continua and Ly$\alpha$ forest simultaneously. Compared to previous methods, QFA achieves state-of-the-art performance for quasar continuum prediction robustly but without the need for predefined training continua. In addition, the generative and probabilistic nature of QFA paves the way to understanding the evolution of black holes as well as performing out-of-distribution detection and other Bayesian downstream inferences.

## 1. Introduction

Powered by the accretion of matter into the supermassive black holes in the galactic nuclei, luminous quasars have played a crucial role in astronomy as the lighthouses in the distant universe since their first discovery in the 1960s. Over the last 30 years, along with the unprecedented breakthroughs in quasar spectral observations (e.g., SDSS Survey, Lyke et al., 2020), absorption systems in quasar spectra have been widely used to probe the state and matter distribution of intergalactic medium (IGM). The intervening gas between the quasar and us creates characteristic absorption superimposed on the quasar spectra, known as the Ly$\alpha$ forest. However, to extract information from the Ly$\alpha$ forest,

one would first need to recover the intrinsic quasar continua. Any imperfection in continuum reconstruction can incur non-negligible uncertainties for any other downstream tasks, from inferring the cosmological parameters to understanding the physics of the IGM and constraining the astrophysics of reionization (Lee & Spergel, 2011; Palanque-Delabrouille et al., 2013; Chabanier et al., 2019; Montero-Camacho & Mao, 2020; Bosman et al., 2021).

Due to its central role, different methods have been proposed to determine the intrinsic quasar continua, such as power-law explorations (Fan et al., 2006), principal component analysis (PCA, Suzuki et al., 2005; Lee et al., 2012; Davies et al., 2018) and deep learning techniques including multi-layer perceptron (MLP, Ďurovčíková et al., 2020; Liu & Bordoloi, 2021) and normalizing flow (Reiman et al., 2020). However, most methods proposed thus far belong to the nature of supervised learning. These methods are based on a training set of quasar spectra with pre-determined continua. As there is no way to know the ground truth continua, such continuum "labels" are usually determined in an ad-hoc way through hand-fitting spectra with a high signal-to-noise ratio (SNR). Consequently, even though SDSS has observed more than 750,000 quasar spectra to date, typically only $0.01\% - 1\%$ spectra with high SNR were used to construct the training set in previous works.

As the high SNR quasar spectra are often a biased subset of the entire quasar population, even with adequately supervised training, previous methods often struggle to generalize to all observed quasar spectra. These limitations invariably call for an unsupervised learning method, directly learning the distribution of the quasar spectra. However, unsupervised learning does not a priori lead to continuum; as quasar spectra are the composite of the quasar continuum, the intervening Ly$\alpha$ forest, and noise, these components are entirely degenerate. Thus, breaking this degeneracy requires us to harness the power of machine "learning" while including sufficient physical prior for the "modeling."

Here we propose Quasar Factor Analysis (QFA), an unsupervised algorithm that resolves these challenges that bottleneck the field: QFA can learn directly from millions of quasar spectra without any training continuum. It provides a fully probabilistic posterior of the continuum given the quasar spectrum. Furthermore, QFA provides physically

[1]Department of Astronomy, Tsinghua University, Beijing, China [2]Research School of Astronomy & Astrophysics, Australian National University, Canberra, Australia [3]School of Computing, Australian National University, Acton, ACT 2601, Australia. Correspondence to: Zechang Sun <sunzc18@mails.tsinghua.edu.cn>.

meaningful spectrum embedding and enables us to perform out-of-distribution detection from noisy data.

## 2. Method

**Statistical quasar model:** Let S be a quasar spectrum. We model the quasar spectrum as the composite of the continuum, C, and the Ly$\alpha$ forest,

$$S = C \circ \exp(-\tau(z)) + \omega(z) + \epsilon, \tag{1}$$

where z is the redshift of the absorption systems, $\epsilon \sim \mathcal{N}(0, \Sigma_\epsilon)$ is the noise and, "$\circ$" denotes element-wise product. In this model, we specifically separate the continuum from the absorption features in the modeling. In this way, through unsupervised learning, we could learn the two components simultaneously, breaking the degeneracy between continuum and absorption.

