

**Ruprecht-Karls-Universität Heidelberg**  
**Institut für Informatik**  
**Lehrstuhl für Artificial Intelligence for Programming**

**Masterarbeit**

Signature of warm dark matter in the  
cosmological density fields extracted using  
Machine Learning

Name: Ander Artola  
Matrikelnummer: Matrikelnummer des Autors  
Betreuer: Name des Betreuers  
Datum der Abgabe: dd.mm.yyyy

Ich versichere, dass ich diese Bachelor-Arbeit selbstständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

Heidelberg, dd.mm.yyyy

---

## **Zusammenfassung**

Dies ist eine Zusammenfassung der Arbeit.

## **Abstract**

This is the abstract.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Cosmological preliminaries . . . . .	1
1.2	Dark matter . . . . .	3
1.2.1	Evidence for the inclusion of dark matter in the cosmological model	3
1.2.2	Manifold dark matter candidates . . . . .	3
1.2.3	Constraining mechanisms for dark matter . . . . .	3
1.3	The intergalactic medium . . . . .	3
1.3.1	The Lyman-alpha forest as a probe of the IGM . . . . .	3
<b>2</b>	<b>Simulating the Lyman-alpha forest: The Sherwood simulation suite</b>	<b>6</b>
2.1	Prelude on cosmological simulation . . . . .	6
2.2	Obtaining mock Lyman-alpha skewers from cosmological simulations . .	6
2.3	Statistical analysis of the effect of dark matter in the flux and density fields	6
2.4	Peculiar velocities and optical depth-weighted quantities . . . . .	6
<b>3</b>	<b>Deep Learning the Lyman-alpha forest</b>	<b>8</b>
3.1	Introduction and motivation for the use of Deep Learning . . . . .	8
3.2	Fundamentals of (Bayesian) Neural Networks . . . . .	14
3.2.1	Dataset generation, data augmentation and overfitting. . . . .	17
3.2.2	Deep learning architecture . . . . .	20
3.2.3	Prediction uncertainty and Bayesian models . . . . .	22
3.2.4	Hyperparameter selection . . . . .	24
3.2.5	Loss function and training . . . . .	26
3.3	Workflow implementation: Recovering IGM conditions from the Lyman-alpha forest . . . . .	29
3.4	Recovered field statistics and uncertainties . . . . .	34
3.4.1	Noisy regions and masking . . . . .	35
3.4.2	Uncertainty in the recovered statistics . . . . .	36

3.5	Model interpretability and limitations . . . . .	37
3.5.1	Saliency analysis . . . . .	37
3.5.2	Covariate shift . . . . .	39
3.5.3	Extreme covariate shift and malicious data . . . . .	41
3.5.4	Model pruning . . . . .	41
<b>4</b>	<b>Constraining Warm Dark Matter at the density level</b>	<b>44</b>
4.1	Inference pipeline: from Lyman-alpha skewers to WDM constraints . . . . .	44
4.2	Inference testing on Sherwood spectra under realistic observational conditions . . . . .	46
4.2.1	Untrained DM models . . . . .	46
4.2.2	Realistic UVES observational conditions . . . . .	48
4.3	Inference on alternative hydrodynamical codes . . . . .	50
4.3.1	The Nyx code . . . . .	51
4.3.2	Inference test on Nyx Lyman-alpha skewers . . . . .	54
4.4	WDM constraints from SQUAD DR1 observational data . . . . .	56
4.5	WDM constraints from GHOST observed spectrum . . . . .	58
4.6	Comparison of the inference pipeline against Information Maximising Neural Networks . . . . .	60
4.6.1	Information Maximising Neural Networks . . . . .	60
4.6.2	IMNN training and non-linear summaries . . . . .	61
4.6.3	Summarising a Gaussian signal . . . . .	62
4.6.4	IMNN inference results on WDM masses . . . . .	66
<b>5</b>	<b>Conclusions</b>	<b>68</b>
<b>Bibliography</b>		<b>69</b>

# 1 Introduction

## 1.1 Cosmological preliminaries

The currently accepted cosmological model describes space-time as a 4-dimensional Lorentzian manifold equipped with the Robertson-Walker metric [1]

$$ds^2 = c^2 dt^2 - a(t)^2 \left( \frac{dr^2}{1 - kr^2} + r^2 d\Omega^2 \right), \quad (1.1)$$

with  $c$  the speed of light in vacuum,  $a$  the scale factor,  $k$  a curvature parameter and  $d\Omega$  the angular volume element in spherical coordinates. The scale factor is taken to be unity at the present time. At time  $t$ , a physical (proper) distance  $l_{\text{phy}}$  is then related to a comoving distance  $l_{\text{cov}}$  by

$$l_{\text{phy}} = a(t)l_{\text{cov}}. \quad (1.2)$$

The physical distance at time  $t$  between an observer at  $r = 0$  and a point at  $r$  is then

$$l_{\text{phy}} = a(t) \int_0^r \frac{dr}{\sqrt{1 - kr^2}} = a(t)\chi(r). \quad (1.3)$$

The Robertson-Walker metric implies that for a radial luminous signal emitted at time  $t_e$  and received at time  $t_0$ , we have

$$ds^2 = 0 \implies \frac{dt_0}{a(t_0)} = \frac{dt_e}{a(t_e)}. \quad (1.4)$$

As a consequence, the received frequency is redshifted according to

$$1 + z = \frac{\lambda_0}{\lambda_e} = \frac{\nu_e}{\nu_0} = \frac{a(t_0)}{a(t_e)}, \quad (1.5)$$

where  $z$  is the redshift.

The time-dependence of physical distances in equation (1.2) implies that an object whose comoving distance  $\chi$  to an observer is constant recedes by following the Hubble

flow according to

$$v(t) = \dot{a}(t)\chi = \frac{\dot{a}}{a}a\chi = H(t)l_{\text{phy}}, \quad (1.6)$$

where  $H(t)$  is known as the Hubble factor. equation (1.6) is known a Hubble's law. At present time,  $H(t_0) = H_0$  is referred to as Hubble's constant. For historical reasons, it is common to work with the reduced Hubble constant  $h = H_0[\text{km/s/Mpc}]/100$ . Note that, according to equation (1.3), and using the Robertson-Walker metric for a radial light signal, we obtain

$$d\chi = \frac{cdt}{a} \implies \chi = \int_a^1 \frac{da}{a\dot{a}} = \int_0^z \frac{cdz}{H(z)}. \quad (1.7)$$

As a consequence, the proper line element satisfies

$$d\chi = \frac{cdz}{H(z)} = \frac{dl}{a(t)} \implies \frac{dl}{dz} = \frac{c}{(1+z)H(z)}, \quad (1.8)$$

which will be useful when integrating quantities along a line of sight. When working with such sightlines in spectroscopy, it is often advantageous to work with velocity units instead of redshifts (or proper distances). Differentiating equation (1.6) and considering a slow varying Hubble factor around a mean redshift  $\bar{z}$ , we obtain the following useful expression:

$$dv = H(\bar{z})dl = H(\bar{z})\frac{cdz}{(1+\bar{z})H(\bar{z})} = \frac{cdz}{1+\bar{z}}. \quad (1.9)$$

The evolution of the scale factor (and hence of the redshift) with time is completely determined by the energy content of the universe through Einstein's field equation, which is known as Friedmann's equation in this context

$$H^2 = H_0^2 (\Omega_M(1+z)^3 + \Omega_R(1+z)^3 + \Omega_\Lambda + \Omega_K(1+z)^2) = H_0^2 E(z)^2, \quad (1.10)$$

where the density parameters  $\Omega$  are related to the physical densities of the components according to

$$\begin{aligned} \Omega_M &= \frac{8\pi G}{3H_0^2} \rho_{M0} \\ \Omega_R &= \frac{8\pi G}{3H_0^2} \rho_{R0} \\ \Omega_\Lambda &= \frac{8\pi G}{3H_0^2} \rho_\Lambda \\ \Omega_K &= -\frac{k}{H_0^2} \end{aligned} \quad (1.11)$$

In equation (1.11),  $\rho_M$  denotes the matter density of the universe,  $\rho_R$  the radiation density, and  $\rho_\Lambda$  the dark energy component. In the following, the values used for the cosmological parameters are  $\Omega_m = 0.308$ ,  $\Omega_\Lambda = 0.692$ ,  $h = 0.678$ ,  $\Omega_b = 0.0482$ ,  $\sigma_8 = 0.829$  and  $n = 0.961$ ,  $\Omega_K \approx 0$ , as obtained from CMB measurements by the Planck Collaboration [2]. With the previous cosmological parameters, the matter and cosmological constant are equal when

$$\Omega_M(1+z)^3 = \Omega_\Lambda \implies z \approx 0.3. \quad (1.12)$$

In consequence, at the redshift of interest for this work,  $z \sim 4 - 5$ , the universe is well-described as a matter dominated universe.

## 1.2 Dark matter

### 1.2.1 Evidence for the inclusion of dark matter in the cosmological model

### 1.2.2 Manifold dark matter candidates

### 1.2.3 Constraining mechanisms for dark matter

## 1.3 The intergalactic medium

### 1.3.1 The Lyman-alpha forest as a probe of the IGM

Let us now describe how an intergalactic cloud (with no peculiar velocity) along the line of sight of a quasar affects its spectrum, allowing for a powerful probing mechanism of the IGM. Consider the situation illustrated in figure 1.1, where a QSO at redshift  $z_{\text{QSO}}$  emits photons, and consider the propagation of an emitted photon with rest-frame frequency  $\nu_e$ . Those photons are redshifted and are absorbed at  $z_\alpha$  by a neutral hydrogen absorber with local number density  $n(z_\alpha)$  producing an absorption feature in the flux at the rest-frame Lyman- $\alpha$  resonance  $\lambda_\alpha \approx 1215\text{\AA}$ . The unabsorbed photons are then redshifted and are detected by an observer at  $z = 0$  and a frequency  $\nu_o$ . The relationship between the frequencies mentioned above is then:

$$\nu_o = \frac{\nu_e}{1 + z_e} = \frac{\nu_\alpha}{1 + z_\alpha} \quad (1.13)$$

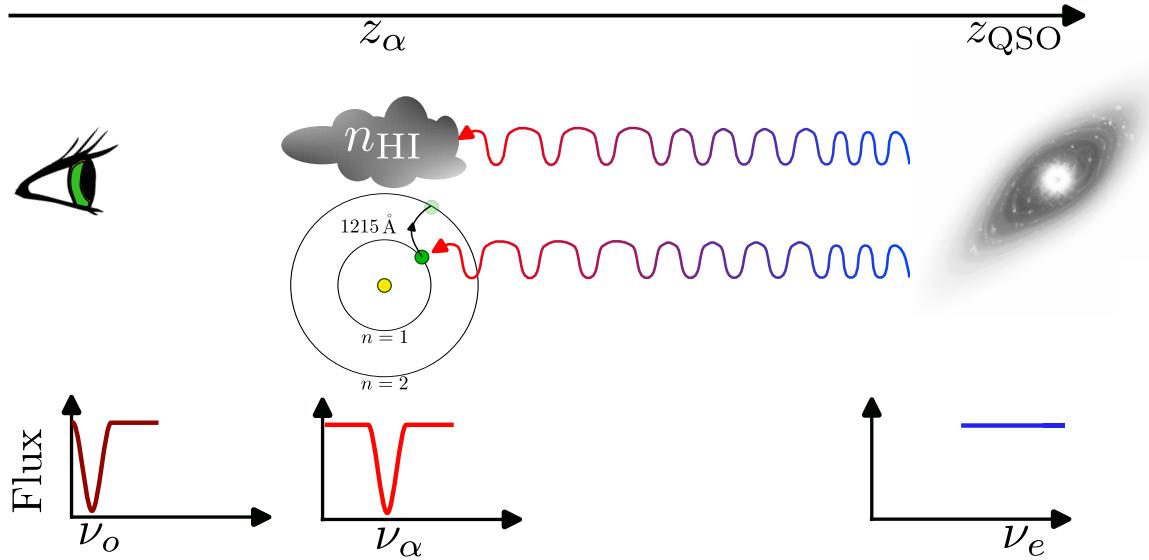


Figure 1.1: Illustration of the Lyman- $\alpha$  absorption by neutral hydrogen at  $z = z_\alpha$  in the line of sight of a QSO at  $z = z_{\text{QSO}}$ . In the observer's rest frame, the observed frequency is  $\nu_o$ . The associated frequency emitted by the QSO is  $\nu_e$ .

We are interested in studying the effect of the Lyman- $\alpha$  absorbed at  $z_\alpha$ . The observed flux attenuation at the observed frequency  $\nu_o$  is then expressed as  $\exp(-\tau_\alpha)$ , with  $\tau_\alpha$  the Lyman- $\alpha$  opacity at the observed frequency, which depends on the observer's density and the Lyman- $\alpha$  cross-section  $\sigma_\alpha(\nu)$ . Observe now that since the Lyman- $\alpha$  cross-section is strongly peaked at the resonance  $\nu_\alpha$ , but can have a non-zero width, a nearby neutral hydrogen cloud might absorb photons at a redshift different to  $z_\alpha$  that would have contributed to the observed flux at frequency  $\nu_o$ . With this consideration, we integrate over the line of sight to obtain the Lyman- $\alpha$  opacity at the observed frequency

$$\tau_\alpha(\nu_o) = \int_o^{z_{\text{QSO}}} n_{\text{HI}}(z) \sigma_\alpha[\nu_o(1+z)] dz. \quad (1.14)$$

If we now take  $\sigma_\alpha(\nu)$  to be a Dirac delta centered at the resonance  $\nu_\alpha$  and we integrate equation (1.14) by using 1.8 we obtain

$$\tau_\alpha(\nu_o) \approx \frac{c n_{\text{HI}}(z_\alpha) \sigma_\alpha}{H_0 \Omega_m^{1/2} (1+z)^{1/3}}, \quad (1.15)$$

where now  $\sigma_\alpha = 4.5 \times 10^{-18} \text{ cm}^2$  is to total Lyman- $\alpha$  cross-section. Equation 1.15 is known as the Gunn-Peterson approximation for the Lyman- $\alpha$  opacity of the IGM [3]. Equation 1.15 demonstrates that quasar spectra are a useful probe of the intergalactic neutral hydrogen density.

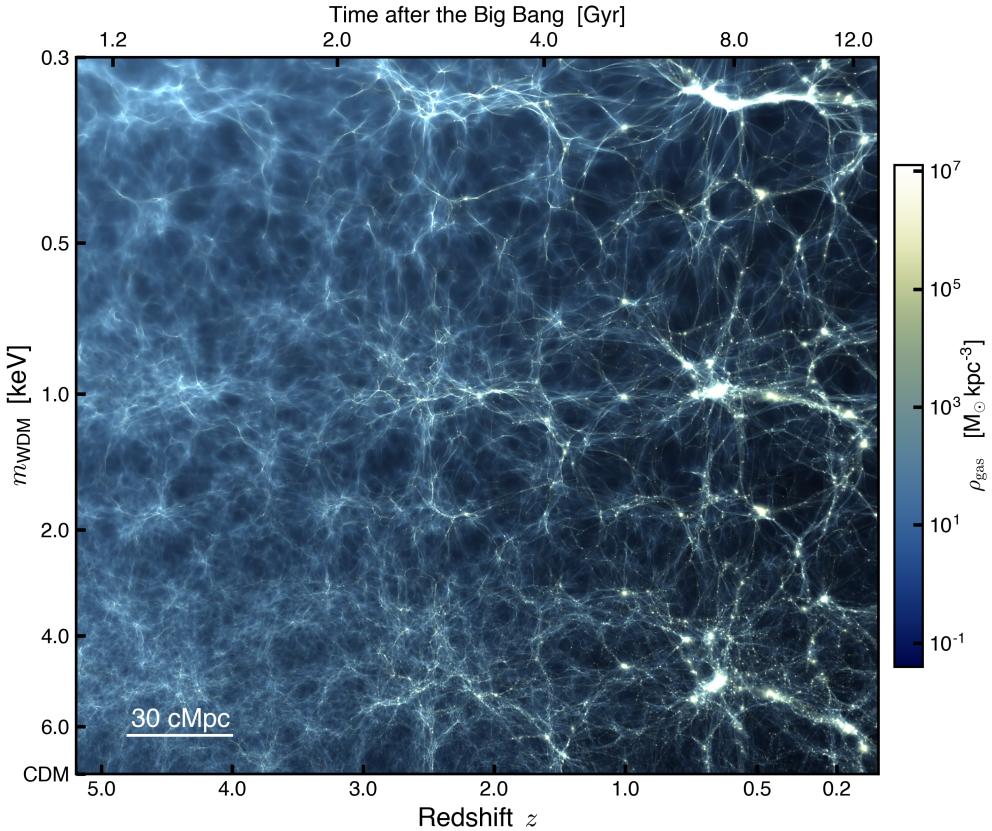


Figure 1.2: Baryonic density plot of the IGM as a function of redshift, and the WDM model mass. On the horizontal axis, the time evolution shows how gravity collapses dense regions into structures. On the vertical axis, the WDM free-streaming length suppress small-scale clustering. Extracted from [6].

The Gunn-Peterson Lyman- $\alpha$  opacity can then be used to estimate the average neutral hydrogen fraction  $x_{\text{HI}}$ . The evolution of the observed Lyman- $\alpha$  optical depth indicates that the IGM is highly ionized at  $z \lesssim 5.5$ , [4], [5].

$$3 + 1 + 4 \quad (1.16)$$

discuss villa plot and pdf and ps plot

ADD notations about overdensities, relationship with omega etc IGM typical density, temp, temp-density relationship, UVB, photoion equiv Gunn Peterson test, and reionization constraints contribution to UVB? ADD IN THE INTRO SOME ABOUT GR AND METRIC ADD SECTION 2.8.2 OF BOOK ADD BOLTON ARTICLE WITH UVB AND INFO

## **2 Simulating the Lyman-alpha forest: The Sherwood simulation suite**

### **2.1 Prelude on cosmological simulation**

### **2.2 Obtaining mock Lyman-alpha skewers from cosmological simulations**

### **2.3 Statistical analysis of the effect of dark matter in the flux and density fields**

DISCUSS plots with color bar of pdf and ps of Density discuss the boeara discrepancy plot the PS of WDM vs thermal models

### **2.4 Peculiar velocities and optical depth-weighted quantities**

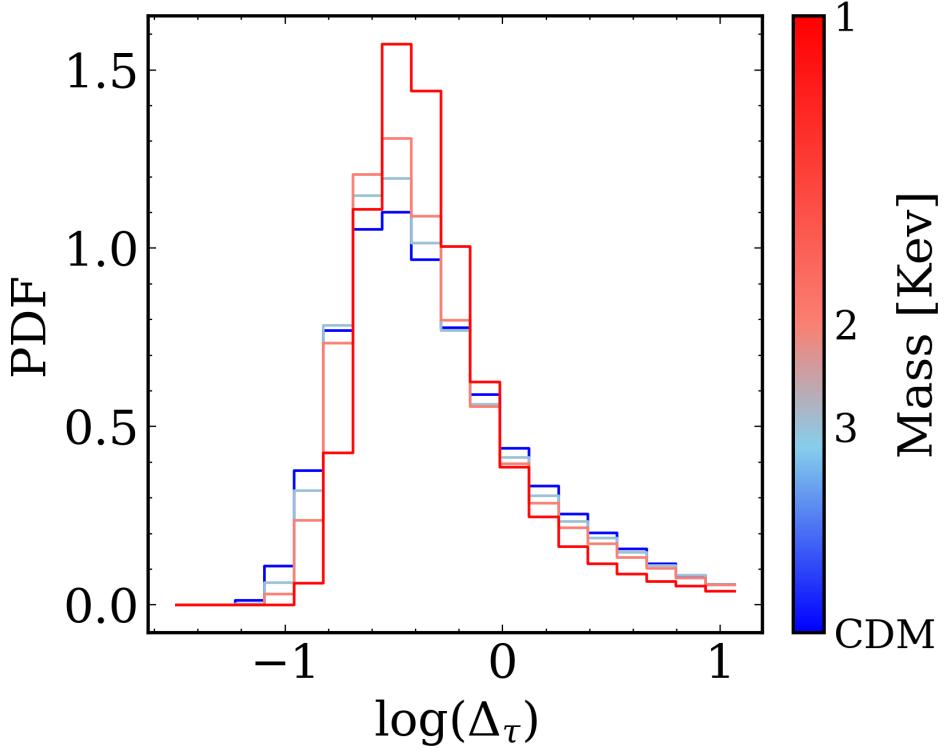


Figure 2.1: The  $\Delta_\tau$  probability distribution function (PDF) for different WDM models in the SHERWOOD suite.

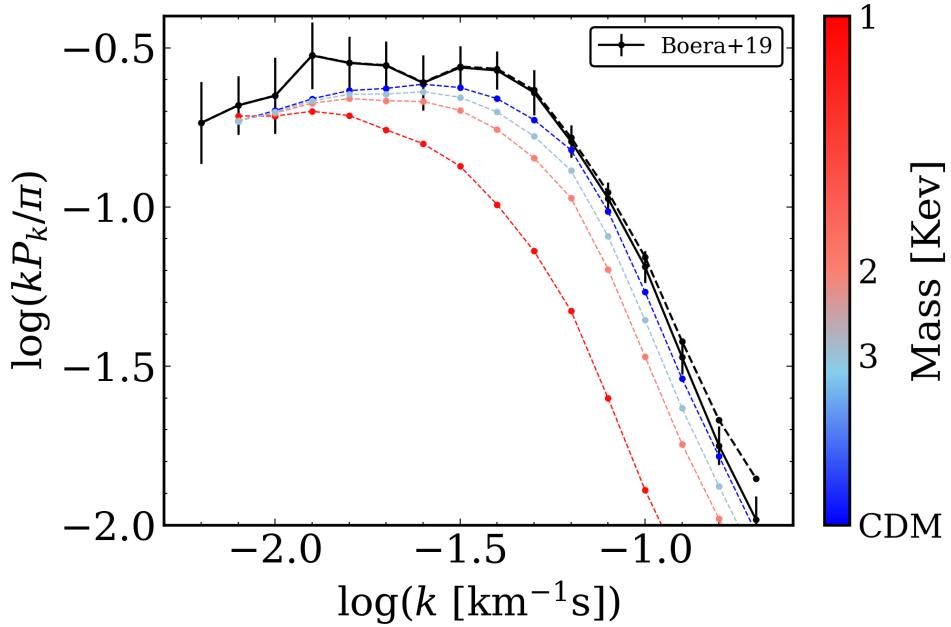


Figure 2.2: The Lyman- $\alpha$  flux power spectrum for different WDM models in the SHERWOOD simulation suite. For reference, we also plot the observed PS by [7].

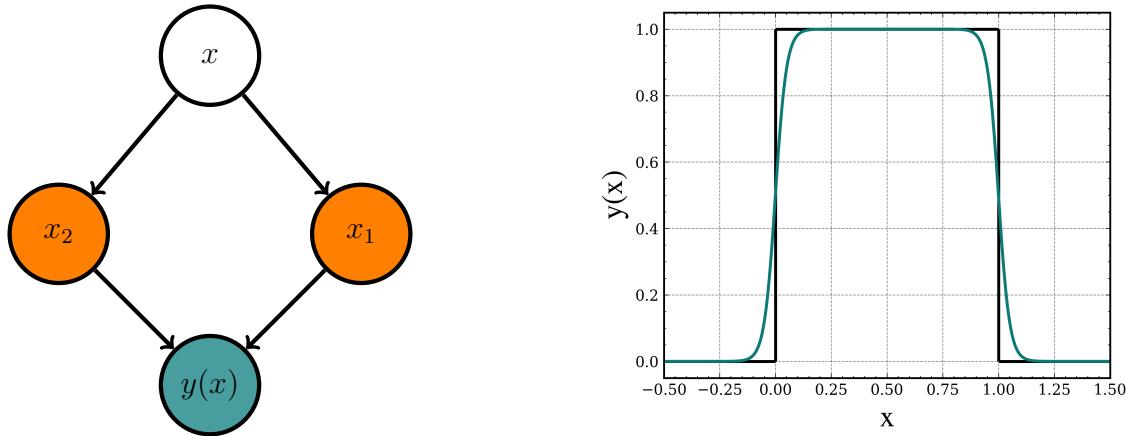
# 3 Deep Learning the Lyman-alpha forest

## 3.1 Introduction and motivation for the use of Deep Learning

Fundamentally, a neural network is a directed and acyclic computational graph [8]. It represents a set of operations that transform input data into an output. In the simplest case of a fully connected network, the nodes of this graph are represented by neurons. Neurons are organized on successive layers, making the information flow from a layer to the following one. In the most basic scenario, the neurons in a given layer are linearly connected to the neurons in the previous layer. Each neuron then adds a bias to the result of the computation and applies a non-linear function to the result, which determines the activation state of the neuron. Consider the graph shown in figure 3.1a, where the input neuron has a value  $x$  and the output neuron has a value  $y(x)$ . The intermediate layers have values

$$\begin{cases} x_1 = \sigma(\alpha_1 x + \beta_1) \\ x_2 = \sigma(\alpha_2 x + \beta_2) \end{cases}, \quad (3.1)$$

where  $\alpha_i$  are the linear weights,  $\beta_i$  the biases, and  $\sigma$  is a non-linear activation function. Typical choices include tanh or ReLU (Rectifier Linear Unit) given by  $\text{ReLU}(x) = \max(0, x)$ . Graphs such as the one shown in figure 3.1a are known as *fully connected layers*. Much more complex architectures have of course been investigated. Depending on the specific problem and dataset, we can incorporate a priori knowledge of the problem in the design of the network. For instance, in Computer Vision, the use of *convolutional layers* especially target at identifying key features in images has proven to be extremely successful [9]. In the analysis of time series, Long short-term memory (LSTM) neurons allow the network to “remember” information from previous inputs [10].