The absorption is characterized by two terms, $\exp(-\tau(z))$ and $\omega(z) \sim \mathcal{N}(0, \Sigma_\omega(z))$. The former is commonly known as the mean transmission field and is fixed using the measurement from Becker et al. (2013). On the other hand, $\omega(z)$ is a learnable system that approximates the unaccounted random fluctuations caused by the absorbers. We further assume $\Sigma_\omega$ is a diagonal matrix and

$$\text{diag}(\Sigma_\omega) = \omega_0 \circ (1 - \exp(-\tau_0(1 + z)^\beta) + c_0)^2, \tag{2}$$

where $\omega_0, \tau_0, \beta, c_0$ are free parameters which we will optimize for.

Previous works largely inspired this particular Ly$\alpha$ model (Ho et al., 2020; du Mas des Bourboux et al., 2020) which show that modeling the absorption field as a Gaussian process with the absorber redshift dependency is a robust description of the quasar absorption. Although not explicitly shown, we assume that absorption terms only apply to wavelength bluer than the quasar Ly$\alpha$ emission as the Universe expands and redshifts the quasar continuum. The Ly$\alpha$ forest only imprints on the blue part of the quasar spectrum.

As for the continuum C, previous PCA based methods (Suzuki et al., 2005; Lee et al., 2012; Davies et al., 2018), have demonstrated that the quasar continua can be robustly captured with a linear latent model. Guided by these successes, we describe the quasar continua with a latent factor model as follows:

$$C = \mu + Fh + \Psi, \tag{3}$$

generalizing PCA to a probabilistic model while going beyond the orthonormal basis as assumed in the PCA formalism. $\mu$ is the mean vector. $\Psi \sim \mathcal{N}(0, \Sigma_\Psi)$ represents the "stochastic" term capturing variances that are not modeled by the latent model, and $F \in \mathbb{R}^{N_s \times N_h}$ is the factor loading matrix. Here $N_s$ is the dimension of spectra, $N_h$ is the dimension of hidden variable and $N_h \ll N_s$. All of these are learnable parameters that will be optimized through maximum likelihood estimation. Moreover, $h \sim \mathcal{N}(0, I)$ is the factor or hidden variable with a much lower dimension (here we chose it to be eight), for which we will infer its posterior.

Finally, let A be the diagonal matrix of $\exp(-\tau(z))$, it follows that the distribution of the quasar spectrum S can be approximated as:

$$S \sim \mathcal{N}(A\mu, A(FF^T + \Sigma_\Psi)A^T + \Sigma_\omega + \Sigma_\epsilon). \tag{4}$$

**Maximum likelihood optimization:** Given our statistical model of the quasar spectra, it then follows that for any observed set of N training quasar spectra with their measurement uncertainties $\mathcal{D} = \{S^{(i)}, \Sigma_\epsilon^{(i)}, z^{(i)}\}$, we can optimize the model parameters $\mathcal{M} = \{\mu, F, \Sigma_\Psi, \omega_0, \tau_0, \beta, c_0\}$ simultaneously through maximizing the likelihood. In particular, the log-likelihood of all the observations can be written as

$$\mathcal{L}(\mathcal{M}) = \frac{1}{N} \sum_{i=1}^{N} \log \Pr(S^{(i)} | \Sigma_\epsilon^{(i)}, z^{(i)}, \mathcal{M}), \tag{5}$$

We note that training such a non-conventional latent factor model with multiple massive moving parts cannot be solved with the classical EM method (Barber, 2012) but is made possible with the modern-day deep learning framework. We develop our codes in Pytorch. The codes are available on Github. The models are optimized via `Adam` algorithm (Kingma & Ba, 2014) to search for a local minimum $\mathcal{M}^*$.