(a) Graph for a simple MLP with a hidden layer (in orange) and an output neuron (in cyan).

(b) Unit impulse (in black) and the output of the MLP shown in cyan approximating the impulse.

Figure 3.1: A simple multilayer perceptron (MLP) with a single hidden layer and a tanh activation function is able to approximate a unit pulse function. From left to right and top to bottom, the three biases are  $\{0, 20, 0\}$  and the four weights are  $\{20, -20, 1/2, 1/2\}$ .

From the theoretical standpoint, neural networks are universal approximations (for sufficiently well-behaved functions), which makes them especially appealing in the modelling of complex systems. A concrete result is as follows [11]:

**Theorem 1** (Universal approximation theorem). *If  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a Lebesgue  $p$ -integrable function and  $\varepsilon > 0$ , then there exists a fully connected ReLU network  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that*

$$\int_{\mathbb{R}^n} \|f(x) - F(x)\|^p dx < \varepsilon.$$

Theorem 1 can be expanded to include tight bounds on the depth or width of the network, which then depend on  $n$  and  $m$ . Note that all the complexity of a fully connected neural network is generated by the non-linear activation function. With a linear activation function, a fully connected network would be an affine transformation, which cannot approximate arbitrary non-linear functions.

This theoretical result can be easily visualized by understanding how a simple fully connected network (a multilayer perceptron, MLP) can approximate a unit impulse function. It is then enough to recall that a linear combination of step functions can approximate any integrable function. In figure 3.1a we show a MLP consisting of an input neuron, and output neuron and a hidden layer with two neurons with a tanh activation function. From left to right and top to bottom, we set the three biases as

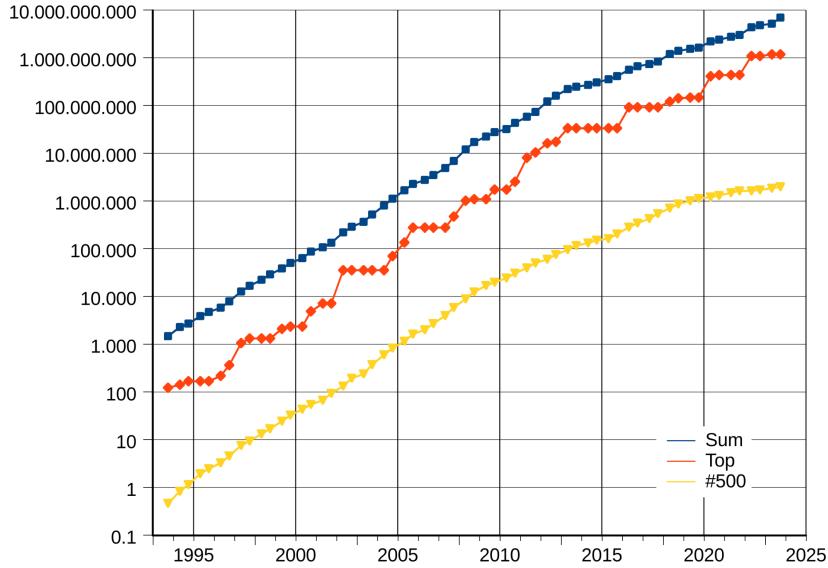


Figure 3.2: Evolution of the largest supercomputers in the TOP500 list in recent years (x-axis). The y-axis shows the peak performance in GFLOPS for the first-ranked computer (red), the last-ranked computer (yellow), and the cumulative power for the top-500 computers (blue). Source: Wikipedia.

$\{0, 20, 0\}$  and the four weights as  $\{20, -20, 1/2, 1/2\}$ . The resulting MLP then produces the approximation shown in figure 3.1b of the unit impulse over  $[0, 1]$ . By expanding the number of neurons in the hidden layer we could approximate impulse functions of arbitrary width and height and centered at every arbitrary value. The high degree of expressivity of neural networks makes them particularly suited to parametrically approximate complex non-linear relationships between variables that would otherwise be intractable.

Theorem 1 establishes a theoretical aspect of neural networks supporting their utility, but it does not provide a method (nor does it guarantee the existence of any) to find networks that approximate a particular function of interest. The process of fitting a statistical model (such as a neural network) to a dataset of interest is known as *supervised learning*, an constitutes one of the main aspects of this work, as we will explain in the rest of this section. In recent years, the main points have allowed the rapid expansion of machine learning, marking a new era in the use of large parametric networks for real-world problems:

- **Hardware development.** The rapid (quasi-exponential) growth in computer power in recent years, as exemplified by Moore’s law [12] means that extremely

deep and complex models can now be trained and effectively used. As a reference, the language models LLAMA by Meta can have over 65 billion parameters [13]. To handle this amount of data, High Performance Computing (HPC) infrastructures (such as supercomputers) are needed. Figure 3.2 shows the rapid evolution for the top supercomputers in the world, as ranked by the TOP500 list<sup>1</sup>. In 1993, the FUJITSU NUMERICAL WIND TUNNEL in Japan topped the list with 124 GFLOPS. In 2023, the list is topped by FRONTIER in the United States, with over 1000 PFLOPS, which corresponds to an 8000-fold improvement. Progress in hardware technology has also been remarkable. For instance, Graphics Processing Unit (GPU) are now common when training machine learning models. As a last example, in May 2017, Google introduced an architecture known as Tensor Processing Unit (TPU) especially designed to accelerate neural network operations [14]. HPC also enables the generation of synthetic datasets from simulations, that can then be used as training datasets. We will leverage this idea in the rest of this work.

- **Computationally-efficient algorithms.** Together with hardware development, we have seen a rapid evolution of computational and mathematical algorithms in the field of statistical learning that have enabled the efficient utilization of HPC resources. The epitome of such algorithm might be *backpropagation* [15] which is the most common algorithm used in the training phase to update the weights when deploying a neural network. Together with backpropagation, a growing set of optimizers for deep learning problems have been developed. A popular choice is ADAM, which was originally published in 2017 [16].
- **Availability of large datasets.** The advent of next-generation instruments and data-collections systems provides the scientific community with increasingly large datasets. Prominent examples in the field of astronomy include Gaia [17], whose third data release (DR3) includes 10TB of data for 1.46 billion sources<sup>2</sup>, or the Euclid telescope, expected to deliver 850 GB of compressed data per day<sup>3</sup>. At last, in figure 3.3 we show the evolution in the number of SNe discoveries per year. The figure has been obtained from [18].

---

<sup>1</sup><https://top500.org>

<sup>2</sup><https://www.cosmos.esa.int/web/gaia/dr3>

<sup>3</sup><https://sci.esa.int/web/euclid/-/46661-mission-operations>

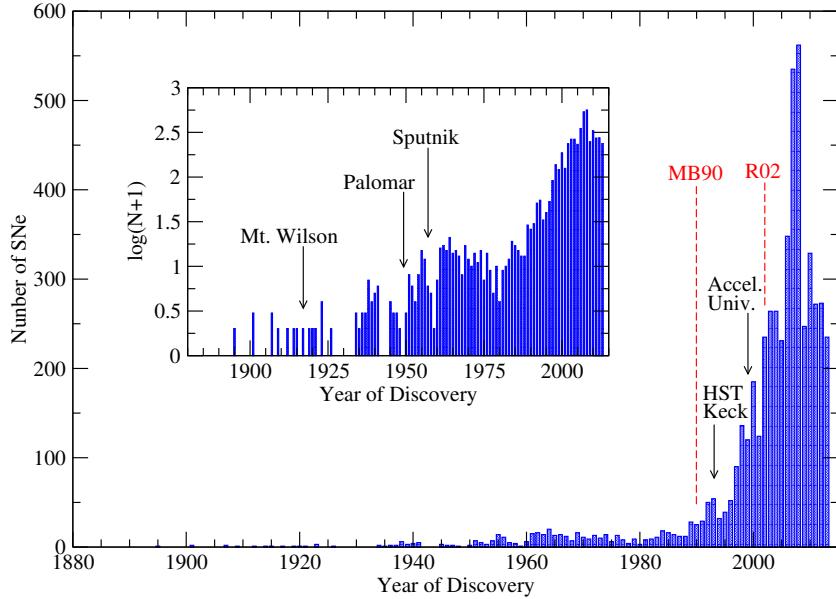


Figure 3.3: Histogram showing the number of SNe discovered each year as given by the Asiago Supernova Catalogue.

Let us demonstrate the relevance, pertinence, and range of applicability of machine learning methods in astronomy and cosmology by examining recent endeavours documented in the literature. While not exhaustive, this examination provides a broad overview of previous efforts in the field, and illustrates some recent successes of machine learning-oriented approaches in data-driven problems.

The abundance of data from surveys covering large regions of the sky aimed at targeting QSOs makes supervised statistical techniques a particularly attractive data analysis technique. For reference, the main quasar sample in the data release number 16 for the Extended Baryon Oscillation Spectroscopic Survey (eBOSS) contains 434820 targets with redshifts in the range  $0.8 < z < 2.2$  [19]. The prediction of the intrinsic Lyman- $\alpha$  emission line from high redshift QSOs is a non-trivial problem that has important implication in the study of IGM damping wings and the reionization of the later [20]. Reconstruction techniques often rely on the correlation of the Lyman- $\alpha$  peak with other observable lines and information redward of the Lyman- $\alpha$  line. Machine learning approaches are then suitable to connect the unattenuated information redward of the Lyman- $\alpha$  peak with its intrinsic profile. In *Blind QSO reconstruction challenge: Exploring methods to reconstruct the Ly $\alpha$  emission line of QSOs* [21], the authors perform an in-depth comparison of different state-of-the-art techniques based on statistical learning. The authors blindly evaluate their performance on two QSOs samples randomly

extracted from X-Shooter and BOSS with  $3.5 < z < 4.5$ , in such a way that selecting samples already used in the training data-set was avoided. The various techniques range from principal component analysis (PCA) approaches, such as [22], to deep learning networks [23]. The authors conclude that the better performing pipelines consistently rely on machine learning approaches. The authors caution against overreliance on machine learning techniques due to their potential lack of statistical uncertainty, which is one of the main aspect that we will develop in this work.

Parameter inference on WDM using deep learning has already been tentatively explored in the recent literature. The paper *Inferring Warm Dark Matter Masses with Deep Learning* [24], which is especially relevant for this work, demonstrates that neural networks can be used to recovered WDM parameters from observed field density images. The authors present a suite of 1500 cosmological N-body simulations with varied WDM mass in the range 2.5 to 30 KeV. Field density images of size  $25h^{-1}\text{Mpc}$ , with varied image resolution, simulation resolution and redshift (in the range  $0 \leq z \leq 5$ ) are extracted from the simulation runs. The images are augmented usual standard techniques (such as image rotation) and then incorporated into training datasets. The authors use a Convolutional Neural Network (CNN) trained to directly predict WDM masses based on an input density field image. Their fiducial convolutional network trained with the set of highest resolution image is able to accurately recover WDM masses with an accuracy of  $\pm 1$  KeV for models up to 10 KeV. After this threshold value, the network is no longer able to recover the true mass as predicts an approximately constant mass, see figure 3.4. Note that in the architecture presented by the authors, the network is also trained to predict an uncertainty estimate. Another interesting insight offered by this paper refers by the capacity of neural network to make use of the full information contained on a field distribution (in this case, the density field). By training a network on a full density field image the models learns the relevant properties of the field that can lead to an accurate parameter prediction. The authors train another model only on summary statistic, in particular, on the density power spectrum, and compare its performance with their fiducial model trained on the full images. The model trained only on the power spectra shows a significantly degraded performance, with higher uncertainties and accurate prediction only up to 5.5 KeV. This illustrate how deep learning techniques can help harvesting the full information present in a field (or generally, in a complex input tensor).

Another recent use of deep learning to analyse IGM data can be found in the paper  *$\text{Ly}\alpha\text{NNA: A Deep Learning Field-level Inference Machine for the Lyman-}\alpha\text{ Forest}$*  [25],

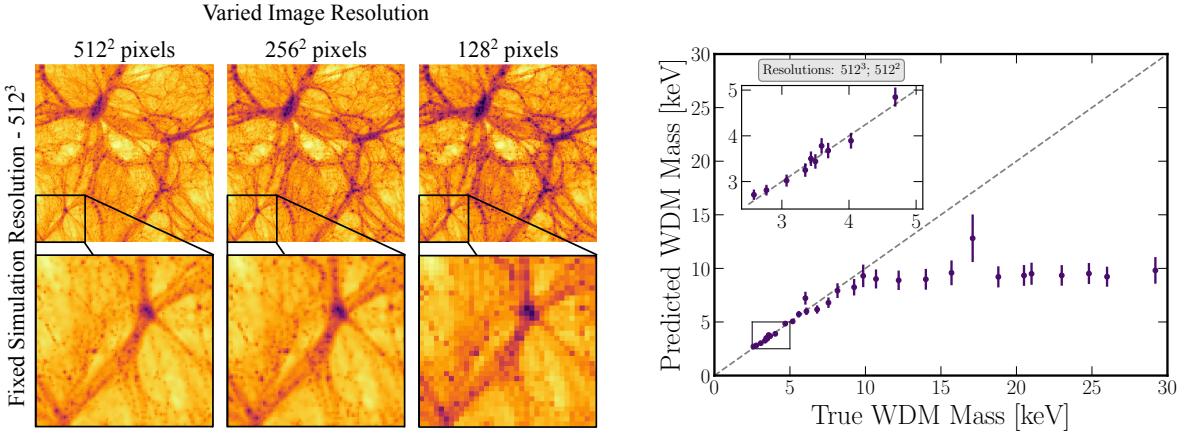


Figure 3.4: Figure extracted from [24]. The left panel shows a sample image density field used in the training data by the authors, with varied image resolution. The right panel shows a sample a predicted WDM masses versus the true WDM mass of the simulation for their fiducial neural network, which can accurately recover the WDM model within a 1 KeV accuracy up to 10 KeV.

where authors harvest the field level potential of residual convolutional networks to perform inference on the thermal parameters of the IGM, namely  $T_0$ , the temperature at mean density, and  $\gamma$ , the slope of the temperature-density relation. Their model is trained on simulation boxes with side-length 120 Mpc from which  $10^5$  sightlines are extracted and processed to produce mock Lyman- $\alpha$  spectra. The simulation boxes are run with different thermal parameters, by sampling 121  $(T_0, \gamma)$  combinations in the parameter space. The network is trained on 24000 labeled spectra from the mix of thermal models, and the architecture is designed to predict a mean value for the thermal parameters as well as an estimate for the parameter covariance matrix. Figure 3.5 shows the scatter in the point predictions for  $(T_0, \gamma)$  for a set of 4000 unseen test spectra. The true parameter values, shown as dashed lines, are recovered by the average of the point predictions, shown as the dark green cross. The authors also perform a comparison with an inference pipeline based on the traditional transmitted flux power spectrum, and find that the posterior constraint using the machine learning field-level approach are 5.65 times tighter.

## 3.2 Fundamentals of (Bayesian) Neural Networks

A very general regression problem in statistical learning [26] arises when we observe a quantity  $Y$  that is assumed to depend on an independent variable  $X$  through a relation

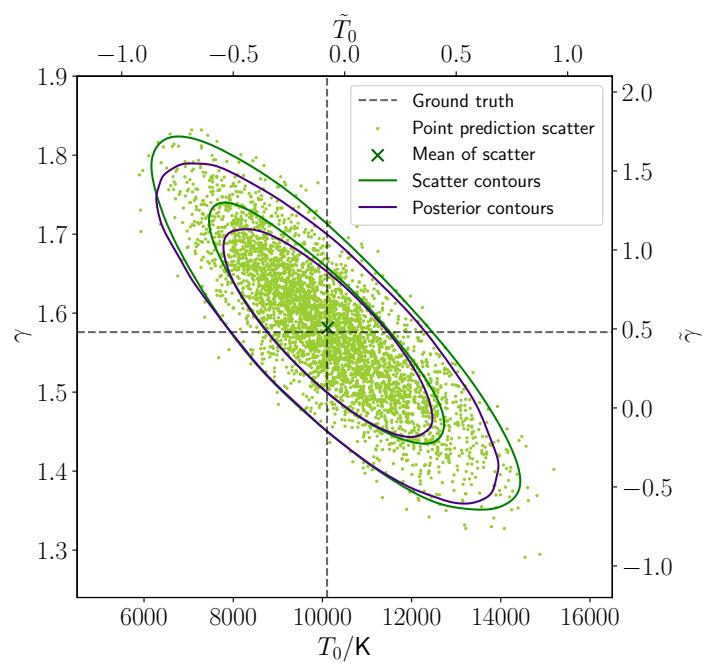


Figure 3.5: Figure extract from [25] showing the performance of the neural network recovering thermal parameters on a set of unseen skewers. The true parameter values, shown as dashed lines, are recovered by the average of the point predictions, shown as the dark green cross.

of the form

$$Y = f(X) + \mathcal{N}(0, \sigma), \quad (3.2)$$

where  $f$  is a function defining the relation  $Y = Y(X)$  and  $\mathcal{N}(0, \sigma)$  is an error term modeling a zero-mean Gaussian noise. The problem then consists of estimating  $f$  given a sample of observations  $\{X_i, Y_i\}_i$  to obtain a functional form  $\hat{f}$  that can then be evaluated to obtain predictions. Following the examples in section 3.1,  $X$  might be the flux spectra redward of the Lyman- $\alpha$  peak, and  $Y$  might then be the intrinsic shape of the Lyman- $\alpha$  peak in the quasar spectra. The noise then is produced by different sources, for instance, associated with the instrumental devices. We could then use the obtained  $\hat{f}$  to predict to intrinsic Lyman- $\alpha$  peak of a QSO,  $\hat{f}(X)$ , when we only have information about the spectrum redward of it,  $X$ . As it is expected, the quality of our prediction depends on how accurate is the approximation  $\hat{f}$  with respect to  $f$ , but also on how noisy our data is. In fact, we see that, if the true quantity associated with  $X$  is  $Y$ :

$$(Y - \hat{f}(X))^2 = (f(x) + \mathcal{N}(0, \sigma) - \hat{f}(X))^2.$$

Taking expected values and noting that the only stochastic component here is the noise, we obtain

$$\mathbb{E}[(Y - \hat{f}(X))^2] = (f(X) - \hat{f}(X))^2 + \sigma^2. \quad (3.3)$$

In equation (3.3), the first term of the right-hand side represents the error produced by approximating  $f$  by  $\hat{f}$ , while the second term is a theoretical limit imposed by the noise properties.

Parametric methods represent a powerful statistical learning tool to approximate the target relation  $f$  between variables. Methods such as linear regression, which approximates  $Y = f(X)$  using a linear functional form on the parameter values, offer limited flexibility in terms of approximating complex relations but allow for a high degree of interpretability of the model. On the other end of the spectrum, deep learning models offer a large degree of expressivity (see 3.1), but interpretation of individual parameters is often not possible. Increasing the number of parameters can also cause the model to learn from spurious structures in the finite data sample, or the learn from possible noise correlation. This process, known as *over-fitting*, can significantly impact the predictive performance of the model when predicting on unseen data. Multiple techniques have been explored to mitigate over-fitting [27], we will adopt some of them in this work and explain them in the following sections.

In this section we will discuss the basic workflow involving a deep learning model in a general and abstract scenario. We will restrict ourselves to the topics that are relevant for this work, and later explain in detail how we implement this workflow for our problem. The general workflow when working with a deep learning statistical model under supervised learning includes the following phases:

1. Collecting/data and processing it to generate a training dataset. This requires specific domain expertise to assess the quality and representability of the data.
2. Designing a deep learning architecture, i.e., a computational graph that depends on a set of parameters and that generates target outputs from the input data.
3. Using the available data to train the model. This is done by optimizing the model parameters by minimizing a selected loss function.
4. Assessing the performance of the network a validation dataset, previously unseen.
5. Using the model on real unseen data to make predictions.

### 3.2.1 Dataset generation, data augmentation and overfitting.

In a general real-world setting only a finite dataset is available to us, for instance, obtained from observational procedures. That is, we have a set of pairs  $\mathcal{D} = \{X_i, Y_i\}_{i=1}^{i=N}$  of  $N$  observations from an input-output pairs, distributed according to some distribution  $F(X, Y)$ :

$$(X_i, Y_i) \sim F, i = 1, \dots, N. \quad (3.4)$$

Of course in general, we don't have access to the distribution of this population (i.e., the function  $f$  describing the data relation  $Y = f(X)$ ), otherwise we would already have a perfectly accurate model, and we would not to invoke any statistical learning tool. The challenge is then to use our sample  $\mathcal{D}$  to infer properties of the population. Since we don't have access to the generating function  $f$ , we cannot assess the general performance of the model on arbitrary input data. For this purpose, the training dataset  $\mathcal{D}$  is typically split into two disjoint subsets  $\mathcal{D} = \mathcal{D}_T \sqcup \mathcal{D}_V$ , where  $\mathcal{D}_T$  represents the subset of the data used for training, and  $\mathcal{D}_V$  represents the subset of the data used to validate the model and assess its performance. The central idea behind this split is to be able to validate the model on data that has not been used in the training process,

allowing for an estimation of the model’s generalization to the population given by:

$$\varepsilon = \mathbb{E}_F[\mathcal{L}(Y, \hat{f}(X|\mathcal{D}_T))], \quad (3.5)$$

where  $\mathcal{L}$  is a loss function. The exact form in which the  $\mathcal{D}$  is split needs chosen a priori, with the aim of having a representative sample of the population. A common easy-to-implement strategy is to randomly select the samples in  $\mathcal{D}$  that will be part of  $\mathcal{D}_V$  and distributing them according to a ratio, that is commonly taken to in the range 60 – 80%, meaning that a larger percentage of  $\mathcal{D}$  is dedicated to training. In this case, an unbiased estimator of  $\varepsilon$  is

$$\hat{\varepsilon} = \frac{1}{|\mathcal{D}_V|} \sum_{i=1}^{|\mathcal{D}_V|} \mathcal{L}(Y_i, \hat{f}(X_i)), \quad (3.6)$$

where the sum runs over the validation dataset. More complex strategies that take into account the topology of the data exist and allow for an optimal training-validation split [28]. The validation subset can also be sequentially cycled through all the available samples in a series of methods called *cross-validation* [29]. Cross-validation uses all available data in a validation stage by retraining a model on different disjoint splits of the data. This allows for a more precise evaluation of the model.

If the training subset  $\mathcal{D}_T$  is not representative of the population, then the global error in equation (3.5) will be large, and the model will not be able to learn the general properties of the data. This can cause the model to overfit and learn from spurious correlation in the data and noise, and can be encouraged by having an excessive number of free parameters. Overfitting is not an intrinsic characteristic of deep-learning models, can also be seen for instance in simple polynomial regression when the number of “independent” points exceeds the order of the fitting polynomial. Figure 3.6 illustrates the simple example of polynomial regression overfitting the training dataset as the order of the polynomial increases. Note how, by construction, the fit becomes increasingly accurate on the training data, but loses generalization power on unseen points.

The training data  $\mathcal{D}_T$  can be artificially extended using *data augmentation* techniques [30] to generate new training points from existing points. Data augmentation aims at presenting the model with data that maintains intrinsic properties of the original data but varies other features that should not affect the prediction. This approach can be understood to be inspired by symmetry considerations<sup>4</sup> and is particularly useful in computer vision tasks. Some common data augmentation techniques for image processing

---

<sup>4</sup>A strawberry is a strawberry regardless of its orientation

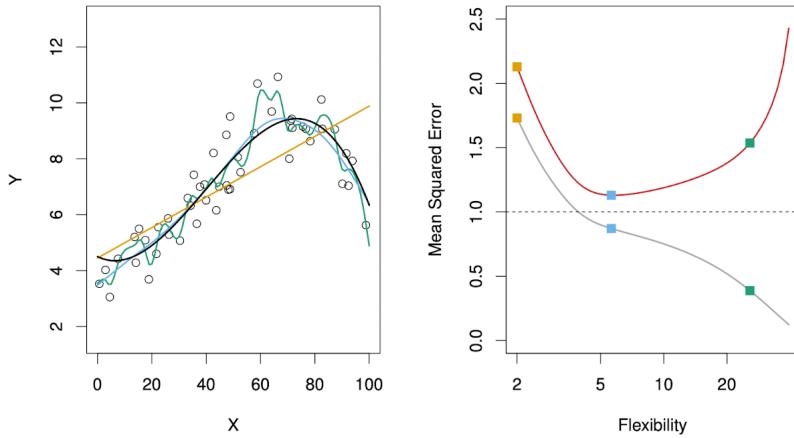


Figure 3.6: The left panel shows the data points (with added noise) generated from a function  $f$  in black. Three generative polynomial models are fitted to the data: linear regression in orange, and two smoothing splines in red and blue. The right panel shows the fitting error as a function of the polynomial degree, i.e., the flexibility of the model. The gray curve shows the error on the training dataset, which is monotonically decreasing. The red curve shows the error on the validation dataset, which initially decreases but then grows as the model overfits the training data. Figure extracted from [26].

include:

- Basic geometric transformations such as rotation, translation or flipping an image (or field) are efficiently implemented. Cropping can also be easily implemented, but changes the dimensionality of the data.
- Basic color space manipulation in colored images, including isolating a single color channel, or modifying the brightness of an image.
- Noise injection during training is particularly useful in making neural networks robust against noisy and corrupted data. Note that this can be relevant if the data we are expecting to use the trained model on has noisy (for example instrumental data from a spectrograph). The noise can be randomly drawn from independent distributions, or drawn from a noise model if we have additional insights on how to model it.
- Convolutional operations with kernels, such as the Gaussian kernel to apply a limited spatial resolution (blurring), or the Sobel kernel for related to edge detection.

### 3.2.2 Deep learning architecture

The network architecture defines the parametrization of the function that approximates the true underlying relation between input and output data points. In this section we briefly explore some basic building blocks used when constructing a deep learning network. A neural network (NN)  $\phi$  is a computational graph with an input tensor  $X$  and an output tensor  $Y$  and parametrized by a set of values  $\omega$  that define a functional model

$$Y = \phi(X|\omega) \equiv \phi_\omega(X). \quad (3.7)$$

For simple linear topologies a NN can be specified given a set of connected *layers* that carry out the elementary operations. These layers can be grouped in *blocks* to form subnetworks inside a NN. Three of the simplest layer architectures are *feedforward* layers (also known as *dense* layers), convolutional layers, and residual layers. Since they will be of use in our machine learning workflow, we will discuss their exact structure.

- In the simplest feedforward network, such as the one illustrated in figure 3.1a, all layers are linearly connected with an additional activation function, leading to a model of the form

$$\begin{aligned} X &= \mathbf{l}_0, \\ \mathbf{l}_i &= s_i(\mathbf{W}_i \mathbf{l}_{i-1} + \mathbf{b}_i) \quad \forall i \in [1, N], \\ Y &= \mathbf{l}_N, \end{aligned} \quad (3.8)$$

where  $\mathbf{l}_i$  are the layers,  $\mathbf{W}_i$  the weights,  $\mathbf{b}_i$  the biases,  $s_i$  the activation functions and  $N$  the number of layers.

- Convolutional layers are especially useful in signal processing and computer vision problem. They incorporate multiple convolution operations on the input array. The parameters of the layers define the exact way the convolution is performed. Consider for illustration purposes the case of an input 2D array of size  $(N, N)$ . The convolutional layer is internally parametrized by a kernel  $\omega$  of shape  $(m, m)$  such that the input  $X$  to the layer is transformed as

$$Y_{i,j} = \sum_{a=0}^{m-1} \sum_{b=0}^{m-1} \omega_{ab} X_{i+a, j+b}. \quad (3.9)$$

The output size will be  $(N-m+1, N-m+1)$ . Lastly, a user-specified non-linearity is applied. The implementation of convolutional layers in machine learning APIs

is usually done by specifying a kernel size and a set of auxiliary parameters. For instance, a single layer can have multiple independent kernels, can apply padding to the input array before convolving, etc.

- Residual layers implement a skipping connection within another layer, typically a convolutional one. A skipping connection adds the input tensor to the output tensor:

$$Y = T(X) + X, \quad (3.10)$$

where  $T$  is other neural layer. A relevant property, both from the theoretical and practical viewpoints, is that the derivative of the output tensor  $Y$  with respect  $X$  to  $X$  always include an identity term that tends to prevent it from being zero. This allows the information to flow easily in deep architectures, speeding up training and limiting overfitting [31].

- Batch-normalization layers [32] address the training difficulty encountered when the input tensors vary significantly from a sample to another. These differences cause different network weights to be adjusted simultaneously in opposite directions, which ultimately impairs training. Batch-normalization layers simply aim at normalization each tensor input feature when we feed a batch of samples to a network. A batch-normalization layer has 2 internal parameters  $\gamma, \beta$  for each input features, that are applied as a linear transformation after normalization. If we consider an input feature  $x$  (that is, a component of  $X$ ), and a batch of samples  $\{x_1, \dots, x_m\}$ , then a batch-normalization layer is implemented as:

$$\begin{aligned} \mu &= \frac{1}{m} \sum_{i=1}^m x_i \\ \sigma^2 &= \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2 \\ \hat{x}_i &= \frac{x_i - \mu}{\sqrt{\sigma^2}} \end{aligned} \quad (3.11)$$

$$\text{BN}_{\gamma, \beta}(x_i) \equiv \gamma \hat{x}_i + \beta$$

- Max pooling layers are usually added after convolutional layers to downsample, add translational invariance and hence make the network more robust against the presence of features in different spatial positions. Similar to a convolution, max pooling is done by sliding a kernel windows onto the input array, but now we select

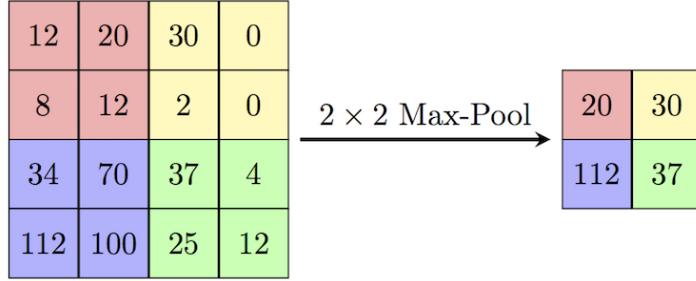


Figure 3.7: Max pooling operation with kernel size (2, 2)

the maximum value within the window.

### 3.2.3 Prediction uncertainty and Bayesian models

In classical supervised deep learning, a neural network with weights  $\theta$  is trained on some dataset  $\mathcal{D}$  to produce a minimal cost estimator  $\hat{\theta}$  for a predefined loss function. From the statistical point of view, this is a point estimate for each one of the network parameter. Point-estimate networks might lack explainability and generalize in overconfident ways to unseen data, and there is no obvious mechanism for such models to express their ignorance in such cases.

This is an analogous situation to classical parameter inference in statistic. From this perspective, the *frequentist* approach to parameter estimation can be compared to point-estimate networks. However, *Bayesian* statistics [33] has flourished in the last decades and is becoming the dominant statistical framework for data analysis inference problems. In the Bayesian paradigm, parameters are treated as random variables to reflect our ignorance about their “true” values. Our prior beliefs about the parameters  $\theta$  are then updated in the presence of new data  $\mathcal{D}$  using Bayes’ theorem:

$$P(H|D) \propto P(D|H)P(H), \quad (3.12)$$

where  $H$  is a certain hypothesis about  $\theta$ ,  $P(D|H)$  is known as the *likelihood* and encodes the relation between the data generation and the parameters, and  $P(H)$  reflect our prior beliefs. Bayesian frameworks have two main advantages. Firstly, they allow us to make our assumptions explicit by setting up the prior  $P(H)$ . This allows us to clarify, discuss and criticize prior knowledge in a clear way. Secondly, they provide a natural approach to quantify uncertainties in the inference parameters.

Bayesian neural networks are models that incorporate stochastic elements and that

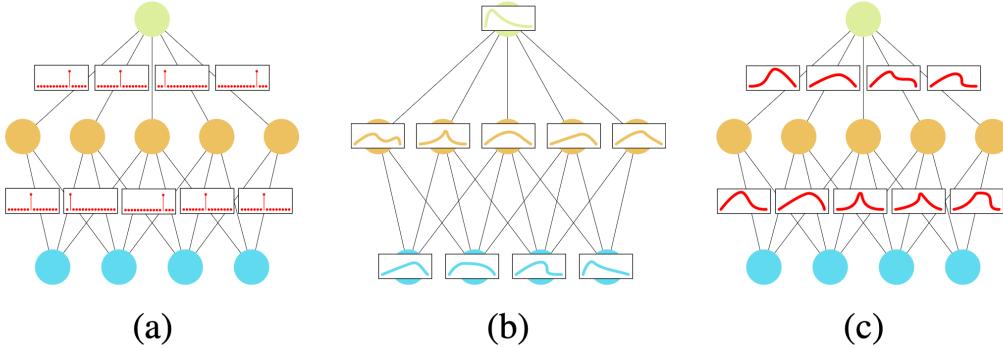


Figure 3.8: Illustration for a non-stochastic neural network a), a network with stochastic activations b), and a network with stochastic weights in c). Source: [34].

are trained using Bayesian inference techniques [34]. This is generally implemented either by considering stochastic weights in the network, or by considering stochastic activations, see figure 3.8.

A Bayesian neural network is defined by a prior distribution over the weights (when they are set to be stochastic), a functional model that forwards the inputs and generates the outputs, and a likelihood model that defines the predictive power of the network  $p(y|x, \theta)$ . The model weights are update in the presence of data  $\mathcal{D}$  according to Bayes' formula

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}_y|\mathcal{D}_x, \theta)p(\theta) \quad (3.13)$$

When the model has seen the data, we can use the posterior distribution on the parameters to generate data or make prediction:

$$p(y|x, \mathcal{D}) = \int_{\theta} p(y|x, \theta)p(\theta|\mathcal{D}). \quad (3.14)$$

Note that equation (3.14) weights the predictions by each possible parametrization of the network using the posterior distribution on the weights. In that sense, Bayesian networks can be understood as training an ensemble of networks and then averaging their predictions. Note that ensembles are also a popular technique with classical neural networks [35]. In fact, since the training of a classical network has also stochastic components (the initial weights, shuffling of the data, ...), it is common to train multiple networks with different initialization seeds and weights their predictions, much tend to outperform even the best-performing network.

Let us illustrate how a committee of networks can outperform even a top-performer

network. Suppose we train 3 independent networks  $A, B, C$  for a classification task with two classes, and that they all have the same error probability  $p$ . In an exercise of pure democracy<sup>5</sup>, consider a classifier  $D$  who picks the class with the majority of votes. Since  $D$  is wrong if and only if at most one of the classifiers is wrong, the probability of error for  $D$  is (considering  $p \rightarrow 0$ )

$$p^3 - 3(1-p)^2 = \mathcal{O}(p^2),$$

and hence  $D$  outperforms  $A, B$  and  $C$ .

### 3.2.4 Hyperparameter selection

The performance and convergence properties of a deep learning model are highly sensitive to its precise architecture, and hence, to the *hyperparameters* that control the latter: number of layers, the exact type of layers (convolutional, dense,...) and other auxiliary parameters that influence the training process. The optimal set of hyperparameters are the ones that produces the best-perfoming model on unseen data. There are three main difficulties when choosing hyperparameters. Firstly, it is not obvious a priori how modifying a certain parameter will affect the performance of a model. For example, if we add an extra layer to a network, will that improve expressivity and performance or will it lead to overfitting? This is a consequence of the low interpretability of deep learning models. Secondly, evaluating the performance of a model requires training it, which is computationally expensive. Finally, there are possibly an infinite number of possible architectures to explore.

Many Python APIs implementing routine to find optimal hyperparameters exists. They include state-of-the-art algorithms for exploring the hyperparameter space, selecting which parameters to explore, and early-stopping the evaluation of unpromising trials. In this work we use `OPTUNA` [36], which a Python optimization API, to select the appropriate set of hyperparameters for our deep learning models. `OPTUNA` implements a wide variety of searching algorithms to explore the hyperparameter space. Among those strategies, we can choose a naive grid search, where the user defines a grid over the hyperparameter space, and all possible combination are tested and ranked in a sequential order. `OPTUNA` can also implement a random search, which randomly selects the next trial over a grid. However, `OPTUNA` also implements more complex searching algorithms. The default search strategy is the *Tree-structured Parzen Estimator* (TPE) approach.

---

<sup>5</sup>“I love democracy. I love the Republic”, Senator Palpatine.

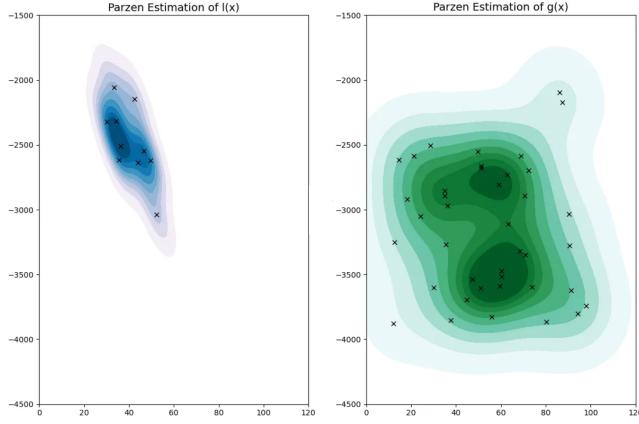


Figure 3.9: Example DE estimation for the distribution  $l(x)$  and  $g(x)$  used in the TPE algorithm.

TPE is designed to maximize the expected improvement when selecting a new sample. If we have already explored a set of points  $x$ , the expected improvement (EI) for a new trial  $y^*$  is defined as

$$\text{EI}_{y^*}(x) = \int_{-\infty}^{y^*} (y^* - y)p(y|x)dy = \int_{-\infty}^{y^*} (y^* - y)p(x|y)\frac{p(y)}{p(x)}dy , \quad (3.15)$$

where we are assuming a one-dimensional problem for clarity, and  $p(y|x)$  represents the probability of choosing a trial  $y$  having observed  $x$ , as obtained from the searching algorithm. TPE optimizes equation (3.15) by using the following routine [37]:

1. We initiate a set  $\mathcal{D}$  of explored trials with  $N_i$  parameters and compute their performance.
2.  $\mathcal{D}$  is split as  $\mathcal{D} = \mathcal{D}^l \cup \mathcal{D}^g$  in the top  $\gamma$  performers (a common value is  $\gamma = 0.2$ )  $\mathcal{D}^l$  and the lower performers  $\mathcal{D}^g$ .
3. We fit a two Gaussian mixture distributions to  $\mathcal{D}^l$ ,  $l(x)$  and to  $\mathcal{D}^g$ ,  $g(x)$ . This is a Kernel Density Estimation (KDE) for the distribution of both sets. See figure 3.9 for an illustration of this step.
4. The new trial parameters are selected and added to  $\mathcal{D}$  to maximize  $l(x)/g(x)$ . In practice, this can be done by generating  $N_s$  samples from  $l(x)$  and then maximizing  $l(x)/g(x)$  over those samples.

We can easily show that the TPE routine maximizes the expected improvement equation (3.15). Firstly, note that  $\gamma = p(y < y^*)$  and that  $p(x) = \int p(x|y)p(y)dy = \gamma l(x) + (1 - \gamma)g(x)$ . Furthermore, if  $y < y^*$ , we have that  $p(y|x) = l(x)$ , and hence,

$$\int_{-\infty}^{y^*} (y^* - y)p(x|y)p(y)dy = \ell(x) \int_{-\infty}^{y^*} (y^* - y)p(y)dy = \gamma y^* \ell(x) - \ell(x) \int_{-\infty}^{y^*} p(y)dy. \quad (3.16)$$

We conclude that

$$EI_{y^*}(x) \propto \left( \gamma + \frac{g(x)}{\ell(x)}(1 - \gamma) \right)^{-1}, \quad (3.17)$$

and so we need to maximize  $\ell(x)/g(x)$ , as previously discussed.

### 3.2.5 Loss function and training

In this section, we further detail to training process of a deep learning model, and describe the alhorithms that will be used to train our model. As we have already discussed in section 3.2.1, the goal of our model is the minimize the generalization error:

$$\varepsilon = \mathbb{E}_F[\mathcal{L}(Y, \hat{f}(X|\mathcal{D}_T))]. \quad (3.18)$$

The loss function  $\mathcal{L}$  should be choosen to encourage the model to produce realisitc outputs. For regression, where the goal is to estimate an output tensor  $Y$  from an input  $X$ , typical choices include the usual Euclidean norm or the absolute difference:

$$\mathcal{L}(Y, \hat{Y}) = \sum_i |\hat{Y}_i - Y_i|^p, p = 1, 2, \quad (3.19)$$

where  $\hat{Y}$  is the model prediction. This choise of loss function is minimzed when the model correctly predicts the target tensor. An important information to note is that equation (3.5) is generally unkown, since we don't have acces to all the possible realizations of the data, but only to a limited training dataset. As a consequence, we typically minimize the loss function computed on small samples of the training data, called *batches*, by averaging the loss over all samples in the batch. Observe that then, computing the average loss over a large batch will yield a closer approximation to the true loss function. The training process then uses these estimations of the loss to update the network parameters,  $\theta$ , using techniques that involve the computation of the loss gradient. This process is normallh repeated multiple times over the training dataset.

Each time the model sees the whole dataset is known as an *epoch*. A training loop with  $N_e$  epochs and batches of size  $B_s$ , with  $N_b$  batches, can then be written as:

---

**Algorithm 1** Classical deep learning training loop

---

```

while epoch <  $N_e$  do
    while batch <  $N_b$  do
        Loss  $\leftarrow \frac{1}{B_s} \sum_i \mathcal{L}(\hat{Y}_i, Y)$ 
        Compute  $\nabla_{\theta}(\text{Loss})$ 
        Update  $\theta \leftarrow \theta(\nabla_{\theta}(\text{Loss}))$ 
    end while
end while

```

---

In algorithm 1 there are two main steps we have not specified. Firstly, the computation of  $\nabla_{\theta}(\text{Loss})$ . This is trivial from the mathematical viewpoint, but it is not trivial to implement<sup>6</sup>. Most APIs implement a *backpropagation* algorithm for this step that we will not discuss here<sup>7</sup>. The second key aspect of a training loop is using the gradient  $\nabla_{\theta}(\text{Loss})$  to update  $\theta$ . This step in algorithm 1 was represented by  $\theta \leftarrow \theta(\nabla_{\theta}(\text{Loss}))$ . The precise way of using the calculated loss gradients to update the weights define the *optimizer*<sup>8</sup>. A naive way to optimize the loss at each iteration is to use the most straightforward implementation of *gradient descent*, where the parameters are updated using

$$\theta \leftarrow \theta - \alpha \nabla_{\theta}(\mathcal{L}), \quad (3.20)$$

where  $\alpha$  is a constant known as the *learning rate*. In general, optimizing the loss function  $\mathcal{L}$  corresponds to optimizing a complex (perhaps non-convex) function on a highly multidimensional space, with millions or billions of parameters. As a consequence, naive gradient descent might not be the best-performing optimizer. Additionally, recall the loss  $\mathcal{L}$  calculated for every batch is just an approximation of the global loss 3.5. The batch size determines how accurate the computed gradient is. Calculating the loss over all possible data points would yield the exact gradient. On the opposite extreme, calculating the gradient on a single data point would generate a biased gradient, and lead to a noisy gradient descent. The loss function can also have a complex structure. In the literature, a multitude of optimizers that improve and generalize the naive gradient descent exist. In this work, we use *Adam* [16]. Adam is a first-order momentum-based optimizer. Adam uses recursive update of the first and second moment of the gradient

---

<sup>6</sup>Note that this is purely a computer science problem.

<sup>7</sup><https://www.tensorflow.org/guide/autodiff>

<sup>8</sup>Note that this is purely a mathematical problem in optimization.

to update the learnable parameters. The Adam algorithm reads as follows:

---

**Algorithm 2** Adam optimizer

---

**Require:**  $\alpha$  : Learning rate  
**Require:**  $\beta_1, \beta_2$ : Moment decay rates  
**Require:**  $f(\theta)$ : Target function  
**Require:**  $\theta_0$ : Initial parameters  
**Require:**  $\varepsilon$ : Numerical tolerance constant

$m_0 \leftarrow 0$ : Initialize first moment  
 $v_0 \leftarrow 0$ : Initialize second moment  
 $t \leftarrow 0$ : Initialize time

**while**  $\theta_t$  not converged **do**

- $t \leftarrow t + 1$
- $g_t \leftarrow \nabla_{\theta} f(\theta_{t_1})$ : Get gradient
- $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$ : Update first moment
- $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ : Update second moment
- $m_t \leftarrow m_t / (1 - \beta_1^t)$ : Correct first moment bias
- $v_t \leftarrow v_t / (1 - \beta_2^t)$ : Correct second moment bias
- $\theta_t \leftarrow \theta_{t-1} - \alpha m_t / (\sqrt{v_t} + \varepsilon)$ : Update parameters

**end while**

---

Adam differs in two key features with respect to the naive gradient descent. Firstly, it includes two decay constants  $\beta_1$  and  $\beta_2$  used to include information about the past state of the optimization into the next current step. Since typical values are  $\beta \sim 0.9$ , recent gradient values contribute more to the current step update than older ones, but the inclusion of a momentum term allows the optimization to potentially escape local minima. Secondly, Adam includes a term  $v_t$  to account for the second moment in the gradient estimation. The running averages in Adam can potentially help to navigate noise functions by smoothing the local gradient.

Most current deep learning implementations are deployed or trained using HPC infrastructure. This becomes a necessity when training on large datasets. A quickly deployable way of parallelize training, and which requires minimal code modification is known as *Synchronous Distributed Training*. Synchronous Distributed Learning aims to harvest the computational power of multiple machines (GPUs,...) in the training stage by having different workers perform parallel calculations in a single batch. This can be done, for instance, by splitting the batch and sending each part to a different worker. In this work, we use the TensorFlow implementation `tf.distribute.MirroredStrategy`. This training strategy mirrors the model and its parameters in every worker (for instance, in every GPUs), slices the training batch and distributes them across all workers. Each

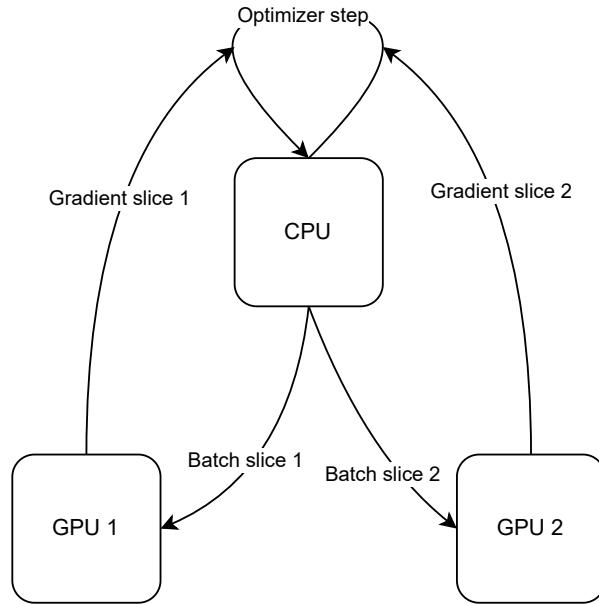


Figure 3.10: Distributed mirrored learning strategy with two workers (GPUs) and a CPU aggregating the gradient and updating the model parameters.

worker then calculates the loss and gradient in the corresponding batch slices. The gradient is then aggregated before updating the parameters and moving onto the next batch.