**Regularization and tricks:** Due to the large number of model parameters ($\sim 2 \times 10^4$) and the strong degeneracy between continuum, absorption, and noise, carefully designed model regularization and optimization procedures are needed to ensure that QFA converges. First, as the variations between different continua are relatively small compared to the mean continuum, we expect the learned continuum components to have small deviations from zero. This inspires us to assign a L2 regularization to model parameters. We further avoid model singularity by setting minimum values of $\Sigma_\Psi$ and $\Sigma_\omega$ as well as early stopping when the average negative log-likelihood is smaller than zero. Finally, to enforce that the predicted continua are smooth radiation profiles, we smooth each continuum component with a running window every 20 epochs during training.

**Continuum inference:** Once the model parameters are optimized, recall that the continuum is modeled as $\mu + Fh$. Therefore, to drive the probabilistic output of the continuum, it suffices to derive the posterior distribution of the hidden variable h given the observed quasar spectrum S, or mathematically,

$$P(h|z, S, \Sigma_\epsilon, \mathcal{M}^*) \propto P(S|z, h, \Sigma_\epsilon, \mathcal{M}^*)P(h). \tag{6}$$

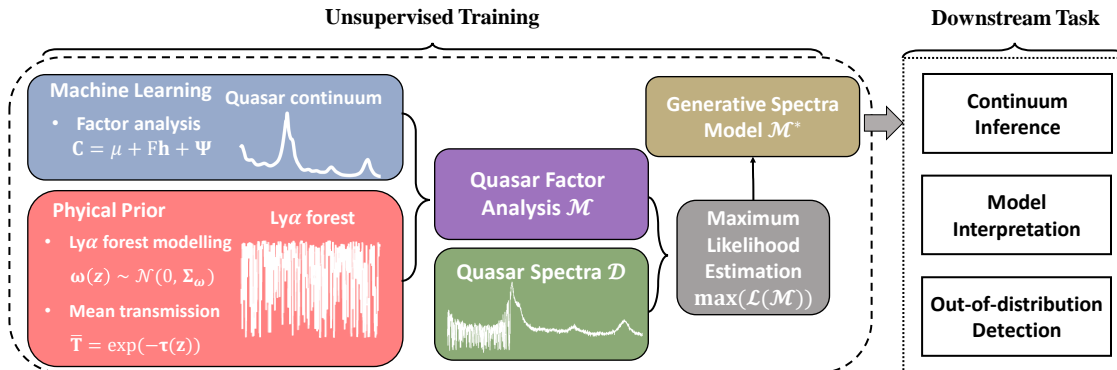Here we assume a standard normal distribution for the prior of h, or $h \sim \mathcal{N}(0, I)$.

*Figure 1.* Model architecture. Quasar spectra consist of the intrinsic quasar continua and the Ly$\alpha$ forest. We show that by assuming a Gaussian process model on the Ly$\alpha$ forest and a linear latent model for the quasar continua, it is possible to model the distribution of the quasar spectra accurately through unsupervised learning, leading to a probabilistic way of performing other downstream tasks.

## 3. Data

We evaluate model performance both on SDSS quasar spectra and mock quasar spectra. For the SDSS data set, we select $90,606$ quasar spectra with SNR greater than 2 and redshift from 2 to 3.5 from the SDSS data release 16 quasar catalogue (Lyke et al., 2020). To numerically evaluate model performance, we also generate mock quasar spectra with which the ground truth continua are known. More specifically, we measure continua of the SDSS data set with PCA components given by Pâris et al. (2011), learn the distribution of the PCA coefficients via a Gaussian Mixture Model and sample new continua from the distribution of the coefficients. As for the mock absorption field, we adopt the well-calibrated SDSS data release 11 quasar-Ly$\alpha$ mock data set (Bautista et al., 2015). We further add Gaussian random noise to mock spectra to mimic observations. Besides, to demonstrate that our model can generalize to other quasar spectra outside the high SNR data with which the PCA model was trained, we also assume $\lesssim 10\%$ random linear perturbations to the mock continua. Our final mock data set consists of $\sim 150,000$ quasar spectra.