### 3.3 Workflow implementation: Recovering IGM conditions from the Lyman-alpha forest

This section describes the deep learning implementation used in this work and built in the ideas that have already been introduced regarding machine learning. The code used is largely based on [38] with very minor modifications and is available at <https://github.com/nicenustian/bh2igm> and based on TensorFlow and TensorFlow Probability.

The ultimate aim is to use a Bayesian neural network to recover the IGM gas density from an input Lyman- $\alpha$  skewer. Note that this is essentially equivalent to invert Equation 1.16, which describes the Lyman- $\alpha$  opacity as a function of the IGM gas conditions. As we have already discussed in 2.4, we will work in the optical depth-weighted space. This avoids two main potential difficulties in the analysis. Firstly, by working in this

new space, the network will not have to learn the relation between peculiar velocities and the Lyman- $\alpha$  opacity, which simplifies both the architecture design and learning phase. Secondly, it breaks the degeneracy introduced by peculiar velocities, since a shift in the physical space can produce the same opacity as a shift in the velocity space.

We specify the architecture design by the hyper-parameter list found in table 3.1. The global architecture that describes the number and size of the layers is specified by two hyper-parameters. Firstly, the parameter "Layers per block" is an integer list whose size is the number of blocks in the network and whose elements are the number of layers in each block. Secondly, the parameter "Filters per block" is also an integer list that specifies the number of convolutional filters of the layers for each block. If the architecture is a simple MLP, this parameter does not have any effect. The layers are chosen among a simple densely connected layer (MLP), and convolutional layer (ConvNet) or a residual layer (ResNet). Every feedforward pass through a convolutional or densely connected layer is followed by a batch normalization layer and by an activation function, see section 3.2.2. We use the PReLU activation function, which is a generalization of the ReLU activation with a learnable weight  $\alpha$  such that:

$$\text{PReLU}(x; \alpha) = \begin{cases} \alpha x & , x < 0 \\ x & , x \geq 0 \end{cases} \quad (3.21)$$

We append the same final block to all three potential architectures, which consists of a flattening layer transforming the tensor being manipulated to a one-dimensional vector, a dense layer and finally a Gaussian probabilistic layer. This final block includes the stochastic components of the neural network (hence the name Bayesian). For each target output density pixel the Gaussian layer outputs a full Gaussian probability distribution parametrized by the expected density and the standard deviation. This standard deviation will later be used as the estimation for the epistemic uncertainty in the network prediction. We illustrate the fiducial architecture in figure 3.11.

We optimize the network hyper-parameters using OPTUNA as discussed in section 3.2.4. Table 3.1 the optimized hyper-parameters and the range considered during the grid search. Note that these values can potentially depend on the nature of the training dataset, and hence may vary with redshift. In table 3.1 we present the fiducial architecture at  $z = 4.4$ . This optimizing process is automatically carried out whenever the redshift is varied.

For each redshift, we train a different network using the Sherwood simulation suite presented in section 2. Recall that we consider two training datasets, SHERWOOD and

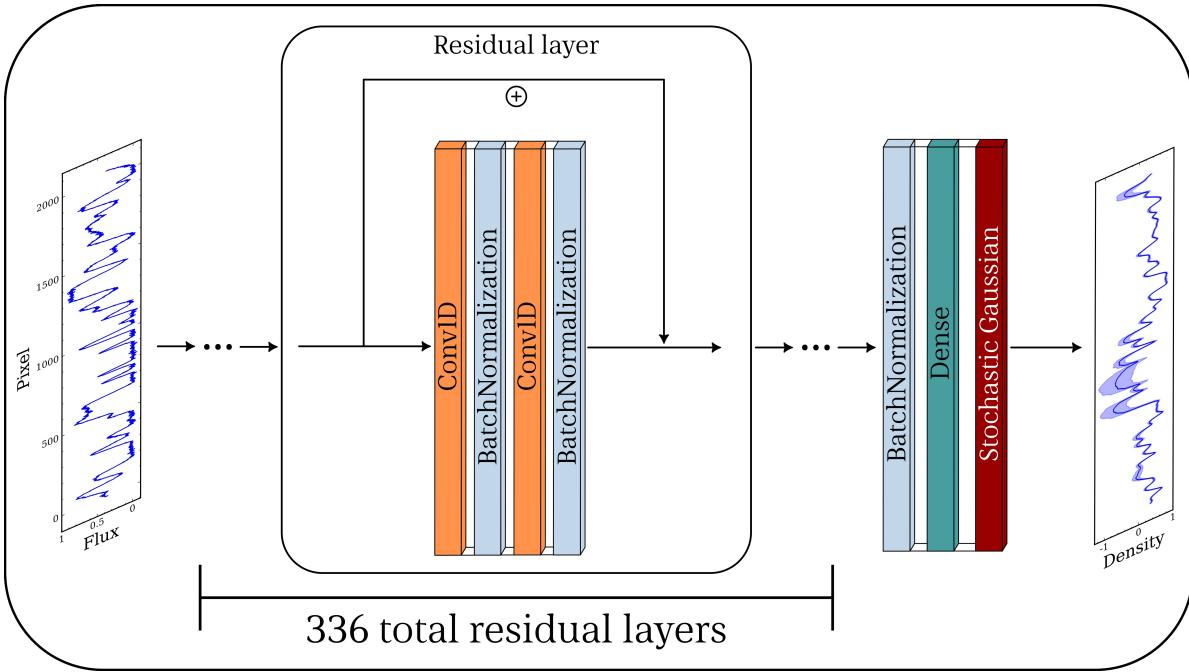


Figure 3.11: Fiducial architecture for our Bayesian neural network trained at  $z = 4.4$  on the Sherwodd simulation suite. The fiducial parameters can be found in table 3.1.

SHERWOOD THERMAL, where the latter includes variations in the thermal parameters and reionization history. See [39] for more details. The inclusion of thermal parameters will lead to a more robust network when used on unseen real data, but its performance will naturally be lower than the network trained on a single thermal hsitory. We use SHERWOOD for an initial model testing, and SHERWOOD THERMAL for our final anaylis on real data. 80 % of the data is used for training, and the rest is used for validation purposes. For reference, the training dataset SHERWOOD with,7 different WDM models has a total size of  $\sim 3.5$  GB. We use the Raven<sup>9</sup> HPC cluster at the Max Planck Computing and Data Facility to train our networks. With 4 Nvidia A100 GPUs, a typical training time is  $\sim 1$  hour, depending on the model's exact architecture and number of parameters.

The network's input consist of a Lyman- $\alpha$  flux skewer, and the network's output consists of the mean density and standard deviation at each pixel. The labelled training data consists of individual Lyman- $\alpha$  flux skewers with their associated density optical depth-weighted density field,  $(F, \log(\Delta_\tau))$ . For each labelled training pair, the loss function is taken to be the negative log-likelihood (NLL) for the normal distribution

<sup>9</sup><https://www.mpcdf.mpg.de/services/supercomputing/raven>

Hyper-parameter	Min.	Max.	Best value
Layers per block (int)	1	4	[1, 2, 4 ,4]
$\log_2(\text{Filters per block})$ (int)	1	5	[4, 5, 5, 5]
Number of blocks (int)	1	4	4
$\log_2(\text{Batch size})$ (int)	3	8	3
Learning rate ( $\log_s$ , float)	$10^{-4}$	0.1	0.004937
Network (MLP, ConvNet, ResNet)	...	...	ResNet

Table 3.1: Hyper-parameter grid search for the fiducial model at  $z = 4.4$ . “ $\log_s$ ” indicates the parameter is sampled in the log domain. “Int” and “float” mean they are sampled as integers or floats, respectively.

that the network parametrizes:

$$-\log \mathcal{L} = \frac{1}{N} \sum_i \left( (Y_i - Y_{i,\text{pred.}})^2 / \sigma_{i,\text{pred.}}^2 + \log \left( \frac{1}{\sigma_{i,\text{pred.}}^2} \right) \right), \quad (3.22)$$

where the sum runs over all skewer pixels,  $Y_i$  is the real density at pixel  $i$ ,  $Y_{i,\text{pred.}}$  is the predicted expected density and  $\sigma_{i,\text{pred.}}^2$  the predicted standard deviation.

The input skewers are processed as follows. Firstly, the input flux is rebinned into a target number of pixels to match the real data that will be used. Downsampling is done by taking the average flux over nearby pixels, while upsampling is done by appending to every pixel a copy of itself. The flux is then convolved using a gaussian kernel to simulate a given instrumental resolution. In most of our analysis the resolution is taken to be 6 km/s, in accordance with state-of-the-art spectrographs, such as UVES [40]. During training, the training data is stacked and randomly shuffled to ensure a correct mixing and representativity of the models. Lastly, we process the rescale input optical depth to match the observed mean flux at the considered redshift [41]. The training data is augmented on-the-fly in two ways. We first roll the input spectra by application translations, this helps the network learning dynamical features, independently of their positions. We also add random uncorrelated gaussian noise to simulate a finite signal-to-noise ratio (SNR). The default SNR per pixel used for testing purposes is 50. When applying our methods to real-data, the network is retrained using a noise model for each target object.

We use the Adam optimizer with fixed moment decay rates as implementation in the TensorFlow API and a learning rate set by the Optuna grid search. The training is evaluated using two main metrics: the loss function NLL, and the mean absolute error, MAE, defined as the Euclidean norm 3.19 in with  $p = 1$ . To prevent over-fitting, the

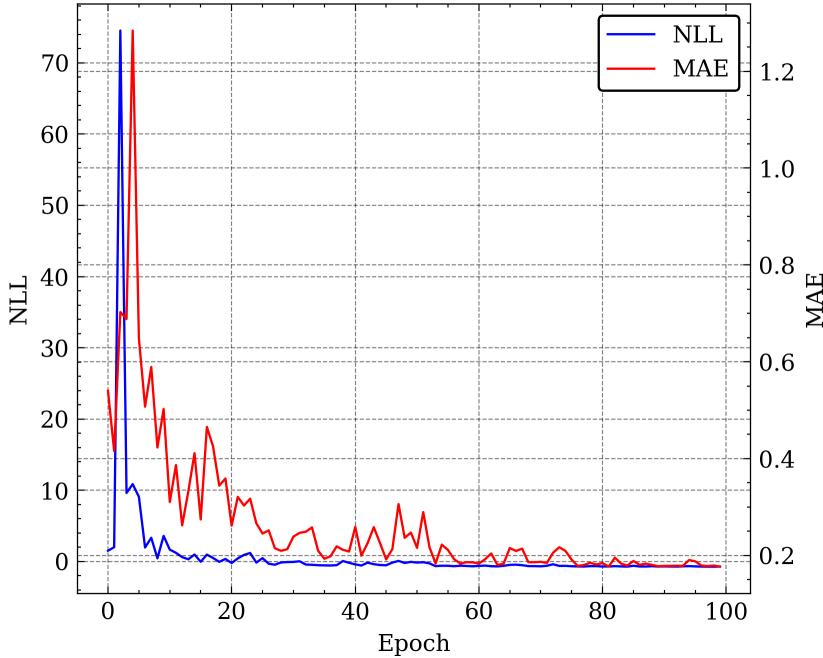


Figure 3.12: Learning curve for our fiducial architecture at  $z = 4.4$  on the **SHERWOOD** dataset. The figure shows the NLL and the MAE on the validation split as a function of the epoch.

network parameters are only updated if the NLL loss improves on the validation split of the data. A policy to halve the learning rate if there is no loss improvement in the test set for 10 epochs is also included. This is the most straightforward adaptation of an adaptive learning rate. Figure 3.12 shows the evolution of the NLL and the MAE on the validation split during training for our fiducial architecture at  $z = 4.4$  and the **SHERWOOD** suite. It is interesting to note that while the MAE is always positive by definition, negative values of the NLL are consistent with its definition. Figure 3.12 demonstrates that, for our problem, the network training converges in  $\sim 100 - 150$  epochs. Recall that we always work with  $\log(\Delta_\tau)$ , which also serves as a regularization step with respect to the density.

In Figure 3.13 we show an example  $20\text{h}^{-1}\text{cMpc}$  Lyman- $\alpha$  validation skewer from the **SHERWOOD** dataset. The top panel shows the input flux to the network. The bottom panel shows the true  $\Delta_\tau$  density fields, the (mean) recovered densities and the  $1\sigma$  envelope predicted by the Bayesian network. Note that, in the spectral regions with large features and variations in the flux, the predicted mean density closely follows the true field. In those regions, there is enough physical information for the network to accurately recover  $\Delta_\tau$ . In contrast, in the saturated regions with low flux, the noise dominates,

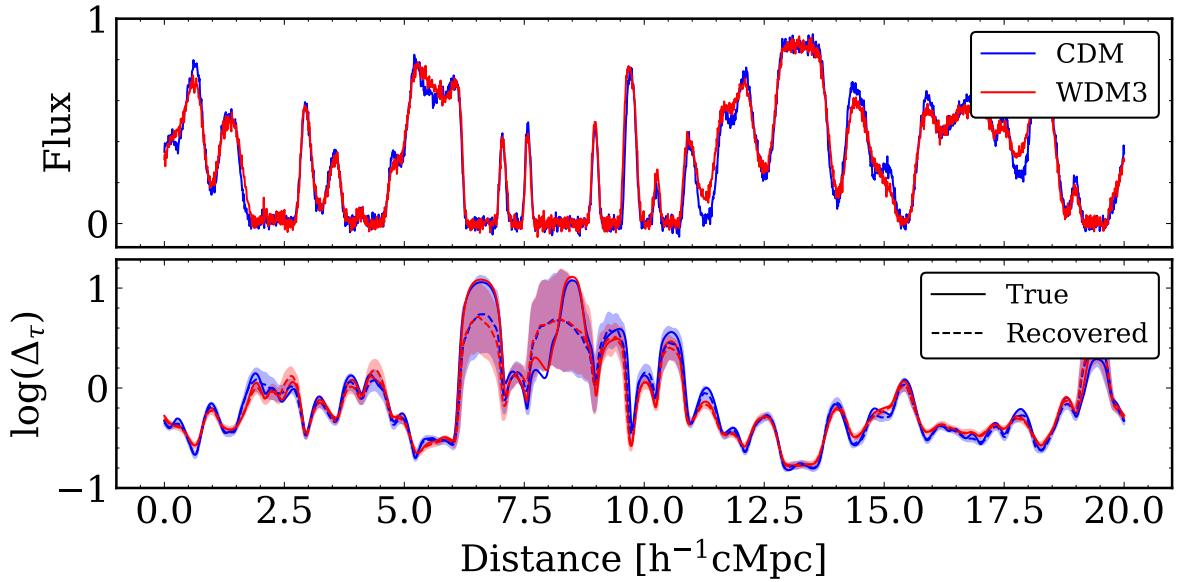


Figure 3.13: An example  $20h^{-1}\text{cMpc}$  Lyman- $\alpha$  validation skewer for the CDM and WDM3 Sherwood models. The top panel shows the input flux to the network. The bottom panel shows the true  $\Delta_\tau$  density fields, the (mean) recovered densities and the  $1\sigma$  envelope predicted by the Bayesian network.

and the network predicts larger uncertainties (observe, for instance, the saturated region in Figure 3.13 around  $8h^{-1}\text{cMpc}$ ). This should be regarded as a strength of Bayesian networks, since they are able to accurately detect and make explicit situations where the predictions should not be trusted. To minimize the mean error in regions where the network cannot make accurate predictions, note how there is a bias towards quasi-constant mean prediction. This is visible around  $8h^{-1}\text{cMpc}$ , where the true density has a steep increasing profile, and the prediction has an almost constant u-shaped profile.

On the SHERWOOD validation split, the network reaches a  $1\sigma$  accuracy of 79%, meaning that 79% of the pixels are correctly predicted within  $1\sigma$ . Note that this is a larger accuracy than expected from purely normally-distributed densities.

,.... violin...compare thermal models and non thermal

### 3.4 Recovered field statistics and uncertainties

Once we have a machinery to recover the IGM density from a Lyman- $\alpha$  skewer, we would like to use this density field information to do statistical inference on the physical parameters of interest. In this work, since the WDM mass directly affects the matter distribution of the Universe (see Figure 1.2) we would like to infer the WDM masses

from the recovered  $\Delta_\tau$  field. Now, note that the exact  $\Delta_\tau$  field of a skewer, such as the one shown in Figure 3.13, not only depends on a set of physical parameters (WDM mass, temperature,...) but also on the initial conditions or equivalently, the random seed in the simulation. As a consequence, we will not compare densities on a sightline by sightline basis, but rather compute and compare aggregated statistics of the fields over multiple realizations that capture global statistical properties, and not simulation-specific characteristics. Statistics that have been amply tested in the literature are the Power Spectrum (PS), the Probability Distribution Function (PDF), the curvature, or the autocorrelation function (see [42] and [43] for more details). In this work we will focus on the  $\Delta_\tau$  PDF. In section 4.6.4 we will address the optimality of this choice.

We will carefully describe the statistical inference pipeline in section 4.1 . The main idea will be to fit the density PDF obtained from the Bayesian neural network using observed skewers to the PDFs from each WDM model in the **SHERWOOD** dataset. An important consideration for a successful pipeline will be to compare *apples to apples*, i.e., to compare statistics that have been obtained through the same process to reduce any potential bias. For that reason, we will always work with the ML-recovered density by our fiducial neural network, and compare the *observed* statistics to the *model* statistics from each WDM run. The model statistics are computed over the whole simulation box using all 5000 skewers from each run. The observed statistics are typically computed over a much smaller number of observed sightline. In this section, we analyse the statistics recovered from the predicted  $\Delta_\tau$  field by the Bayesian network, including how to calculate their uncertainties, which will be crucial for the rest of the inference process.

### 3.4.1 Noisy regions and masking

Recall Figure 3.13 are the central example in this regression problem. The flux-saturated regions where the flux is close to zero and dominated by noise have no physical information to recover the density. Potentially, any density higher than the truth could also produce a saturated profile. The model correctly detects these regions and produces larger uncertainties, but the mean predictions have no physical information and should not be considered in any posterior analysis. Hence, we will mask these regions when computing the model and observed  $\Delta_\tau$  PDFs. This has the effect of modifying the shape of the PDF by removing part of the high-density tail. We show in Figure 3.14 the model  $\Delta_\tau$  PDFs for the **SHERWOOD** runs computed using the recovered density by our neural network and explicitly highlight the effect of masking the saturated regions that have a flux  $\leq 3\frac{1}{\text{SNR}}$ . As can be noted, the artefact that occurs at high densities

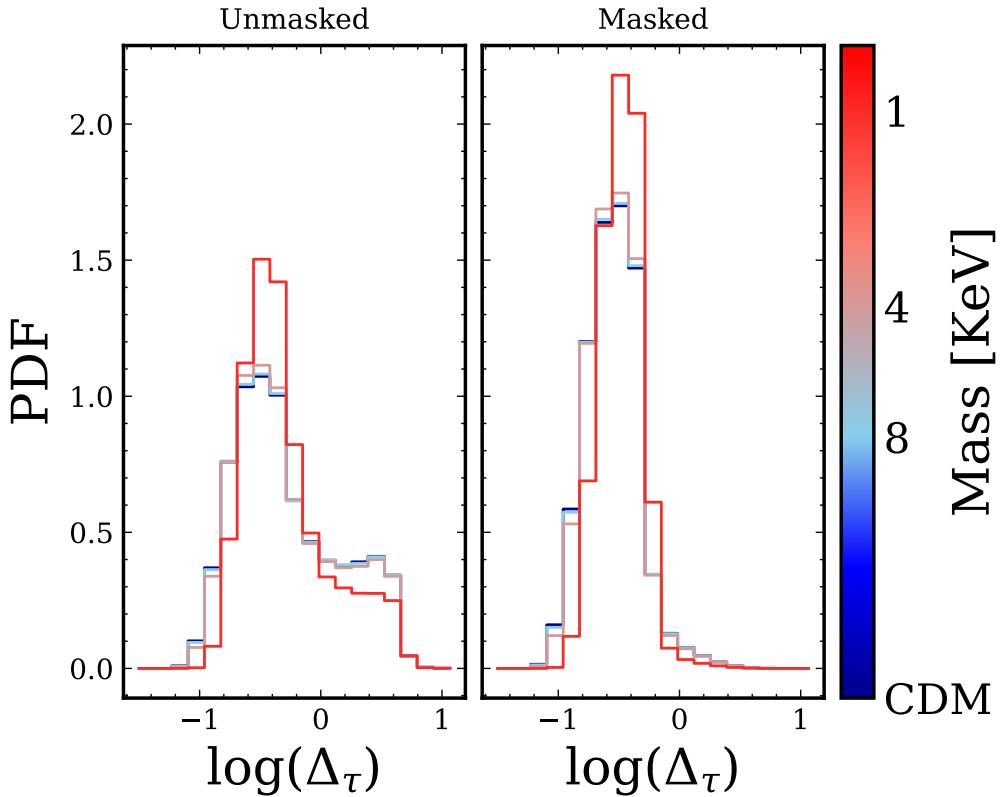


Figure 3.14: Model  $\Delta_\tau$  PDFs for the SHERWOOD runs computed using the recovered density by our neural network. We explicitly highlight the effect of masking the saturated regions that have a flux  $\leq 3\frac{1}{\text{SNR}}$ .

( $\log(\Delta_\tau) \sim 0.5$ ) disappears when masking. It can be interpreted as a bias introduced in the saturated regions.

### 3.4.2 Uncertainty in the recovered statistics

When estimating uncertainties for the recovered PDF, there are two main sources to consider. Firstly, the Bayesian ML predictions have a certain uncertainty associated to them, which transfer to the density PDF as an additional scatter. We include such uncertainties by considering the predicted density field  $\Delta_\tau$  resample 1000 times the posterior distribution produced by the network. We do this process for every observed skewer, compute the 1000 PDFs for each sample, and then the standard deviation between the 1000 statistics. The resampling process requires careful consideration.

ADD bootstrapping and PDF OF PDF

PDF:

$$p(\lambda|f(x), N) \propto \frac{\Gamma(N+1)}{\Gamma(N\lambda+1)\Gamma(N-N\lambda+1)} f(x)^{N\lambda} (1-f(x))^{N(1-\lambda)} \quad (3.23)$$

ADD MATRIX CORR of cdm vs wdm and argue it always use cdm and resampling and recoveries using finite skewers

$$r_i = \frac{\Delta_{\tau,i} - \mu_i}{\sigma_i} \quad (3.24)$$

## 3.5 Model interpretability and limitations

Deep learning models tend to have, by definition, numerous parameters. As a consequence, giving an interpretation for individual model parameters is far from being a trivial task. On top of the large number of parameters, the nonlinearities and the potentially biased and incomplete datasets can lead to complex training behaviors. In that regard, deep learning models have classically been regarded as “black-box” models. They are often more accurate than simpler statistical models, but lack explainability. Great efforts have been recently made in understanding the learning dynamics of neural networks [44], [45], [46]. Due to the complexity of this interpretation task, here we choose to only give a qualitative analysis of some aspects that can help gain intuition on how the network operates internally.