## 4. Results

**Comparison to classical methods:** We first compare QFA's performance on the mock data set with the more widely adopted PCA algorithm (hereafter PCA, Pâris et al., 2011). By fitting all $\sim$150,000 mock spectra, we found that when fitting mock spectra are generated from the PCA basis without perturbation, both QFA and PCA reach a median $\sim 2\%$ relative error (integrated over all pixels). Our result demonstrates that even without any training continuum, QFA can perform on par with PCA, recovering continua from noisy data based entirely on unsupervised learning.

More importantly, as we perturb the mock continua, PCA

fails to generalize and incurs a relative error of $\sim 4\%$, while QFA maintains the same level of performance, at a $\sim 2\%$ relative error on the blue side. Furthermore, on the red side, QFA can reach a relative error of $\lesssim 1\%$, but PCA often incurs a non-negligible $\sim 3-4\%$ relative error.

Fig. 2 further demonstrates that QFA also performs well on actual SDSS quasar spectra. As shown in Fig. 2, PCA fails to generalize to some low SNR quasar spectra; the biased training based on high SNR leads to a nonphysical continuum. However, QFA performs well and achieves a more sensible continuum, consistent with the mock test.

**Model interpretability:** As QFA assumes a linear latent model for the continuum, it decomposes quasar continua into the key contributing components. Applying QFA to the actual SDSS DR 16 quasar spectra appears to lead to physically sensible components. As shown in Fig. 2, these individual components include Ly$\alpha$ emission, CIV emission, and the power-law slope components. Our model learns these components directly from the quasar spectra without training continua, thus not biasing us to only the high SNR and hence, low-redshift quasar data. We found that, within the redshift range covered by SDSS DR16 (z $= 2-3.5$), the individual physical components of the quasar continua have no detectable evolution with redshift, demonstrating that the quasar population has not evolved much during this period. Additionally, QFA reveals also a well-documented correlation known as the Baldwin effect (Baldwin, 1977). The Baldwin effect posits that the component corresponding to CIV emission strength negatively correlates with the quasar monochromatic luminosity, which QFA bears out.

**Outlier detection:** Finally, QFA's probabilistic nature enables also out-of-distribution detection. Specifically, for any observed spectrum, the likelihood can be evaluated with Equation 4. We applied our method to the SDSS DR16 data and performed outlier detection following Zhao et al.
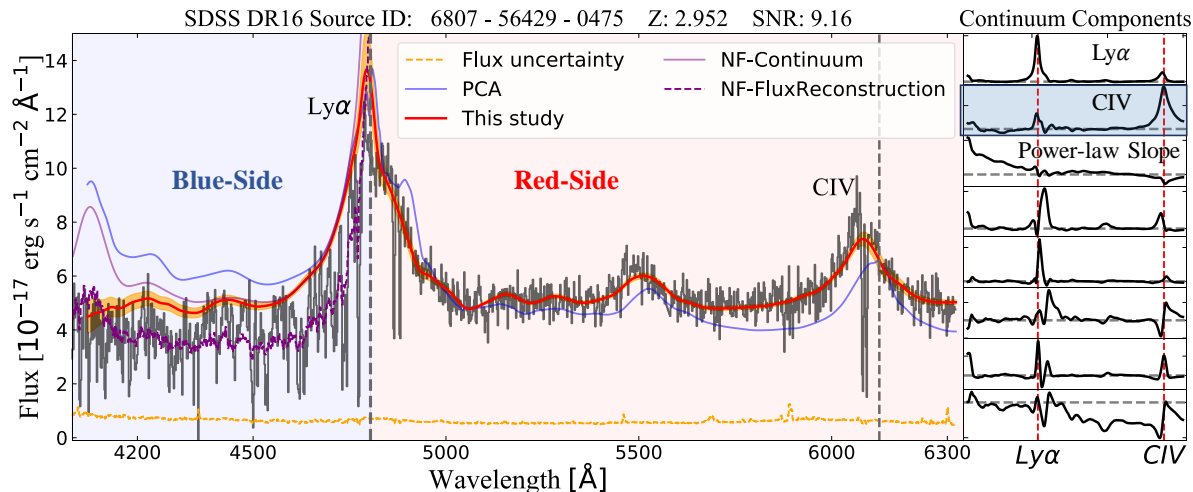
*Figure 2.* **Left:** QFA provides a probabilistic output and describes the continuum more sensibly than PCA. Shown in black is an SDSS quasar spectrum. The Orange region denotes the 95% confidence interval for QFA posterior prediction. The PCA model trained on high SNR training spectra fails to generalize to lower SNR samples, whereas QFA recovers the continuum accurately. In addition, the unsupervised deep learning method (normalizing flows or NF) fails as it suffers from mode collapse, and the reconstructed flux converges to mean spectrum. Such failure leads to problematic continuum inference with NF. **Right:** The continuum components learned by QFA are physically sensible, including one probing the Lyα emission and the other a power-law slope reminiscent of Fan et al. (2006). The CIV component shaded in blue shows the Baldwin effect (see text).
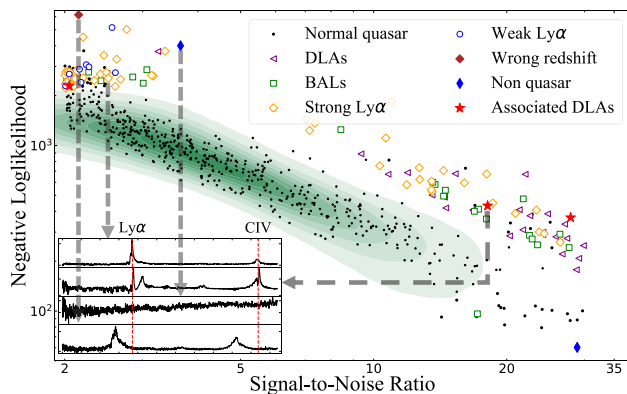