### 3.5.1 Saliency analysis

Although many open-source libraries, such as *Captum*<sup>10</sup>, implement popular methods for deep learning model visualization and interpretation, in this work we use TensorFlow’s built-in differentiation capabilities. In particular, we use Automatic Differentiation to compute the *saliency* score of the model, defined as the gradient of the model output with respect to the model input. As a consequence, for every target density pixel, the saliency score measures which flux pixel variation contributes the most to a change in density. Saliency is a simple score giving us insights into how the model uses flux to reconstruct densities. Additionally, note that calculating such gradients do not require any numerical finite difference approximation, since TensorFlow’s GradientTape class can build an exact computational graph with all the operations performed on the input

---

<sup>10</sup><https://github.com/pytorch/captum>

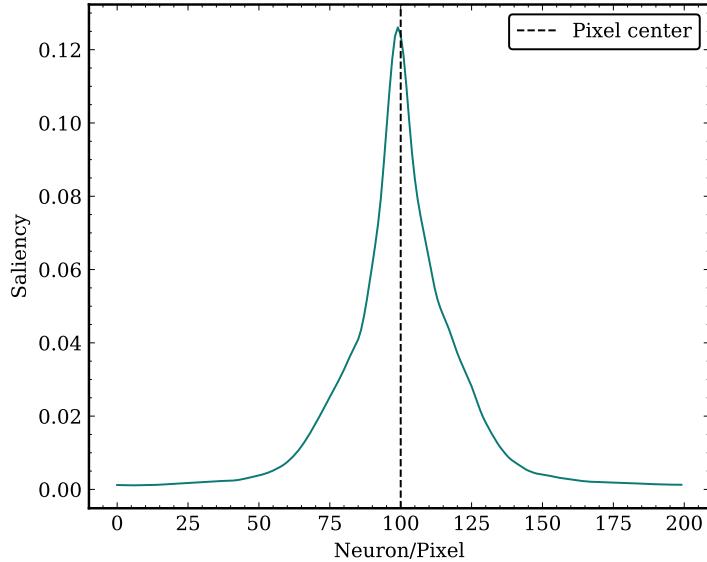


Figure 3.15: Saliency of the CDM and WDM3 Sherwood network at  $8\text{h}^{-1}\text{cMpc}$  Lyman- $\alpha$  flux.

flux. Saliency maps are a popular technique in applied science to explore deep learning model. In fact, in many physics research articles, it tends to be the only tool used for this purpose. This highlights the need to develop more robust and easily-deployable frameworks for deep learning models.

The saliency at density pixel  $i$  and input flux pixel  $j$  is simply

$$\text{Saliency}(i, j) = \frac{\partial \mu_i}{\partial f_j}. \quad (3.25)$$

Figure 3.15 shows the saliency profile averaged over all output density pixels. This gives a global average metric on the relevant flux pixels to predict the density at a certain pixel. The typical saliency “length”, that is, the typical number of pixels from the pixel center that are useful to predict its density is  $\sim 30$ . Since at redshift  $z = 4.4$  our pixel scale on the **SHERWOOD** data is 1.3 km/s, we obtain that in velocity space a window of  $\sim 40$  km/s is needed to recover the density at a pixel. This particular dynamic is coherent with the underlying physical process, which is a strong robustness sign of our machine learning model. In fact, for the typical IGM gas in the **SHERWOOD** suite, the Lyman- $\alpha$  absorption cross-section profile has a typical scatter of (see Equation 1.16)  $b(T) \sim 20$  km/s.

### 3.5.2 Covariate shift

As we have already mentioned when discussing Figure 3.13, saturated regions lead to larger uncertainties in the network’s predictions, reflecting the fact that noise dominates over the physical signal. Observe again the saturated region around  $7\text{h}^{-1}\text{cMpc}$  and  $8\text{h}^{-1}\text{cMpc}$  in Figure 3.13. Both of these regions have completely different density profiles, but the network predicts the same u-shape profile with a similar peak density. This peak density is slightl above the mean density in the simulation box and is similar for the CDM and WDM3 models. We can interpret this as the network leaning a unique mean “high-density” value for saturated regions and the whole dataset. The way the model predicts those mean densities based on the training dataset is, again, far from being trivial. However, it is clear that the training data contains all the information leading to any possible bias/characteristic in the deep learning model predictions. An important effect regarding the avaible data is *covariate shift*. Covariate shift occurs when the training and validation data are sampled from different distributions. Covariate shift can be difficult to detect and can lead to biased predictions. In this work, we try to test the generalization capabilities of our model by adding random noise and validation on unseen data from multiple hydrodynamical codes with different specification, to make sure the model is learning the relevant physics. This is a pragmatic approach to address covariate shift, since more formal techniques are out of the scope of this work. Nonetheless, we are aware an insufficient dataset can significantly degrade performance. As an illustration of the undesirable effects of covariate shift, we retrain the fiducial architectures on a subset of the **SHERWOOD** dataset where we iteratively eliminate each one of the WDM models. We call this set of models **NOTRAIN**. For instances, **NOTRAIN-WDM1** has not been train on the WDM1 data. In Figure 3.16 and Figure 3.17 we show the models **NOTRAIN-WDM1** and **NOTRAIN-WDM4** predicting, respectively, on WDM1 and WDM4 skewers from the **SHERWOOD** dataset.

Note the clear covariate shift effect in Figure 3.16, where the **NOTRAIN-WDM1** have a systeatic bias towards lower densities. In comparison, in Figure 3.17 the model **NOTRAIN-WDM1** shows no clear indication of any bias, and is almost indistinguishable from the fiducial model. In the former case, the model has to extrapolate on data generated from an unseen model, which gives clearly an incorrect predidction. In the later case, the model has been trained on CDM and WDM3, so predicting on WDM4 is an interpolation task.

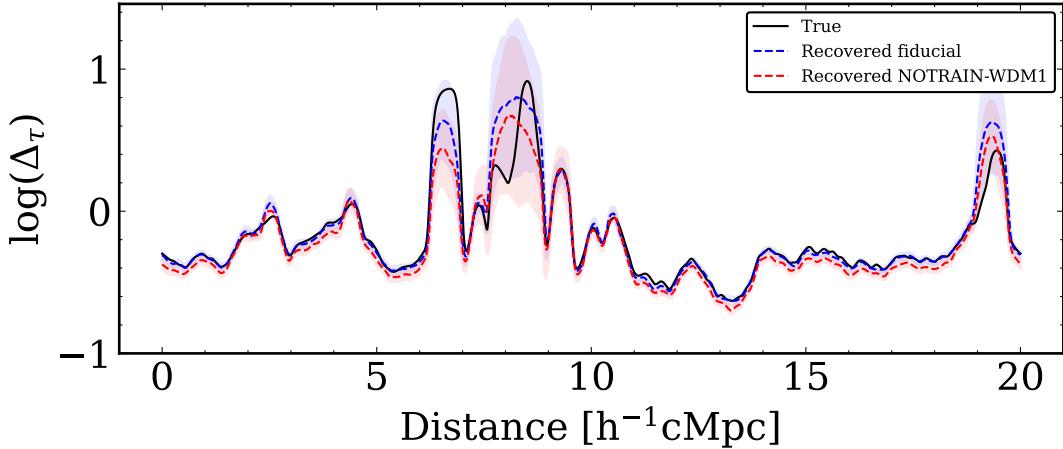


Figure 3.16: The NOTRAIN-WDM1 model predicting on a WDM1 skewer, compared to the fiducial model trained on the whole SHERWOOD dataset.

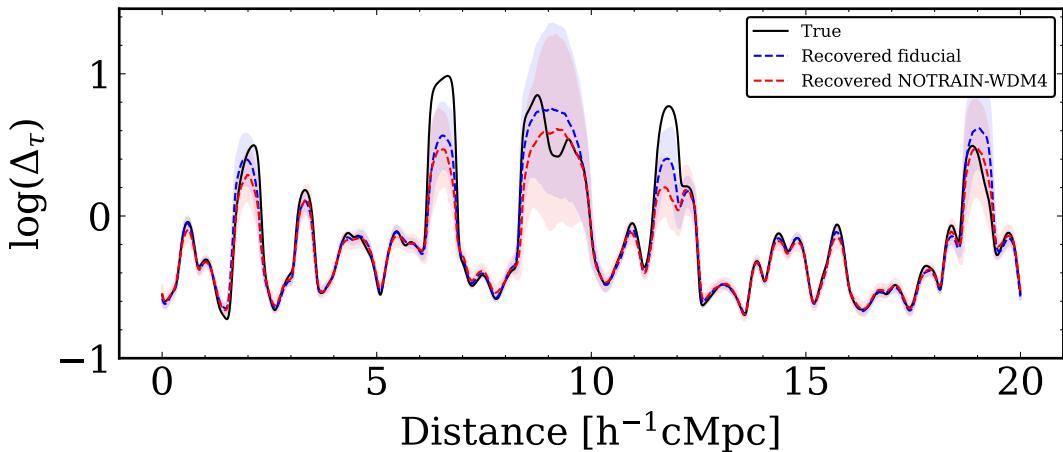


Figure 3.17: The NOTRAIN-WDM4 model predicting on a WDM4 skewer, compared to the fiducial model trained on the whole SHERWOOD dataset.

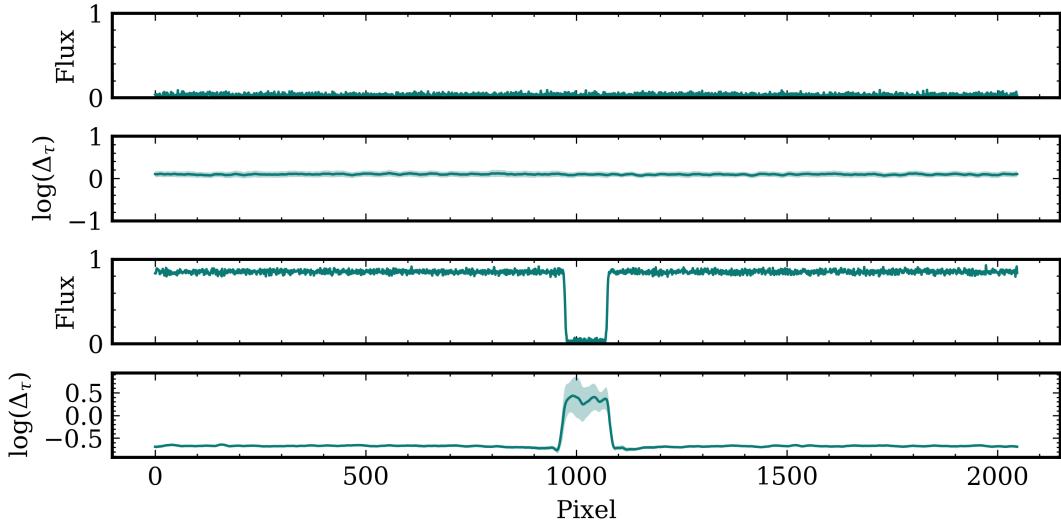


Figure 3.18: The NOTRAIN-WDM4 model predicting on a WDM4 skewer, compared to the fiducial model trained on the whole SHERWOOD dataset.

### 3.5.3 Extreme covariate shift and malicious data

In regression tasks related to computer vision, it is customary to mask certain parts of the input to explore how the model reacts to specific features. Following this idea, we have designed a set of malicious Lyman- $\alpha$  skewers that include artefacts and unnatural features that are not present in the training dataset. Therefore, this can be considered as an extreme case of covariate shift. In Figure 3.18 we show the fiducial model predictions for two such skewers. In the first case, we consider a fully saturated skewer with zero flux and only a small amount of noise. Observe how the model makes an educated guess and predicts the mean density for the whole skewer. However, the model is very confident in the predictions and produces minimal uncertainties. This is not correct from the physical standpoint, since any density higher than the truth can produce a saturated flux. In the second skewer we have complete transmission except a saturated absorption feature. The predictions for the pixels with complete transmissions are below the mean density, and on the saturated pixels it is above the mean density with a higher uncertainty. This is closer to a “real” absorption feature with more accurate prediction.

### 3.5.4 Model pruning

Pruning refers to the idea of eliminating certain components of a deep neural network, or equivalently, considering a sub-network. The motivation to prune a network is generally to reduce the number of parameters, which can improve interpretability, and allow

users to better trace the flow of information through the network. Since most deep network are over-parametrised, pruning can improve training while maintaining a very similar performance to the original architecure. Pruning is an active area of research, and there are efforts being made to use domain-specific expertise to prune networks and make informed choices about their topology [47]. In this work, we have used OPTUNA in section 3.2.4 to tune the hyperparameters of the model on the sole basis of a performance metric function. This naturally over-parametrises the network, with connections and filters that might be redundant. This means that there is a potential room for improvement when it comes to finding more interpretable networks with fewer parameters and a more stable trainig.

To illustrate how pruning can affect the performance of our model, we briefly focus on the first convolutional layer of the architecutre in Figure 3.11. In Figure 3.19 we show the activation of the first 4 convolutional filters for a trained network with architecute [2, 2], [4, 16] for a randomly selected skewer, following the notation in Table 3.1. Note that out of the four filters, the blue and green ones are scaled versions of each other, and similarly, the red and orange filters are also scaled versions of one another. To test of this redundancy affects the network, the prune half of the filters in this first layer to an architecture [2, 2], [2, 16]. Figure 3.20 shows the training NLL curve for both architecures. As can be seen, they reach a simialar performace in terms of the NLL, but the pruned architecure has a more stable learning curve and convergence in a smoother way.

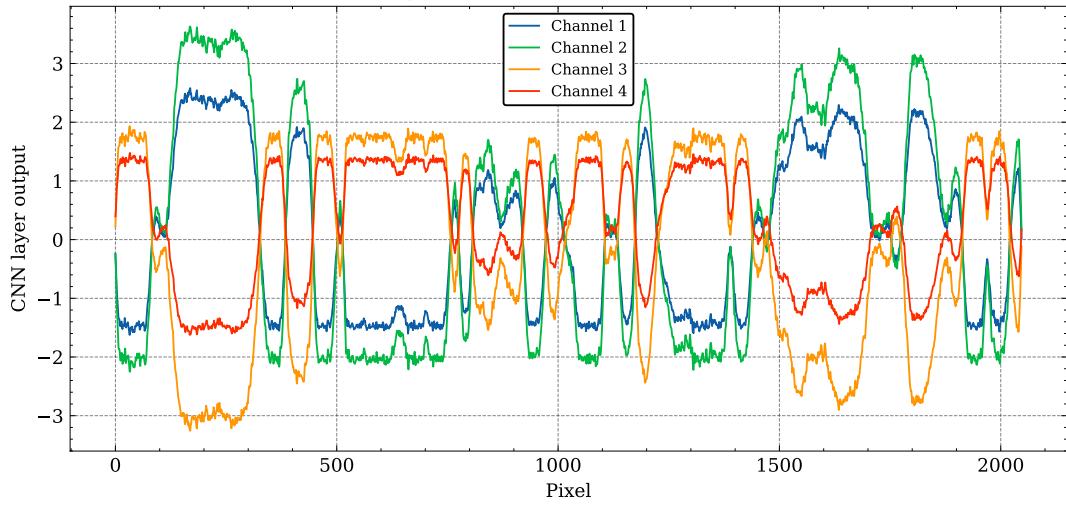


Figure 3.19: First layer activation of the convolutional filters for a trained network with architecture  $[2, 2], [4, 16]$  for a randomly selected CDM skewer from SHERWOOD.

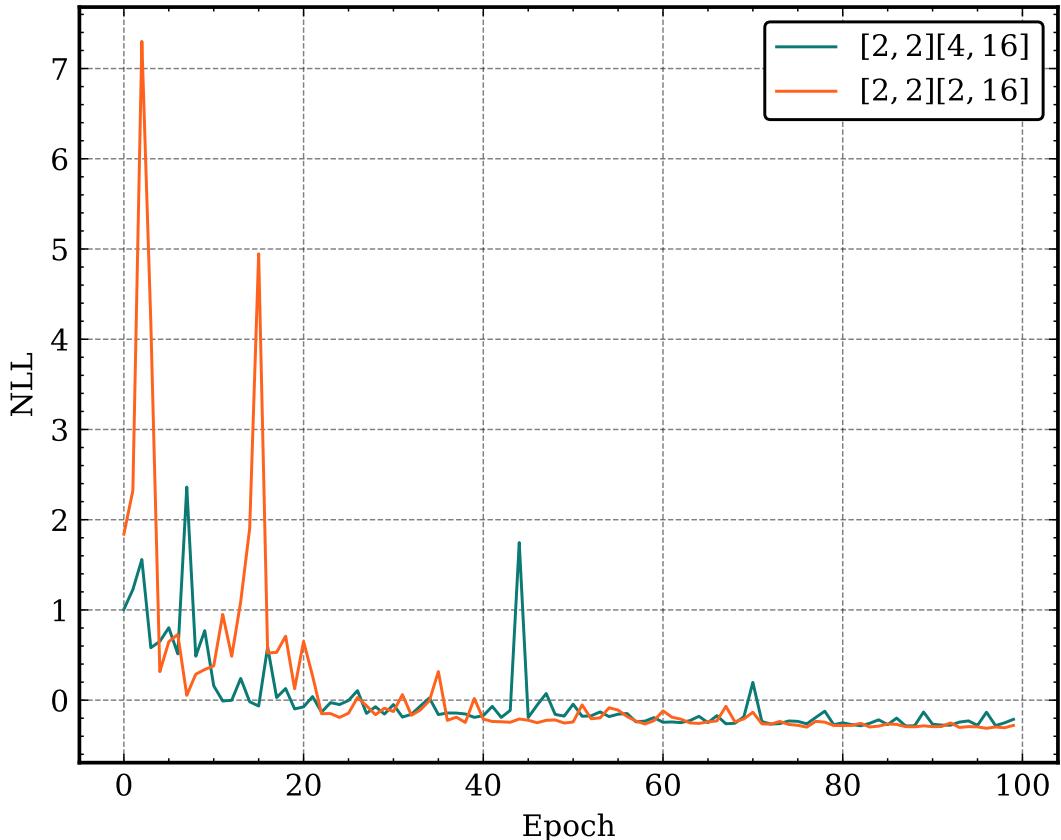


Figure 3.20: Training curves for the NLL comparing the model  $[2, 2], [4, 16]$  and the prunned architecutre  $[2, 2], [2, 16]$  with half of the convolutional filters in the first layer.

# 4 Constraining Warm Dark Matter at the density level

## 4.1 Inference pipeline: from Lyman-alpha skewers to WDM constraints

In section 3 we have given a detailed analysis of our Bayesian deep-learning algorithm to recover IGM densities from a Lyman- $\alpha$  skewer. The baryonic density of the IGM is sensitive to the WDM mass through a clear physical mechanism related to gravitational clustering. In this section, we use the recovered IGM density fields to constrain WDM candidates. Note that the natural observable quantity related to the Lyman- $\alpha$  forest is the flux. As a consequence, an almost omnipresent choice in the literature has been to work directly with summary quantities on the flux, which is a proxy of the underlying density. Such summaries include the power spectrum (PS) [6], the curvature [48], the probability distribution function (PDF), etc. The deep-learning approach introduced in section 3 allows us to directly recover the baryonic density field along a line of sight, thus having full access to the field level IGM properties. In this section, we use the recovered  $\Delta_\tau$  fields by our neural network to constraint WDM directly at the density field level. Recall Figure 1.2 and Figure 2.1 showing how different WDM models affect the density field. We strive to capture that difference in the WDM models to constrain which free-streaming length are compatible with QSO observations. Note that for a given line of sight, the precise value of the density field not only depends on the WDM masses (and other possible physical parameters), but most crucially it depends on the random density fluctuation that have seeded the gravitational collapse process. Equivalently, in a simulation setting, the obtained densities would depend on the seed used to initialize each simulation. This makes it infeasible to compare a given Lyman- $\alpha$  skewer to a simulated one, and means that we must use aggregated summaries over multiple skewers that capture the global properties of the field. In this work, we perform the inference using

the density PDF as the summary statistic of choice. This is a well-tested and robust statistic [42]. In section 4.6.4 we give an additional argument, based on Information Maximising Neural Networks, to support this choice of summary statistic.

The basic working principle of inference pipeline is to fit the observed  $\Delta_\tau$  PDF, with its associated uncertainties, to the corresponding  $\Delta_\tau$  PDF produced by each WDM model. To compare similar quantities, we always work with the recovered field by our neural network from section 3. Note that in the **SHERWOOD** only a finite number of DM models are available, due to the computational cost of running this simulation. In more detail, we only have access to the models CDM, and WDM1,2,3,4,8,12. We smoothly and linearly interpolate the PDF by interpolating each PDF bin to generate a  $\Delta_\tau$  PDF in the range of WDM masses from 0 to 1 KeV. Note that, since we expect the real observations to fall close to the CDM model, and have multiple simulations close to CDM, we expect this interpolation not to limit the inference pipeline. See Figure 2.1 again and observe how similar are the PDF for CDM and WDM3. For more massive models in **SHERWOOD** the PDF converge to the CDM PDF.

For each DM model of inverse mass  $m$  we denote by  $\text{PDF}(m)$  the  $\Delta_\tau$  PDF computed over the recovered densities by our NN network over all available sightlines in the **SHERWOOD** or **SHERWOOD THERMAL** datasets. We refer to this as the model PDF. Let us denote by  $\widehat{\text{PDF}}$  the recovered PDF for a target set of (observed) skewers. Then, we fit  $\widehat{\text{PDF}}$  to the model PDFs using a simple  $\chi^2$  fit:

$$\chi^2(m) = \sum_i \frac{(\text{PDF}(m)_i - \widehat{\text{PDF}}_i)^2}{(\sigma_i)^2}, \quad (4.1)$$

where the index  $i$  refers to each PDF bin and  $\sigma_i$  are the uncertainties on the observed data. If the data is normally independently and normally distributed, the quantity in Equation 4.1 follows a  $\chi^2$  distribution [49]. The model that minimises the quantity is the best-fit model, on which we can compute uncertainties and obtain a confidence region of compatible models with the observed data. Since we are only fitting a single parameter model, the 1 and 2-sigma confidence regions on the WDM mass are given, respectively, by the boundaries of

$$\chi^2(m) - \chi^2_{\min} = 1, 4, \quad (4.2)$$

where  $\chi^2_{\min}$  is the best-fit  $\chi^2$  value. In the following, we will be interested in the  $2\sigma$  confidence regions. This region can be interpreted as the set of WDM models that guarantee

to contain the “true” model with a 95% probability. In the current literature, WDM constraints are often reported as the  $2\sigma$  upper limit, where the lower limit typically corresponds to CDM. The current more stringent  $2\sigma$  WDM limit constraints are  $\sim 3$  KeV, see [6] and [39]. Note that this fitting procedure is non-Bayesian, in the sense that we don’t include any prior knowledge or use Bayes’ theorem. Again, this procedure is compared in section 4.6.4 to a IMNN Bayesian fit, leading to similar results.

## 4.2 Inference testing on Sherwood spectra under realistic observational conditions

In this first section we run our inference pipeline from section 4 using a set of toy observed skewers. More precisely, we use our neural netowrk trained on different subsets of the SHERWOOD dataset and use validation SHERWOOD skewers as the “observed” skewers.

### 4.2.1 Untrained DM models

We begin by testing the robustness and inter(extra)polation capabilities of the neural networks by considering the **NOTRAIN** models that are trained on data that iteratively excludes each one of the WDM models. For each one of those trained neural networks, for instance **NOTRAINWDM4** (which was not trained on WDM4), we predict on the WDM4 sightlines, compute the recovered  $\Delta_\tau$  with its uncertainties according to 3.4 and run the inference pipeline. We also perform additional variations by running the pipeline only on the PDF bins whose value is greater than a fixed constant. This has the effecto of only fitting the peak of the PDF, and neglecting the low-information tails. Figure 4.1 summarises this inference tests. Each plot corresponds to a different fit combining predictions from each **NOTRAIN** model with each of the masks applied to the PDF when fitting. The light and dark blue regions correspond to the 1 and 2 sigma confidence regions. The red line corresponds to the true DM model mass, and the black line to the best-fit model that minimises the  $\chi^2$ . The blue curve is the  $\chi^2$  metric. Note that we are using all 5000 sightlines on the inference step. This is not a realistic sample size, but rather a test to the inter(extra)polation of the pipeline.

Recall that on WDM2,3,4 the models are interpolating. Observe that as a consequence, the recovered mass is consistently recovered within the  $1\sigma$  region. In contrast, with the models CDM and WDM1, the neural networks have to extrapolate on unseen DM models. As expected, the recovered model mass might not even be included in

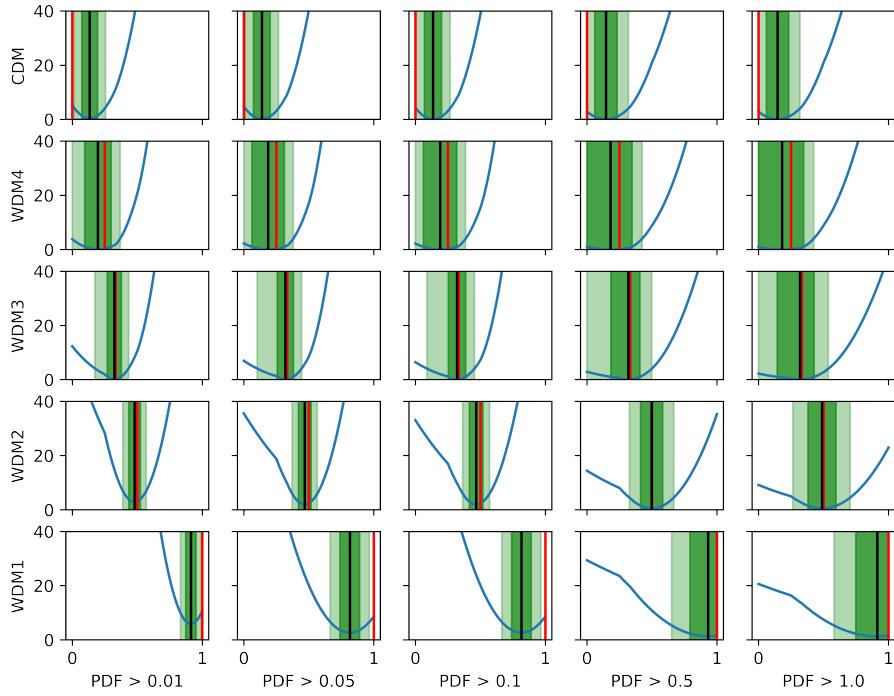


Figure 4.1: Inference results on the NOTRAIN neural netowrks. Each plot corresponds to a different fit combining predictions from each NOTRAIN model with each of the masks applied to the PDF when fitting. The light and dark blue regions correspond to the 1 and 2 sigma confidence regions. The red line corresponds to the true DM model mass, and the black line to the best-fit model that minimises the  $\chi^2$ . The blue curve is the  $\chi^2$  metric. Note that we are using all 5000 sightlines on the inference step.

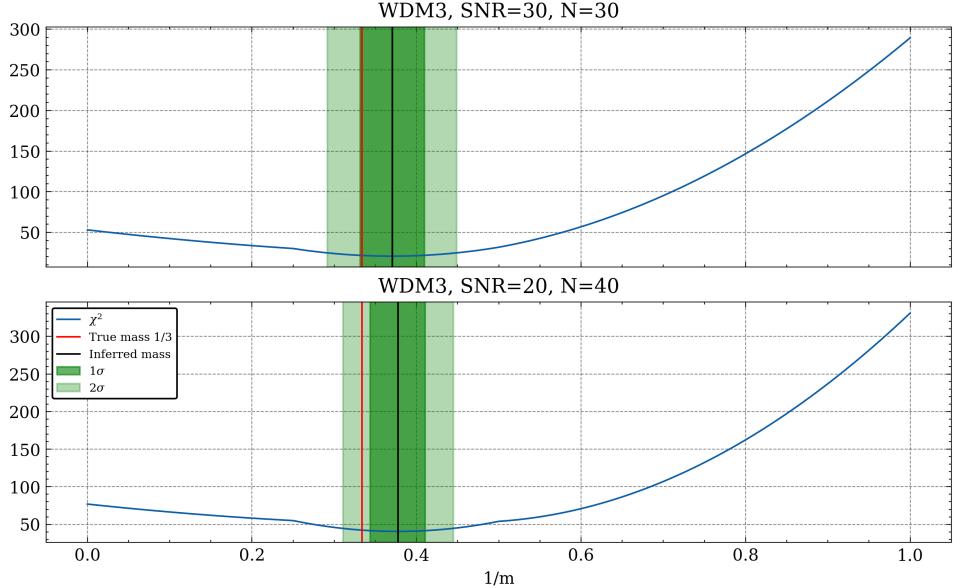


Figure 4.2: Inference predictions on WDM3 for different combinations of SNR and number of targets N.

the  $2\sigma$  regions, meaning that the pipeline fails to correctly recover the true mass if we interpolate models. This does not affect our prediction with real data, since, as we have already mentioned, current WDM constraints favour a lower mass limit of  $\sim 3$  KeV. Lastly, observe in Figure 4.1 that the mask applied on the horizontal axis does not significantly affect the recovered masses.

#### 4.2.2 Realistic UVES observational conditions

In this section, we explore the effect of realistic observational conditions, such as the number of observed quasars, the signal-to-noise ratio (SNR, or the instrumental resolution, in the inferred DM constraints. For that purpose, we use typical parameters for the Ultraviolet and Visual Echelle Spectrograph (UVES) on the European Southern Observatory's Very Large Telescope [50]. We consider a spectral resolution of 6 km/s per pixel, variable SNR in the range  $20 - 30$  and a variable number of targets in the range  $30 - 40$ . Note that the skewers in the **SHERWOOD** dataset are  $20h^{-1}cMpc$  in length, while the spectral range in the UVES instrument expands multiple times that range. In particular, since measurements can extend up to redshift differences  $\Delta z \sim 1$ , we assume that each observed spectrum can be decomposed into  $\sim 15$  of our **SHERWOOD** skewers.

A common compromise in an observational program with a fixed observational time

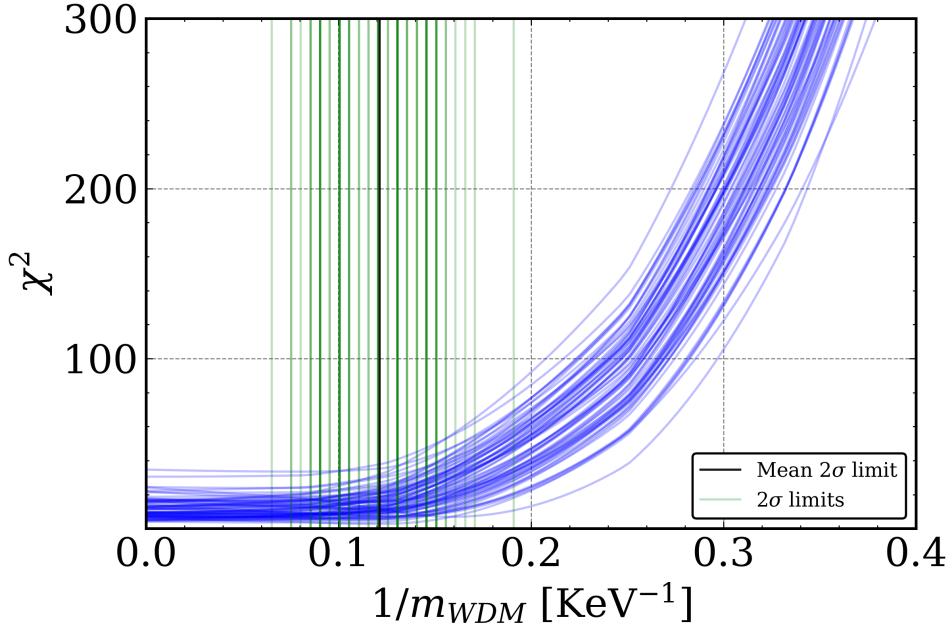


Figure 4.3: The figure shows 100 different  $\chi^2$  fits on 450 SHERWOOD CDM skewers and the  $2\sigma$  constraints distribution as we vary the exact observed draw.

is between number of targets and exposure time per target, which determines the SNR. In Figure 4.2 we show how prioritising SNR or the number of targets affects to inferred WDM masses on the WDM3 model. In general, we find that increasing the target number leads to slightly tighter confidence regions, while increasing the SNR leads to more accurate constraints. Most crucially, observe how the true model mass is, in both cases, recovered within  $2\sigma$ .

We now evaluate the constraining power of the approach developed in this work. For that purpose, we assume CDM to be the true DM model and use our fiducial neural network trained on SHERWOOD. We then draw 450 SHERWOOD skewers, corresponding to 30 observed UVES spectra, post-process them with a resolution of 6 km/s, add random Gaussian noise with  $\text{SNR} = 30$  and use them to run our inference pipeline from section 4. Since the fit depends on the exact draw of “observed” skewers, we repeat this process 100 times with a random draw each time to obtain the  $2\sigma$  limit distribution.

Figure 4.3 shows the distribution of  $2\sigma$  limits and the mean  $2\sigma$  constraint produced by this process. The mean  $2\sigma$  constraint that we report for the inverse mass is  $\sim 0.12 \text{ KeV}^{-1}$ , or  $\sim 8.3 \text{ KeV}$  for the WDM model mass.

To confirm that the fitting process works as expected, we plot in Figure 4.4 the best fit PDF, which corresponds to the CDM model according to Figure 4.3 and an example recovered PDF from a set of 450 observed skewers. Recall that the uncertainties in

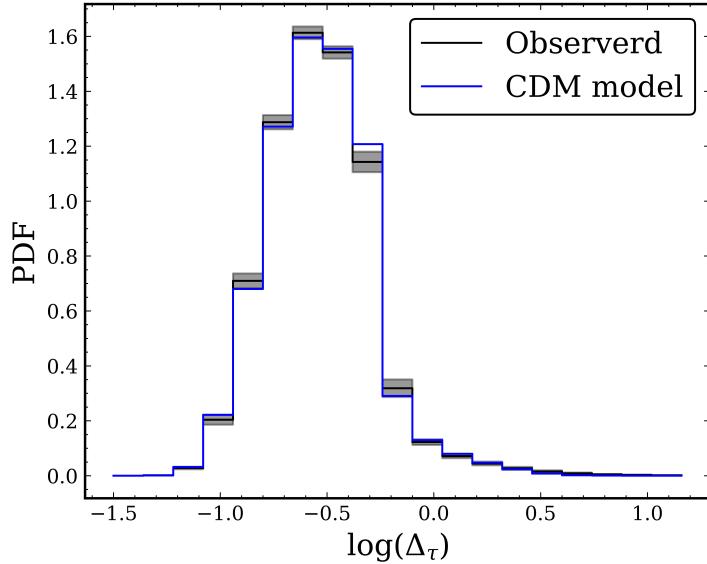


Figure 4.4: An example observed  $\Delta_\tau$  PDF recovered using 450 skewers and its uncertainties in black, plotted against the model CDM PDF, which is the best-fit model in the  $\chi^2$  test.

the recovered PDF include the sample scatter using bootstrapping as well as the machine learning uncertainties, as we have discussed in section 3.4. As expected, the observed PDF is within a  $2\sigma$  distance of the model CDM PDF.

In Figure 4.5 we summarise in black the current  $2\sigma$  state-of-the-art constraint in the literature, using a non-ML approach. In orange, we compare the forecasted constraints from the non-machine learning approach in [39] to our approach in an equivalent dataset to the one used in section 4.2.2. Compared to current limits, our forecasted constraint is a twofold improvement, from  $\sim 4$  to  $\sim 8$  KeV. On the same dataset, we forecast our machine learning technique to also outperform the current pipeline in [39]. As a significant caveat, note that the work in [39] is a joint analysis not only on WDM but also on thermal parameters of the IGM, cosmological parameters, etc. The aforementioned paper encompasses a larger number of parameters with a more complex and refined approach than this work.

### 4.3 Inference on alternative hydrodynamical codes

In this section we test our inference pipeline on Nyx, a different hydrodynamical code. We start by describing in a broad way the differences between the Nyx run used and the Sherwood simulations, and then we use our fiducial neural network trained on the

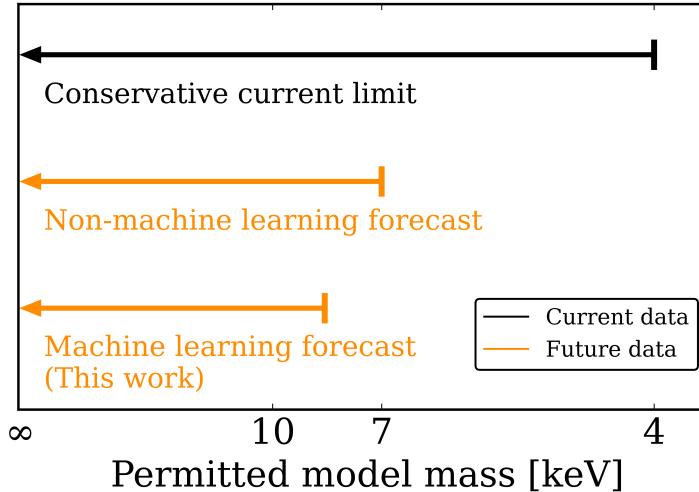


Figure 4.5: Summary  $2\sigma$  relic WDM constraints based on [39] and this work.

SHERWOOD THERMAL suite to recover the Nyx densities and obtain the corresponding constraints.

### 4.3.1 The Nyx code

Nyx [51] is an N-body and gas dynamics code for large-scale cosmological simulations. Nyx uses Adaptive Mesh refinement (AMR) in time and space based on the Eulerian formulation of hydrodynamics, as opposed to the Lagrangian formulation used in the GADGET code employed by the Sherwood simulations. We expect Sherwood and Nyx runs to intrinsically show these non-physical differences related to the different hydrodynamical solvers.

In the Nyx code, dark matter is modeled as discrete Lagrangian particles, allowing the code to follow their evolution under gravity effectively. The evolution of its phase space distribution  $f$  is given by the collisionless Boltzmann equation

$$\frac{\partial f}{\partial t} + \frac{1}{ma^2} \mathbf{p} \cdot \nabla f - m \nabla \phi \cdot \frac{\partial f}{\partial \mathbf{p}} = 0 \quad (4.3)$$

where  $m$  and  $\mathbf{p}$  are mass and momentum and  $\phi$  is the gravitational potential.  $a$  is the scale factor, obtained by using a second-order Runge-Kutta solver. Nyx solves this phase space evolution of  $f$  by sampling its distribution and evolving the particles as an N-body system. The gravitational potential is obtained by solving the Poisson equation

$$\nabla^2 \phi(\mathbf{x}, t) = \frac{4\pi G}{a} (\rho_b + \rho_{dm} - \rho_0) \quad (4.4)$$

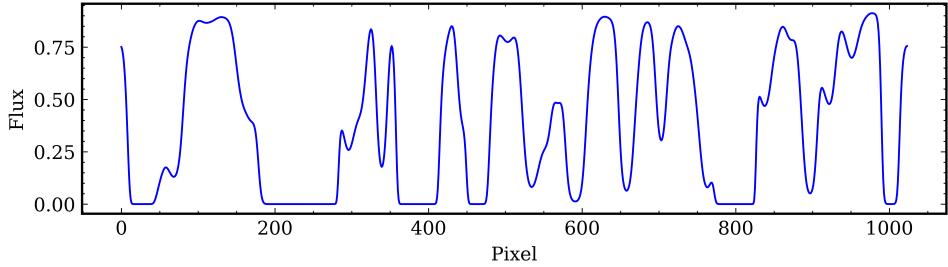


Figure 4.6: A typical Lyman- $\alpha$  skewer obtained from the Nyx runs with  $z_{\text{re}} = 6$  at  $z = 4.4$ .

where  $\rho_0$  is the mean density,  $\rho_b$  the baryonic density and  $\rho_{dm}$  the dark matter density. Dark matter particles are gravitationally coupled to a baryonic fluid, which is treated as an inviscid ideal gas. The gas is described by a state vector  $\mathbf{U} = (\rho_b, a\rho_b U, a^2 \rho_b E, a^2 \rho_b e)$  where  $U$  is the peculiar proper baryonic velocity,  $e$  the internal energy, and  $E$  the total energy. The hydrodynamical equations are approximate by a Riemann solver and can be written in the form

$$\frac{\partial \mathbf{U}}{\partial t} = -\nabla \cdot \mathbf{F} + S_e + S_g + S_{HC}, \quad (4.5)$$

where  $F$  is the flux vector,  $S_g$  the gravity source term,  $S_{HC}$  the heating and cooling term, and  $S_e$  the internal energy flux.

In the rest of this section, we consider 3 Nyx runs at  $z = 4.4$  using CDM and  $20h^{-1}\text{cMpc}$  boxes. The 3 runs different in the different reionization history, and are labelled by the end of reionization redshifts of  $z_{\text{re}} = 6, 7, 8$ . Each skewer has 1024 pixels. In Figure 4.6 show an example Lyman- $\alpha$  skewer for the Nyx run with  $z_{\text{re}} = 6$ .

Since we want to test our inference pipeline, in this test we want to be as agnostic as possible about the nature of our ‘‘observed’’ Nyx spectra. If real data is observed, we would, *a priori*, have no information on the exact thermal history that has led to the observed field. A similar situation occurs with the Nyx runs. Our **SHERWOOD THERMAL** dataset constrains thermal models, but we have *a priori* no guarantee that they match the Nyx runs that we are analysing. In fact, we know that this is not the case. We visually explore the difference between the Nyx and **SHERWOOD THERMAL** runs, we plot the 2D distribution of pixels in the temperature-density plane. We show the result in Figure 4.7, comparing the Nyx CDM run  $z_{\text{re}} = 6$  to the **SHERWOOD THERMAL** CDM runs. The top panel shows the 95% contours in red and black colors. The bottom panel shows the temperature distribution at the mean density value. As can be observed, the Nyx does not fit any of the runs in our training dataset.

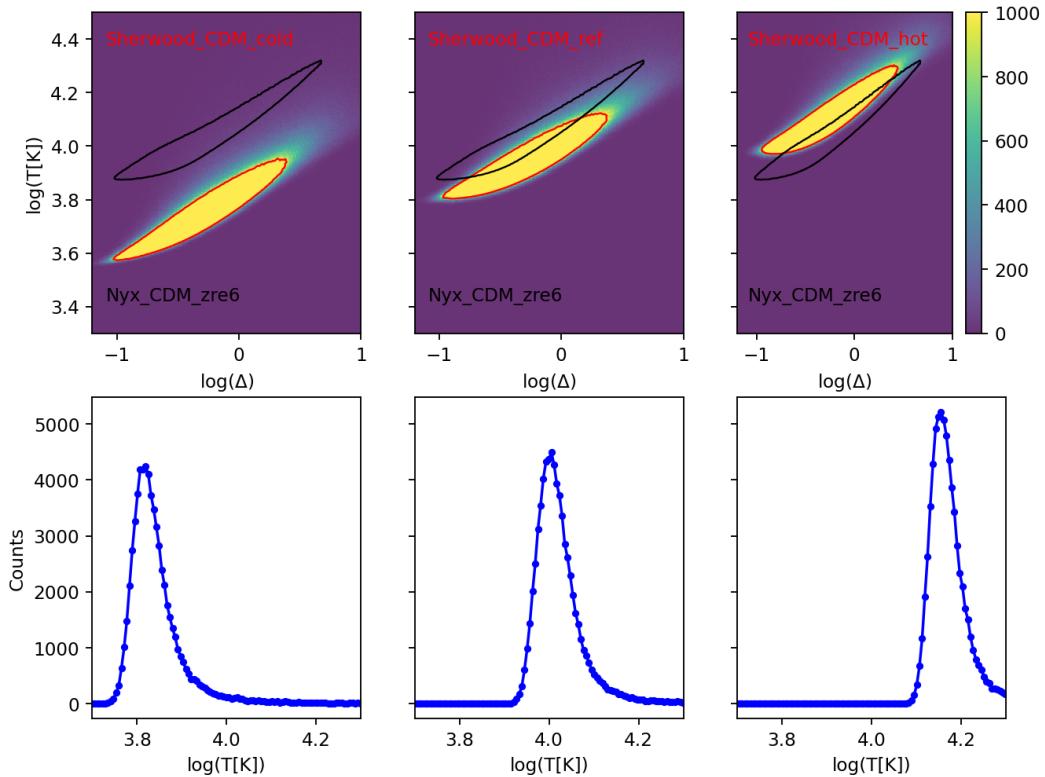


Figure 4.7: The 2D distribution of pixels in the temperature-density plane comparing the Nyx CDM run  $z_{\text{re}} = 6$  to the SHERWOOD THERMAL CDM runs. The top panel shows the 95% contours in red and black colors. The bottom panel shows the temperature distribution at the mean density value. As can be observed, the Nyx does not fit any of the runs in our training dataset.

### 4.3.2 Inference test on Nyx Lyman-alpha skewers

We begin by evaluating the performance of our fiducial neural network with frozen weights trained on **SHERWOOD THERMAL** when predicting of skewers generated by the Nyx code. In Figure 4.8 we show a violin plot with the  $1\sigma$  residues distribution as defined in Equation 3.24. Values in the range  $[0, 1]$  correspond to pixels that have been successfully recovered within a  $1\sigma$  accuracy. Observe how the models with earlier reionization ( $z_{\text{re}} = 7$  and  $z_{\text{re}} = 8$ ), which have lower temperatures, have a higher recovery rate. This is likely due to the fact that our training data set contains more WDM models close to CDM, and we know that low WDM masses and low temperature have a degenerate effect on the Lyman- $\alpha$  forest. In general, we note that the  $\geq 75\%$  of the pixels are correctly recovered. Even if the performance is slightly degraded compared to the **SHERWOOD THERMAL** validation, as expected, this is a strong indication that the neural network has learnt the relevant physical relations, and not potential simulation-specific correlations.

We then run the same inference pipeline for the  $z_{\text{re}} = 6$  model as we did in section 4.2.2 with a single major difference. Since we know that the Nyx runs do not fit any of our thermal models (and the same circumstance will occur when using real observations), we need to utilise the fiducial neural network trained on multiple thermal modes, **SHERWOOD THERMAL**. Additionally, when fitting the recovered  $\Delta_\tau$  PDF to each model PDF, we will have 3 different  $\chi^2$  curves for each one of the thermal models {ref, hot, cold}. To avoid a joint optimization problem and constraining the thermal history on top of the WDM mass (which is out of the scope of this work, and also unrealisitic since we are not using a fine thermal model grid), we will select the thermal model that minimises the  $\chi^2$  value, that is

$$\min_{t,m} \chi^2(t,m), \quad (4.6)$$

where  $t \in \{\text{ref}, \text{hot}, \text{cold}\}$  and  $m$  labels the continuously interpolatted inverse warm dark matter mass. We compute  $\chi^2(t, m)$  and find that the ref and hot model produce similar  $\chi^2$  values, while the cold model has  $\chi^2 \sim 1000$ , as expected from Figure 4.7. Both the ref and hot models produce a similar  $2\sigma$  lower bound on the WDM mass:  $m_{\text{WDM}} \gtrsim 10$  KeV. Compared to section 4.2.2, the results are fairly similar, showcasing the robustness of our approach.

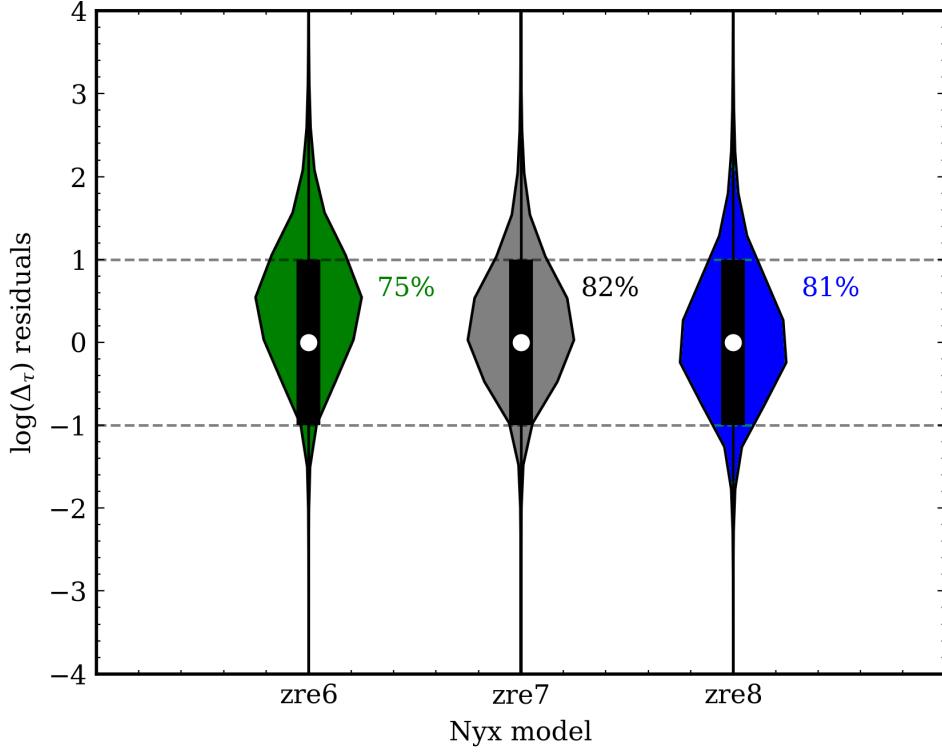


Figure 4.8: iolin plot with the  $1\sigma$  residues distribution as defined in Equation 3.24. Values in the range  $[0, 1]$  correspond to pixels that have been successfully recovered within a  $1\sigma$  accuracy. Observe how the models with earlier reionization ( $z_{\text{re}} = 7$  and  $z_{\text{re}} = 8$ ), which have lower temperatures, have a higher recovery rate. This is likely due to the fact that our training data set contains more WDM models close to CDM, and we know that loss WDM masses and low temperate have a degenerate effect on the Lyman- $\alpha$  forest. In general, we note that the  $\geq 75\%$  of the pixels are correctly recovered.