*Figure 3.* The generative and probabilistic nature of QFA allows it to perform outlier detection effectively. The background contour shows the likelihood of all the SDSS DR16 spectra. Color-coded at the low-likelihood peripheric regions shows some of the outlier objects of interest, classified by their respective classes. The bottom left corner exemplifies the spectra of some of these outliers.

(2019). As shown in Fig. 3, QFA unearthed multiple outliers, a few of them were previously unknown, including (a) undetected damped Lyα absorbers (DLAs); (b) associated damped Lyα absorbers (associated DLAs); (c) broad absorption lines (BALs); (d) Type II quasar feature – strong Lyα emission but weak continuum; (e) wrong redshift estimation; (f) misclassified non-quasar spectra. These outliers will be explored in detail in our forthcoming paper.

## 5. Discussion and Conclusion

**Deep learning?** As our method is rooted in describing the distribution of quasar spectra, a natural question arises: Can more sophisticated deep learning generative models better describe the quasar spectra distribution? Applying deep learning was the first attempt in this study. In particular, we assumed two separate normalizing flows (Rezende & Mohamed, 2015) to capture the continua and the absorption features. We pre-trained each normalizing flow to learn their respective domain and then perform variational inference on real spectra, with the hope that the two normalizing flows will domain-adapt to any synthetic gap in their respective domain. However, the lack of physical prior often leads to the reconstruction collapsing to the mean spectrum, as shown in Fig. 2. The inability to break the degeneracy subsequently leads to poor continuum reconstruction. Thus, the lack of rigidity on the physical prior and the high model complexity render unsupervised deep learning infeasible in this particular context.

**In Summary**, we demonstrate that we could recover the quasar continua robustly without any training continua. QFA harnesses the entire data set with heteroscedastic noise through unsupervised learning and can reach state-of-the-art performance for quasar continuum recovery. With more quasar spectra being collected from ongoing spectroscopic surveys such as DESI and SDSS-V, QFA might prove critical. It does not rely on continuum labels and can automatically adapt to latent variations in the data set. But perhaps

more interesting, in a time where deep learning has reigned supreme in many areas of astronomy, our work sheds light that classical machine "learning" with physical "modeling" can remain powerful, harnessing the best of both worlds.

# References

Baldwin, J. A. Luminosity Indicators in the Spectra of Quasi-Stellar Objects. *The Astrophysical Journal*, 214: 679–684, June 1977. doi: 10.1086/155294.

Barber, D. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012. doi: 10.1017/CBO9780511804779.

Bautista, J. E., Bailey, S., Font-Ribera, A., Pieri, M. M., Busca, N. G., Miralda-Escudé, J., Palanque-Delabrouille, N., Rich, J., Dawson, K., Feng, Y., Ge, J., Gontcho, S. G. A., Ho, S., Goff, J. M. L., Noterdaeme, P., Pâris, I., Rossi, G., and Schlegel, D. Mock quasar-lyman-$\alpha$ forest data-sets for the SDSS-III baryon oscillation spectroscopic survey. *Journal of Cosmology and Astroparticle Physics*, 2015(05):060–060, may 2015. doi: 10.1088/1475-7516/2015/05/060. URL https://doi.org/10.1088/1475-7516/2015/05/060.

Becker, G. D., Hewett, P. C., Worseck, G., and Prochaska, J. X. A refined measurement of the mean transmitted flux in the Ly$\alpha$ forest over $2 < z < 5$ using composite quasar spectra. *Monthly Notices of the Royal Astronomical Society*, 430(3):2067–2081, April 2013. doi: 10.1093/mnras/stt031.

Bosman, S. E. I., Ď urovčíková, D., Davies, F. B., and Eilers, A.-C. A comparison of quasar emission reconstruction techniques for $z \geq 5.0$ ly$\alpha$ and ly$\beta$ transmission. *Monthly Notices of the Royal Astronomical Society*, 503(2):2077–2096, feb 2021. doi: 10.1093/mnras/stab572. URL https://doi.org/10.1093%2Fmnras%2Fstab572.

Chabanier, S., Palanque-Delabrouille, N., Yèche, C., Le Goff, J.-M., Armengaud, E., Bautista, J., Blomqvist, M., Dawson, K., Etourneau, T., Font-Ribera, A., et al. The one-dimensional power spectrum from the sdss dr14 ly$\alpha$ forests. *Journal of Cosmology and Astroparticle Physics*, 2019(07):017, 2019.

Davies, F. B., Hennawi, J. F., Bañ ados, E., Lukić, Z., Decarli, R., Fan, X., Farina, E. P., Mazzucchelli, C., Rix, H.-W., Venemans, B. P., Walter, F., Wang, F., and Yang, J. Quantitative constraints on the reionization history from the IGM damping wing signature in two quasars at z ¿ 7. *The Astrophysical Journal*, 864(2):142, sep 2018. doi: 10.3847/1538-4357/aad6dc. URL https://doi.org/10.3847%2F1538-4357%2Faad6dc.

du Mas des Bourboux, H., Rich, J., Font-Ribera, A., de Sainte Agathe, V., Farr, J., Etourneau, T., Le Goff, J.-M., Cuceu, A., Balland, C., Bautista, J. E., and et al. The completed sdss-iv extended baryon oscillation spectroscopic survey: Baryon acoustic oscillations with ly$\alpha$ forests. *The Astrophysical Journal*, 901 (2):153, Oct 2020. ISSN 1538-4357. doi: 10.3847/1538-4357/abb085. URL http://dx.doi.org/10.3847/1538-4357/abb085.

Fan, X., Strauss, M. A., Becker, R. H., White, R. L., Gunn, J. E., Knapp, G. R., Richards, G. T., Schneider, D. P., Brinkmann, J., and Fukugita, M. Constraining the Evolution of the Ionizing Background and the Epoch of Reionization with z~6 Quasars. II. A Sample of 19 Quasars. AJ, 132(1):117–136, July 2006. doi: 10.1086/504836.