Table 4.1: List of the SQUAD DR1 sightlines used, see [50] for the reduction details, together with their emission redshift and the average continuum SNR. All sightlines are  $20\text{h}^{-1}\text{cMpc}$  and centered at  $z = 4.4$ .

SQUAD DR1 name	$z_{\text{em.}}$	SNR
J004054	4.976	33
J021043	4.65	25
J025019	4.77	12
J030722	4.728	50
J033829	5.032	14
J145147	4.763	100

## 4.4 WDM constraints from SQUAD DR1 observational data

We begin applying our density recovery and WDM mass inference pipeline to a set of 6 observed quasar sightlines from the SQUAD DR1 survey [50]. Since we are working with observational data, we will always train the model with the complete **SHERWOOD THERMAL** dataset that includes varied thermal histories and WDM masses. Note that since we are not trying to constrain the thermal history (or other parameters that can affect the Lyman- $\alpha$  forest), we should ideally use a training set that includes as much variation as possible to make sure the neural network can perform in a scenario where we ignore the true thermal history.

Our SQUAD DR1 data consists of 6 Lyman- $\alpha$  sightlines of size  $20\text{h}^{-1}\text{cMpc}$  with varied SNR (see table 4.1), observed with the Ultraviolet and Visual Echelle Spectrograph (UVES) on the European Southern Observatory's Very Large Telescope, which has an average resolution of  $\text{FWHM} \approx 6\text{km s}^{-1}$ . We consider sightlines centred at  $z = 4.4$  for this specific application. Since each quasar has its own noise level, we retrain the same fiducial architecture with the corresponding noise level before the prediction step.

We are then set to predict the recovered density fields for each sightline. Figure 4.9 shows all our SQUAD DR1 skewers together with the recovered  $\Delta_\tau$  field.

We now compute the  $\chi^2(t, m)$  for all 3 thermal models, which are minimised for the respective CDM run as expected, and find that the thermal model producing a minimal  $\chi^2$  is the cold. In Figure 4.10 we show, in the left panel, all 3  $\chi^2$  curves as a function of the WDM mass. The right panel shows the recovered  $\Delta_\tau$  PDF from the SQUAD DR1 sample together with the best-fit model, corresponding to the CDM cold **SHERWOOD**

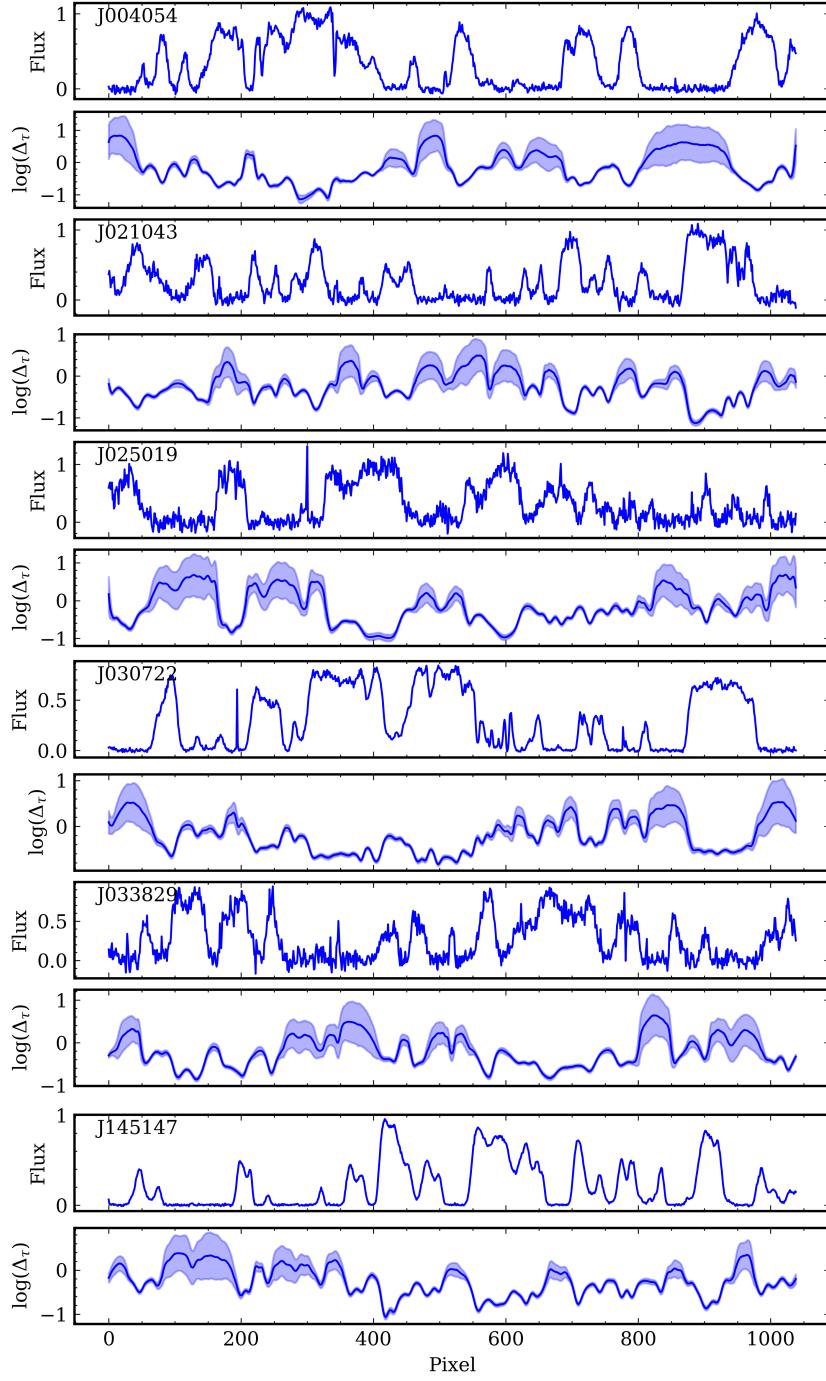


Figure 4.9: All 6 sightlines from the SQUAD DR1 sample in Table 4.1. As can be seen, the noise levels vary, depending on the exposure time to the target. All sightlines are  $20\text{h}^{-1}\text{cMpc}$  and centered at  $z = 4.4$ . We show the recovered density field by our fiducial architecture trained on `SHERWOOD THERMAL` and retrained with the noise specifications of each target.

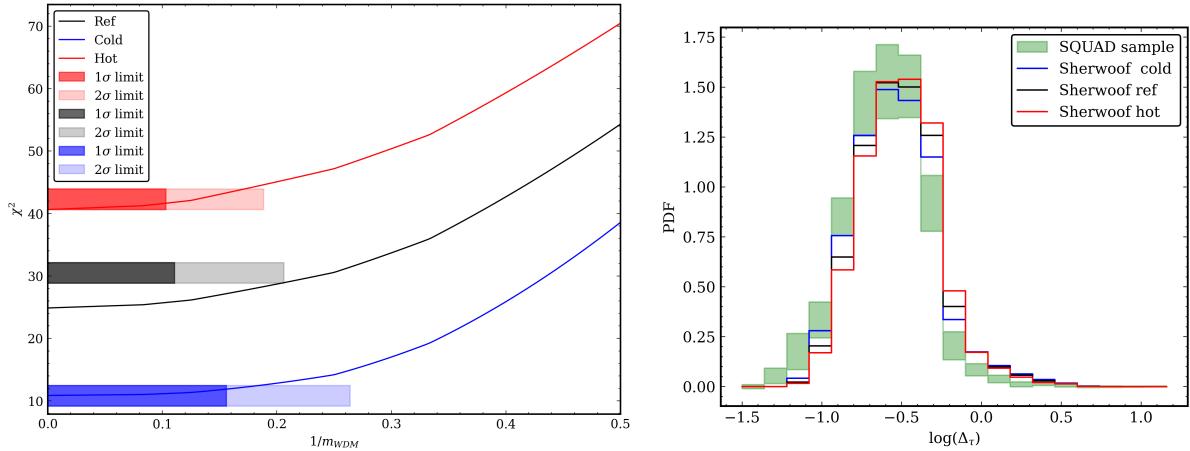


Figure 4.10: In the left panel, we show all 3  $\chi^2$  curves as a function of the WDM mass, together with the  $1, 2\sigma$  confidence regions. The right panel shows the recovered  $\Delta_\tau$  PDF from the SQUAD DR1 sample with the symmetric uncertainty envelop, together with the best-fit model, corresponding to the CDM cold SHERWOOD THERMAL model.

THERMAL model. Using the best-fit thermal model, we find a lower bound on the WDM mass of  $m_{WDM} \gtrsim 3.8$  KeV.

## 4.5 WDM constraints from GHOST observed spectrum

We also consider a Lyman- $\alpha$  skewer obtained from the GHOST instrument, which corresponds to the ultra-luminous quasar J0306+1853 [52] with emission redshift  $z = 5.363$ . For this spectrum, we have a continuum reconstruction in the range [971, 1210] Å in the emission rest frame, with an average signal-to-noise ratio of  $\text{SNR} \approx 150$ . We extract skewers of length  $20\text{h}^{-1}\text{cMpc}$  and consider them independent in order to run the neural network predictions. In total, we obtained 13 such sightlines and discarded the range  $\approx [1120, 1162]$  Å that contains a DLA.

Figure 4.12 shows an example  $20\text{h}^{-1}\text{cMpc}$  portion of the spectrum together with the recovered density field. We split the original spectrum into 13 such skewers of the same length as the ones of the Sherwood dataset, and consider them independent.

We apply our WDM inference pipeline to the 13 segments of the J0306+1853 spectrum. The  $\chi^2$  is minimised for the CDM model within each thermal history, and the best-fit thermal model is the SHERWOOD THERMAL CDM reference run. The corresponding

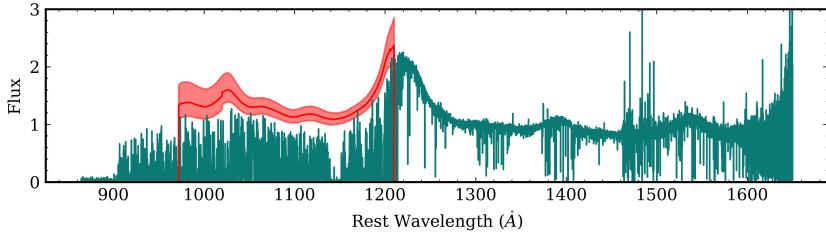


Figure 4.11: THE GHOST spectrum for J0306+1853 [52]. The red curve shows the reconstructed continuum together with  $1\sigma$  uncertainties, obtained with a PCA technique based on the spectrum to the right side of the Lyman- $\alpha$  line. Observe the DLA at  $\sim 1150$  Å, which we mask when analysing the Lyman- $\alpha$  forest.

$2\sigma$  constraint on the WDM mass is  $m_{\text{WDM}} \gtrsim 4.4$  KeV at  $2\sigma$  confidence.

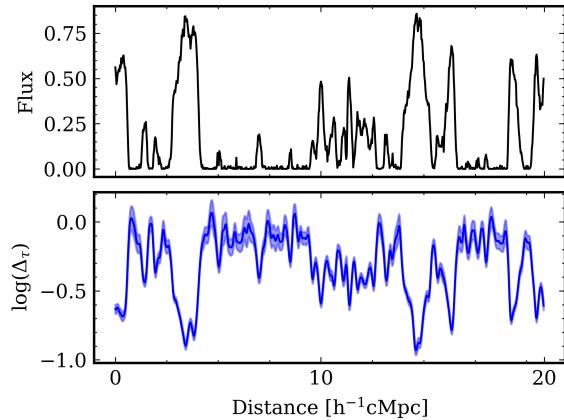


Figure 4.12: An example  $20h^{-1}\text{cMpc}$  portion of the J0306+1853 spectrum together with the recovered density field. We split the original spectrum into 13 such skewers of the same length as the ones of the SHERWOOD THERMAL dataset, and consider them independent.

In Table 4.2 we list the current state-of-the-art  $2\sigma$  lower bounds on  $m_{\text{WDM}}$  thermal relics constraints in the literature obtained from the Lyman- $\alpha$  power spectrum. The previous efforts are based on a Bayesian inference framework to compare the observed power spectrum with the one obtained from simulated data. In contrast, in our work, we do the inference directly on the non-observable density field level. The resulting bounds produced in our work are comparable to previous efforts but with the advantage of requiring substantially less observational data. For reference, in [39], the constraints are obtained from 15 spectra measured along  $327h^{-1}\text{cMpc}$ , while our GHOST data consists of 13  $20h^{-1}\text{cMpc}$  skewers. The tightening of the constraint from the SQUAD DR1 sample to the GHOST sample is also expected since the latter is a larger sample

with lower noise.

Table 4.2: List of current state-of-the-art  $2\sigma$  lower bounds on  $m_{\text{WDM}}$  thermal relics constraints in the literature obtained from the Lyman- $\alpha$  power spectrum. We compare them to the results of this work, obtained doing inference directly at the density field level recovered by our Bayesian neural network.

Source	$2\sigma m_{\text{WDM}}$ lower bound [KeV]
Iršič et al. (2024) [39]	4.1
Villasenor et al. (2023) [6]	3.1
This work (SQUAD DR1 sample)	3.8
This work (GHOST spectrum)	4.4

## 4.6 Comparison of the inference pipeline against Information Maximising Neural Networks

The inference pipeline presented so far in this section is based on a simple  $\chi^2$  fit of the  $\Delta_\tau$  recovered PDFs from our fiducial NN model. This pipeline relies on a series of contingent choices, most notably the use of the density PDFs as the summary statistics of the fields. In this section, we explore, in an agnostic way, the possibility of using other summaries different from the  $\Delta_\tau$  PDFs to perform the inference (note that other summaries such as the density power spectrum, curvature,... could potentially be used). More concretely, we will introduce Information Maximising Neural Networks (IMNNs) and use them to perform the inference within a Bayesian framework. We then compare the results of this procedure with the inference pipeline discussed in section 4.

### 4.6.1 Information Maximising Neural Networks

Information Maximising Neural Networks aim at obtaining optimal summaries of data [53]. Neural networks are used to parametrise these summaries in an agnostic way by maximising the information of the summaries with respect to the model parameters of interest.

Consider a data-generating procedure depending on some model parameters  $\theta$ , generating data realization  $d_i(\theta)$  where  $i$  labels a realization or initial seed of the simulation. We want to obtain a function  $f: d \mapsto x$  that maps each simulation to a summary vector

of the same size as  $\theta$ . This is, essentially, a compression algorithm. IMNNs work by transforming the original likelihood of the data, which is a priori not known, into the gaussian form

$$-2 \ln \mathcal{L}(x|\theta) = (x - \mu(\theta))^T C^{-1} (x - \mu(\theta)), \quad (4.7)$$

where  $C$  is the covariance matrix of the calculated summaries with a set of  $n_s$  simulations, and  $\mu$  the summary mean depending on the model parameters. The information of the observed summaries with respect to  $\theta$  is then the Fisher information matrix [54]:

$$F_{\alpha\beta} = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta_\alpha \partial \theta_\beta} \log \mathcal{L}(x; \theta) \middle| \theta \right] = \frac{\partial \mu}{\partial \theta_\alpha} C^{-1} \frac{\partial \mu}{\partial \theta_\beta}, \quad (4.8)$$

whose determinant we note as  $|F|$ . The goal is to obtain summaries that maximise the Fisher information while maintaining a minimum covariance condition to generate independent summaries. The summaries produced by the network can then be used to perform inference on them. Since the Fisher information is a quantity that depends on the model parameters  $\theta$ , the quantity in Equation 4.8 needs to be evaluated at some fiducial model parameters in order to obtain a numerical result.

IMNNs have been successfully leveraged in the IGM community. Recent papers have explored the possibility of using them to perform IGM thermal parameter inference from Lyman- $\alpha$  skewers, see [55] for instance, where authors find IMNNs to yield tighter and more robust constraints than classical Markov Chain Monte Carlo approaches. Despite these promising results, many challenges arise when using IMNNs on real data, primarily related to the correct identification and interpretation of model parameters.

### 4.6.2 IMNN training and non-linear summaries

In this section we consider a simple MLP architecture with linear layers followed by PReLU( $\alpha$ ) activation functions and a dropout layer that randomly (with probability  $p$ ) sets to 0 the any layer weight during each epoch to prevent over-fitting. The network takes as input a simulated data vector  $d$  and produces a summary vector  $x$  of the same size and the parameter vector  $\theta$ . We use the Adam optimiser to maximise  $|F|$  by minimising the following loss function

$$\mathcal{L}_{IMNN} = -\log(|F|) + \lambda \frac{\mathcal{N}}{\mathcal{N} + \exp(-\mathcal{N})} \mathcal{N}, \quad (4.9)$$

where  $\mathcal{N} = ||C - I|| + |C^{-1} - I||$  measure the deviation from independt summaries and  $\lambda$  is a coupling constant. In Equation 4.9, the second term sets a scale for the Fisher information by producing summaries whose covariance approaches the identity matrix. Once this is achieved, the term containing the exponential factor vanishes and the network will maximise  $|F|$ . Note that there is not a unique set of potential optimal sumamries. In fact, any bijective function of a sufficient statistic for a certian likelihood is also a sufficient statistic.

For each paramter update in the training procedure, we generate a batch of data at the fidicual parameters  $\theta_f$ . The derivatives in Equation 4.8 are numerically approximated with finite differences by running simulations at parameters  $\theta_f \pm \Delta\theta_\alpha$ , where  $\Delta\theta_\alpha$  are a small parameter variation, and then calcualting

$$\frac{\partial x}{\partial \theta_\alpha} = \frac{x(d(\theta_f + \Delta\theta_\alpha)) - x(d(\theta_f - \Delta\theta_\alpha))}{2\Delta\theta_\alpha}. \quad (4.10)$$

We then calculate  $C$ , the covariance of the summaries at fiducal paramters, and use it to compute the Fisher information in Equation 4.8 and the loss function in Equation 4.9. Note that, since the covariance matrix and the derivative in the Fisher information matrix are computed at the data summaries, they implicitye depend on the NN parameters.

### 4.6.3 Summarising a Gaussian signal

We implement IMNNs using Pytorch<sup>1</sup>, a deep-learning Python framework. We test the implementation first by exploring its behavior on a toy model, where we generate ramdon samples from a Gaussian distribution  $\mathcal{N}(\mu, \sigma)$ . The sufficient statistic for the model parameters  $\theta = (\mu, \sigma)$  are, in this case, the sample mean and standard deviation:

$$\hat{\mu} = \frac{1}{n_d} \sum d_i \quad \hat{\sigma}^2 = \frac{1}{n_d - 1} \sum (d_i - \hat{\mu})^2. \quad (4.11)$$

Note that the statistic for  $\sigma$  is non-linear. For this example, we select fiducial parameters  $\theta_f = (\mu = 0, \sigma = 1)$  and  $\Delta\theta = (0.1, 0.1)$  and generate random fields with 100 pixels. In total, 5000 fields are generated for each parameter set, including a validation dataset. Note that testing the network performance in the validation dataset is crucial in performing early-stopping during the training process. Indeed, since the training dataset is limited, it will contain spurious correlations that the network will use to infer a higher

---

<sup>1</sup><https://pytorch.org>

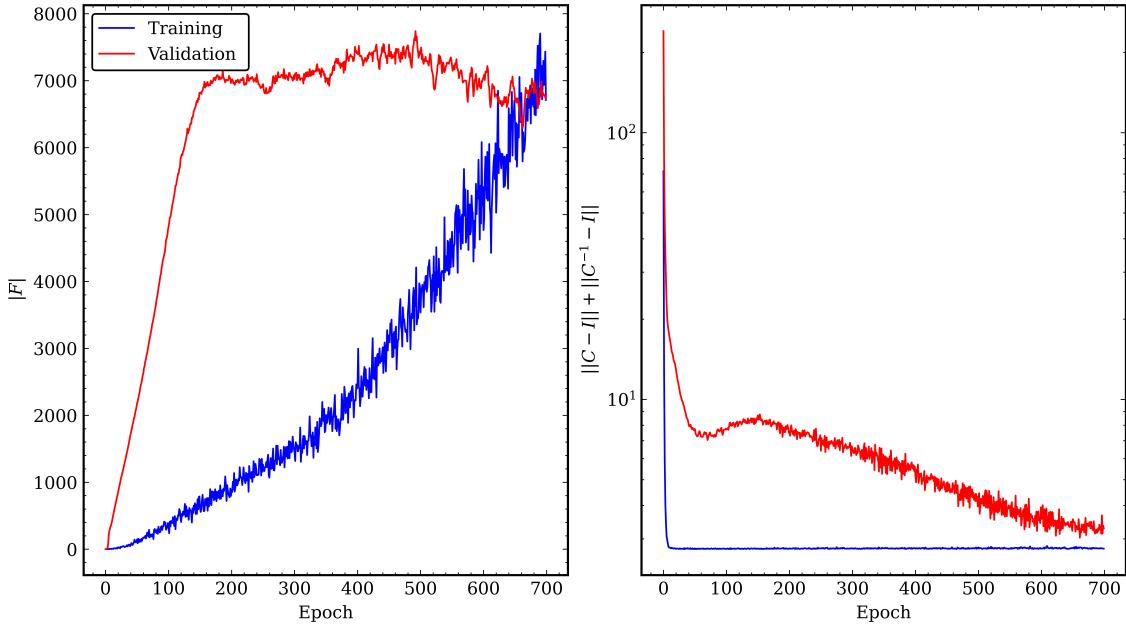


Figure 4.13:  $|F|$  and  $\|C - I\| + \|C + I\|$  as a function of the epoch for the training and validation sets during the training of the IMNN on a normal field. The goal is to find summaries to optimally extract information about the mean and variance of the field.

information than expected. By stopping the training when the network information on the validation set saturates, we can avoid this problem. We use a simple architecture, with layers  $[128, 128, 128, 2]$ , learning rate of 0.001, dropout rate of  $p = 0.5$ , and batch size of 500. Observe the training evolution in Figure 4.13, where we show  $|F|$  and  $\|C - I\| + \|C + I\|$  as a function of the epoch for the training and validation sets. As can be seen, the validation information quickly saturates in  $\sim 100$  epochs, and then slowly decreases as the network over-fits.

To better interpret the network output and to understand its behavior, we generate samples of the same size with the zero mean but a standard deviation randomly sampled from  $(0, 12)$ . We then compute the exact statistic (the sample standard deviation) and plot it against the second IMNN summary output. The result is Figure 4.14. Observe that the exact statistic in the  $x$ -axis is highly correlated to the network summary in the  $y$ -axis. Since the relation between the two quantities is clearly bijective, the model has successfully learnt to extract all the possible information for the field covariance. The natural scatter is due to having a simple NN model.

Note that the network has not seen any normal field with such variances  $\sigma \in (0, 12)$  during training, yet it is able to extract the correct summary. This is of crucial impor-

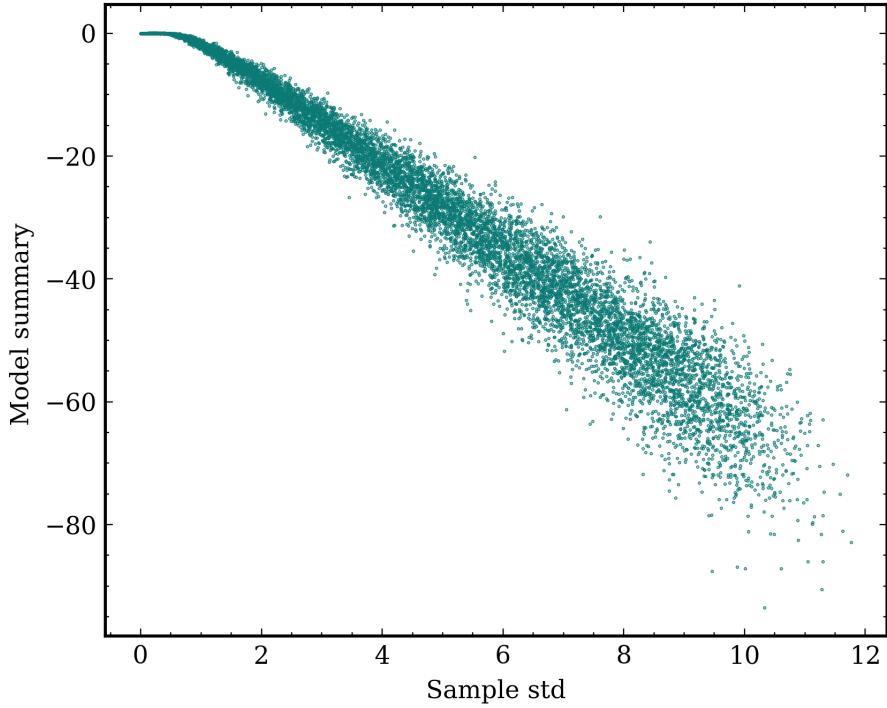


Figure 4.14: The IMNN summary plotted against the exact sufficient statistic for the standard deviation using multiple samples with  $\sigma \in (0, 12)$ .

tance, since it means that we could use the IMNN summaries to do inference on a field with parameter values slightly different from the fiducial ones used during training.