Ho, M.-F., Bird, S., and Garnett, R. Detecting multiple dlas per spectrum in sdss dr12 with gaussian processes. *Monthly Notices of the Royal Astronomical Society*, 496 (4):5436–5454, Jun 2020. ISSN 1365-2966. doi: 10.1093/mnras/staa1806. URL http://dx.doi.org/10.1093/mnras/staa1806.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Lee, K.-G. and Spergel, D. N. Threshold probability functions and thermal inhomogeneities in the ly$\alpha$ forest. *The Astrophysical Journal*, 734(1):21, may 2011. doi: 10.1088/0004-637x/734/1/21. URL https://doi.org/10.1088%2F0004-637x%2F734%2F1%2F21.

Lee, K.-G., Suzuki, N., and Spergel, D. N. Mean-flux-regulated principal component analysis continuum fitting of sloan digital sky survey ly$\alpha$ forest spectra. *The Astronomical Journal*, 143(2):51, Jan 2012. ISSN 1538-3881. doi: 10.1088/0004-6256/143/2/51. URL http://dx.doi.org/10.1088/0004-6256/143/2/51.

Liu, B. and Bordoloi, R. A deep learning approach to quasar continuum prediction. *Monthly Notices of the Royal Astronomical Society*, 502(3):3510–3532, Jan 2021. ISSN 1365-2966. doi: 10.1093/mnras/stab177. URL http://dx.doi.org/10.1093/mnras/stab177.

Lyke, B. W., Higley, A. N., McLane, J., Schurhammer, D. P., Myers, A. D., Ross, A. J., Dawson, K., Chabanier, S., Martini, P., Des Bourboux, H. D. M., et al. The sloan digital sky survey quasar catalog: Sixteenth data release. *The Astrophysical Journal Supplement Series*, 250(1):8, 2020.

Montero-Camacho, P. and Mao, Y. Ly $\alpha$ forest power spectrum as an emerging window into the epoch of reionization and cosmic dawn. *Monthly Notices of the Royal Astronomical Society*, 499(2):1640–1651, sep 2020. doi:

10.1093/mnras/staa2918. URL https://doi.org/10.1093%2Fmnras%2Fstaa2918.

Palanque-Delabrouille, N., Yèche, C., Borde, A., Le Goff, J.-M., Rossi, G., Viel, M., Aubourg, É., Bailey, S., Bautista, J., Blomqvist, M., et al. The one-dimensional ly$\alpha$ forest power spectrum from boss. *Astronomy & Astrophysics*, 559:A85, 2013.

Pâris, I., Petitjean, P., Rollinde, E., Aubourg, E., Busca, N., Charlassier, R., Delubac, T., Hamilton, J.-C, Le Goff, J.-M., Palanque-Delabrouille, N., and et al. A principal component analysis of quasar uv spectra atz 3. *Astronomy & Astrophysics*, 530:A50, May 2011. ISSN 1432-0746. doi: 10.1051/0004-6361/201016233. URL http://dx.doi.org/10.1051/0004-6361/201016233.

Reiman, D. M., Tamanas, J., Prochaska, J. X., and Ďurovčíková, D. Fully probabilistic quasar continua predictions near lyman-$\alpha$ with conditional neural spline flows, 2020. URL https://arxiv.org/abs/2006.00615.

Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows, 2015. URL https://arxiv.org/abs/1505.05770.

Suzuki, N., Tytler, D., Kirkman, D., O'Meara, J. M., and Lubin, D. Predicting QSO continua in the ly$\alpha$ forest. *The Astrophysical Journal*, 618(2):592–600, jan 2005. doi: 10.1086/426062. URL https://doi.org/10.1086/426062.

Zhao, Y., Nasrullah, Z., and Li, Z. Pyod: A python toolbox for scalable outlier detection. *Journal of Machine Learning Research*, 20(96):1–7, 2019. URL http://jmlr.org/papers/v20/19-011.html.

Ďurovčíková, D., Katz, H., Bosman, S. E. I., Davies, F. B., Devriendt, J., and Slyz, A. Reionization history constraints from neural network based predictions of high-redshift quasar continua. *Monthly Notices of the Royal Astronomical Society*, 493(3):4256–4275, Feb 2020. ISSN 1365-2966. doi: 10.1093/mnras/staa505. URL http://dx.doi.org/10.1093/mnras/staa505.