To conclude the exploration of this toy model, we use the network to perform Bayesian inference. We implement a simple Approximate Bayesian Computation (ABC) [56] algorithm that obtains approximate posterior samples from a given set of prior samples and observed data. As the observed data, we take 50 Gaussian fields simulated at fiducial parameter ( $\mu = 0, \sigma = 1$ ) values. As priors, we take 5000 samples from a non-informative uniform distribution in  $(-5, 5)$  for  $\mu$  and  $(0, 10)$  for  $\sigma$ . The ABC rejection algorithm is described in Algorithm 3.

In Figure 4.15 we show the posterior samples and Gaussian Kernel Density Estimation (KDE) for the distributions of  $\mu$  and  $\sigma$ . The dashed vertical lines show the true parameter values. As expected, and even with a non-informative prior, the IMNN summary contain sufficient information to produce tight posterior around the true model parameters. Note that the posterior scatter on the non-linear summary  $\sigma$  is larger than on the linear summary  $\mu$ .

---

**Algorithm 3** Approximate Bayesian Computation Rejection Algorithm

---

```

1: Input: Observed data  $\mathbf{y}$ , threshold  $\epsilon$ , number of simulations  $N$ , prior distribution  $\pi(\theta)$ 
2: Output: Accepted parameter values  $\{\theta_i\}_{i=1}^M$ 
3: Initialize  $M \leftarrow 0$ 
4: for  $i = 1$  to  $N$  do
5:   Sample  $\theta^*$  from the prior distribution  $\pi(\theta)$ 
6:   Simulate data  $\mathbf{y}^*$  from the model using  $\theta^*$ 
7:   if  $d(\mathbf{y}, \mathbf{y}^*) \leq \epsilon$  then
8:     Accept  $\theta^*$ :  $\theta_{M+1} \leftarrow \theta^*$ 
9:     Increment  $M \leftarrow M + 1$ 
10:  end if
11: end for
12: return  $\{\theta_i\}_{i=1}^M$ 

```

---

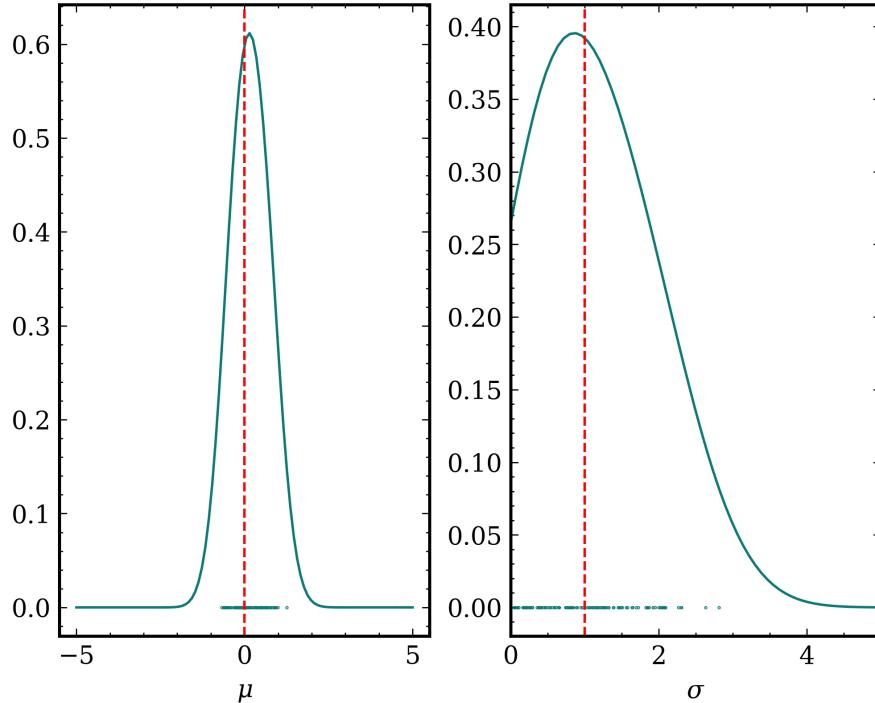


Figure 4.15: Posterior samples and KDE for the ABC rejection algorithm applied to the normal toy model, where we infer the mean and variance of a Gaussian field with flat priors and the summaries output of a IMNN. The dashed vertical lines show the true parameter values.

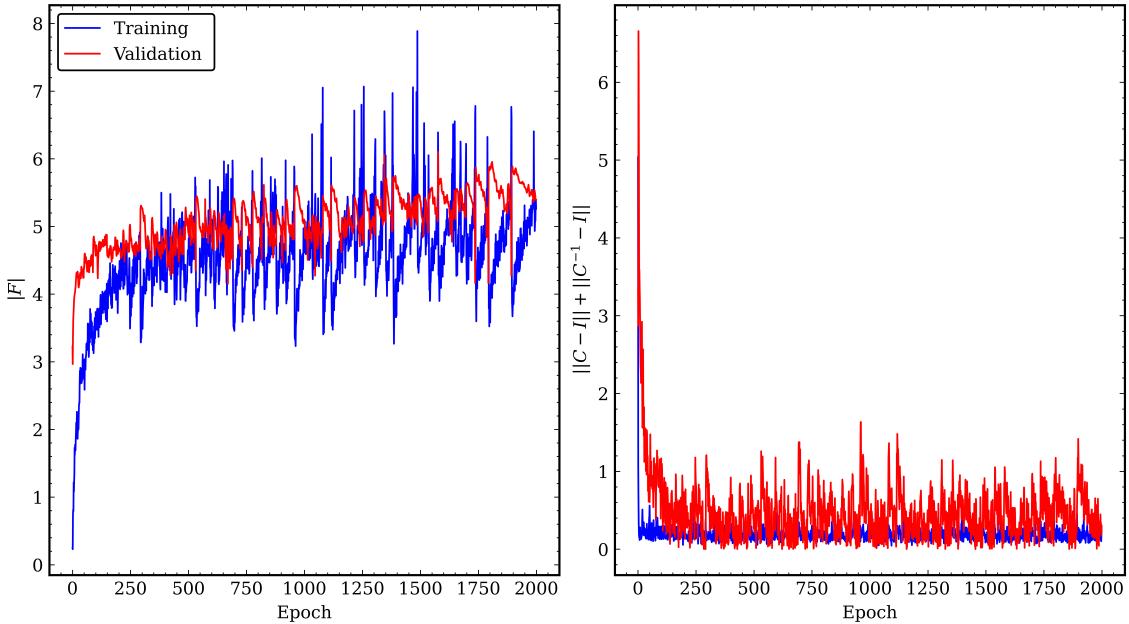


Figure 4.16:  $|F|$  and  $\|C - I\| + \|C^{-1} - I\|$  as a function of the epoch for the training and validation sets during the training of the IMNN on the **SHERWOOD** dataset Lyman- $\alpha$  skewers.

#### 4.6.4 IMNN inference results on WDM masses

We can now deploy a simple IMNN as an alternative way of constraining WDM models. We follow section 4.6.3 and consider a similar architecture but now with 4 dense layers of size [512, 512, 256, 2]. As input to the NN, we consider Lyman- $\alpha$  flux skewers. Since in flux space the skewers has many simulation-specific and prominent features that can be picked-up by a NN, we work in Fourier space. More precisely, the input to the network are

$$\sqrt{k}|\delta_F(k)|, \quad (4.12)$$

where  $\delta_F$  is the flux contrast of the skewer. We use the **SHERWOOD** simulation suite with varied WDM mass to train the IMNN. Note that this means that we are assuming that WDM is the only model parameter affecting the Lyman- $\alpha$  forest property. We ignore thermal parameters variations for this demonstration. We train our model on the fiducial CDM mass corresponding to  $0 \text{ Kev}^{-1}$ , and use the WDM3 model corresponding to  $0 \text{ Kev}^{-1}$  to calculate the summary derivatives. The choice of WDM3 is due to the flux skewers showing sufficient variation with respect to CDM.

In Figure 4.16 we show the training progress of the IMNN as a function of the epoch. The information extracted on the validation split quickly saturates at  $\sim 250$  epochs. It

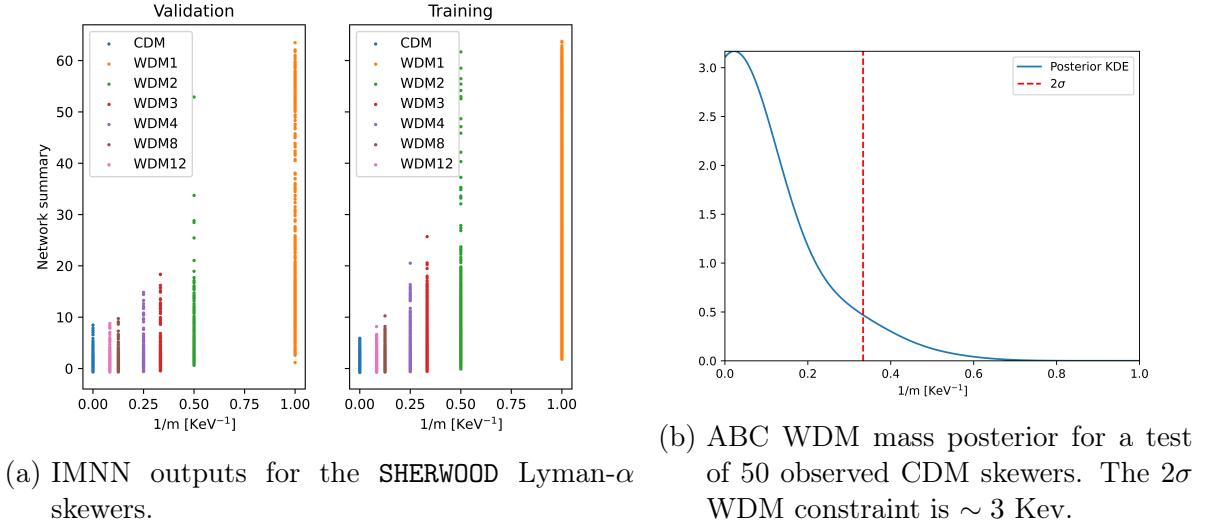


Figure 4.17

is clear that the network is able to learn a map from Lyman- $\alpha$  skewers in Fourier space into a one-dimensional parameter space. Since the SHERWOOD suite has a fix number of simulations, interpreting the network output summaries is a challenging task.

In Figure 4.17a) we show the summaries of the trained IMNN for all the aviable SHERWOOD Lyman- $\alpha$  skewers is Fourier space. Observe how the summaries show a large scatter for every model, corresponding to the large simulation variability within each skewer. However, the mean summaries show a clear bijective trend, manifesting that the IMNN has learnt an informative summary. Note that this can be interpreted as the necessity of a large number of observed samples when constraining WDM models. We use the IMNN summaries to perform an inference test and obtain a Bayesian posterior as follows. First, we select 50 CDM skewers from the validation dataset and obtain their corresponding summaries by passing them through the IMNN. We now use all validation skewers from the SHERWOOD suite and obtain their summaries. We use the ABC rejection algorithm 3 to generate to posterior distribution in Figure 4.17b). The  $2\sigma$  limit for the WDM mass is  $\sim 3$  KeV. Recall that this is a comparable constraint to the one obtained in section 4.4. We interpret this not only as a robustness sign of our original pipeline involving a  $\chi^2$  fit the  $\Delta_\tau$  PDFs, but most notably as a sign that it optimally extracts the majority of the information of the Lyman- $\alpha$  skewers with respect to the WDM mass parameter.

## 5 Conclusions

# Bibliography

- [1] Sean M. Carroll. *Spacetime and Geometry: An Introduction to General Relativity*. Cambridge University Press, 2019.
- [2] Planck Collaboration, P. A. R. Ade, N. Aghanim, C. Armitage-Caplan, M. Arnaud, et al. “Planck 2013 results. XVI. Cosmological parameters”. In: *APP* 571, A16 (Nov. 2014), A16. DOI: [10.1051/0004-6361/201321591](https://doi.org/10.1051/0004-6361/201321591). arXiv: [1303.5076](https://arxiv.org/abs/1303.5076) [astro-ph.CO].
- [3] James E. Gunn and Bruce A. Peterson. “On the Density of Neutral Hydrogen in Intergalactic Space.” In: *APJ* 142 (Nov. 1965), pp. 1633–1636. DOI: [10.1086/148444](https://doi.org/10.1086/148444).
- [4] Robert H. Becker, Xiaohui Fan, Richard L. White, Michael A. Strauss, Vijay K. Narayanan, et al. “Evidence for Reionization at  $z \sim 6$ : Detection of a Gunn-Peterson Trough in a  $z = 6.28$  Quasar”. In: *The Astronomical Journal* 122.6 (Dec. 2001), pp. 2850–2857. ISSN: 0004-6256. DOI: [10.1086/324231](https://doi.org/10.1086/324231). URL: <http://dx.doi.org/10.1086/324231>.
- [5] Ian D. McGreer, Andrei Mesinger, and Valentina D’Odorico. “Model-independent evidence in favour of an end to reionization by  $z \approx 6$ ”. In: *MNRAS* 447.1 (Feb. 2015), pp. 499–505. DOI: [10.1093/mnras/stu2449](https://doi.org/10.1093/mnras/stu2449). arXiv: [1411.5375](https://arxiv.org/abs/1411.5375) [astro-ph.CO].
- [6] Bruno Villasenor, Brant Robertson, Piero Madau, and Evan Schneider. “New constraints on warm dark matter from the Lyman-alpha forest power spectrum”. In: *Physical Review D* 108.2 (July 2023). ISSN: 2470-0029. DOI: [10.1103/physrevd.108.023502](https://doi.org/10.1103/physrevd.108.023502). URL: <http://dx.doi.org/10.1103/PhysRevD.108.023502>.
- [7] Elisa Boera, George D. Becker, James S. Bolton, and Fahad Nasir. “Revealing Reionization with the Thermal History of the Intergalactic Medium: New Constraints from the Ly-alpha Flux Power Spectrum”. In: *The Astrophysical Journal* 872.1 (Feb. 2019), p. 101. ISSN: 1538-4357. DOI: [10.3847/1538-4357/aafee4](https://doi.org/10.3847/1538-4357/aafee4). URL: <http://dx.doi.org/10.3847/1538-4357/aafee4>.

## Bibliography

---

- [8] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *Nature* 521.7553 (May 2015), pp. 436–444. ISSN: 1476-4687. DOI: 10.1038/nature14539. URL: <http://dx.doi.org/10.1038/nature14539>.
- [9] Xia Zhao, Limin Wang, Yufei Zhang, Xuming Han, Muhammet Deveci, and Milan Parmar. “A review of convolutional neural networks in computer vision”. In: *Artificial Intelligence Review* 57.4 (Mar. 2024). ISSN: 1573-7462. DOI: 10.1007/s10462-024-10721-6. URL: <http://dx.doi.org/10.1007/s10462-024-10721-6>.
- [10] Greg Van Houdt, Carlos Mosquera, and Gonzalo Nápoles. “A review on the long short-term memory model”. In: *Artificial Intelligence Review* 53.8 (May 2020), pp. 5929–5955. ISSN: 1573-7462. DOI: 10.1007/s10462-020-09838-1. URL: <http://dx.doi.org/10.1007/s10462-020-09838-1>.
- [11] Sejun Park, Chulhee Yun, Jaeho Lee, and Jinwoo Shin. “Minimum Width for Universal Approximation”. In: *International Conference on Learning Representations*. 2021. URL: <https://openreview.net/forum?id=0-XJwyoIF-k>.
- [12] R.R. Schaller. “Moore’s law: past, present and future”. In: *IEEE Spectrum* 34.6 (1997), pp. 52–59. DOI: 10.1109/6.591665.
- [13] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, et al. *LLaMA: Open and Efficient Foundation Language Models*. 2023. arXiv: 2302.13971 [cs.CL].
- [14] Yu Emma Wang, Gu-Yeon Wei, and David Brooks. *Benchmarking TPU, GPU, and CPU Platforms for Deep Learning*. 2019. arXiv: 1907.10701 [cs.LG].
- [15] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. “Learning representations by back-propagating errors”. In: *Nature* 323.6088 (Oct. 1986), pp. 533–536. ISSN: 1476-4687. DOI: 10.1038/323533a0. URL: <http://dx.doi.org/10.1038/323533a0>.
- [16] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. 2017. arXiv: 1412.6980 [cs.LG].
- [17] A. Vallenari, A. G. A. Brown, T. Prusti, J. H. J. de Bruijne, F. Arenou, et al. “GaiaData Release 3: Summary of the content and survey properties”. In: *Astronomy and Astrophysics* 674 (June 2023), A1. ISSN: 1432-0746. DOI: 10.1051/0004-6361/202243940. URL: <http://dx.doi.org/10.1051/0004-6361/202243940>.

- 
- [18] Dean Richardson, Robert L. Jenkins III, John Wright, and Larry Maddox. “ABSOLUTE-MAGNITUDE DISTRIBUTIONS OF SUPERNOVAE”. In: *The Astronomical Journal* 147.5 (Apr. 2014), p. 118. ISSN: 1538-3881. DOI: 10.1088/0004-6256/147/5/118. URL: <http://dx.doi.org/10.1088/0004-6256/147/5/118>.
  - [19] Romina Ahumada, Carlos Allende Prieto, Andrés Almeida, Friedrich Anders, Scott F. Anderson, et al. “The 16th Data Release of the Sloan Digital Sky Surveys: First Release from the APOGEE-2 Southern Survey and Full Release of eBOSS Spectra”. In: *APJS* 249.1, 3 (July 2020), p. 3. DOI: 10.3847/1538-4365/ab929e. arXiv: 1912.02905 [astro-ph.GA].
  - [20] Jordi Miralda-Escude. “Reionization of the Intergalactic Medium and the Damping Wing of the Gunn-Peterson Trough”. In: *The Astrophysical Journal* 501.1 (July 1998), pp. 15–22. ISSN: 1538-4357. DOI: 10.1086/305799. URL: <http://dx.doi.org/10.1086/305799>.
  - [21] Bradley Greig, Sarah E. I. Bosman, Frederick B. Davies, Dominika Ďurovčíková, Hassan Fathivavsari, et al. *Blind QSO reconstruction challenge: Exploring methods to reconstruct the Ly $\alpha$  emission line of QSOs*. 2024. arXiv: 2404.01556 [astro-ph.CO].
  - [22] Sarah E I Bosman, Dominika Ďurovčíková, Frederick B Davies, and Anna-Christina Eilers. “A comparison of quasar emission reconstruction techniques for  $z > 5.0$  Lyman-alpha and Lyman-beta transmission”. In: *Monthly Notices of the Royal Astronomical Society* 503.2 (Feb. 2021), pp. 2077–2096. ISSN: 1365-2966. DOI: 10.1093/mnras/stab572. URL: <http://dx.doi.org/10.1093/mnras/stab572>.
  - [23] Bin Liu and Rongmon Bordoloi. “A deep learning approach to quasar continuum prediction”. In: *Monthly Notices of the Royal Astronomical Society* 502.3 (Jan. 2021), pp. 3510–3532. ISSN: 1365-2966. DOI: 10.1093/mnras/stab177. URL: <http://dx.doi.org/10.1093/mnras/stab177>.
  - [24] Jonah C. Rose, Paul Torrey, Francisco Villaescusa-Navarro, Mark Vogelsberger, Stephanie O’Neil, Mikhail V. Medvedev, Ryan Low, Rakshak Adhikari, and Daniel Angles-Alcazar. *Inferring Warm Dark Matter Masses with Deep Learning*. 2023. arXiv: 2304.14432 [astro-ph.CO].

## Bibliography

---

- [25] Parth Nayak, Michael Walther, Daniel Gruen, and Sreyas Adiraju. *Ly $\alpha$ NNA: A Deep Learning Field-level Inference Machine for the Lyman- $\alpha$  Forest*. 2023. arXiv: 2311.02167 [astro-ph.CO].
- [26] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning: with Applications in R*. Springer US, 2021. ISBN: 9781071614181. DOI: 10.1007/978-1-0716-1418-1. URL: <http://dx.doi.org/10.1007/978-1-0716-1418-1>.
- [27] Shaeke Salman and Xiuwen Liu. *Overfitting Mechanism and Avoidance in Deep Neural Networks*. 2019. arXiv: 1901.06566 [cs.LG].
- [28] V. Roshan Joseph and Akhil Vakayil. “SPLit: An Optimal Method for Data Splitting”. In: *Technometrics* 64.2 (June 2021), pp. 166–176. ISSN: 1537-2723. DOI: 10.1080/00401706.2021.1921037. URL: <http://dx.doi.org/10.1080/00401706.2021.1921037>.
- [29] Michael W Browne. “Cross-Validation Methods”. In: *Journal of Mathematical Psychology* 44.1 (Mar. 2000), pp. 108–132. ISSN: 0022-2496. DOI: 10.1006/jmps.1999.1279. URL: <http://dx.doi.org/10.1006/jmps.1999.1279>.
- [30] Connor Shorten and Taghi M. Khoshgoftaar. “A survey on Image Data Augmentation for Deep Learning”. In: *Journal of Big Data* 6.1 (July 2019). ISSN: 2196-1115. DOI: 10.1186/s40537-019-0197-0. URL: <http://dx.doi.org/10.1186/s40537-019-0197-0>.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. *Deep Residual Learning for Image Recognition*. 2015. arXiv: 1512.03385 [cs.CV].
- [32] Sergey Ioffe and Christian Szegedy. *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. 2015. arXiv: 1502.03167 [cs.LG].
- [33] Udo von Toussaint. “Bayesian inference in physics”. In: *Rev. Mod. Phys.* 83 (3 2011), pp. 943–999. DOI: 10.1103/RevModPhys.83.943. URL: <https://link.aps.org/doi/10.1103/RevModPhys.83.943>.
- [34] Laurent Valentin Jospin, Hamid Laga, Farid Boussaid, Wray Buntine, and Mohammed Bennamoun. “Hands-On Bayesian Neural Networks—A Tutorial for Deep Learning Users”. In: *IEEE Computational Intelligence Magazine* 17.2 (May 2022), pp. 29–48. ISSN: 1556-6048. DOI: 10.1109/mci.2022.3155327. URL: <http://dx.doi.org/10.1109/MCI.2022.3155327>.

- 
- [35] Thomas G. Dietterich. “Ensemble Methods in Machine Learning”. In: *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2000, pp. 1–15. ISBN: 9783540450146. DOI: 10.1007/3-540-45014-9\_1. URL: [http://dx.doi.org/10.1007/3-540-45014-9\\_1](http://dx.doi.org/10.1007/3-540-45014-9_1).
  - [36] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. “Optuna: A Next-generation Hyperparameter Optimization Framework”. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2019.
  - [37] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. “Algorithms for hyper-parameter optimization”. In: *Proceedings of the 24th International Conference on Neural Information Processing Systems*. NIPS’11. Granada, Spain: Curran Associates Inc., 2011, pp. 2546–2554. ISBN: 9781618395993.
  - [38] Fahad Nasir, Prakash Gaikwad, Frederick B. Davies, James S. Bolton, Ewald Puchwein, and Sarah E. I. Bosman. *Deep Learning the Intergalactic Medium using Lyman-alpha Forest at  $4 \leq z \leq 5$* . 2024. arXiv: 2404.05794 [astro-ph.CO].
  - [39] Vid Iršič, Matteo Viel, Martin G. Haehnelt, James S. Bolton, Margherita Molaro, et al. “Unveiling dark matter free streaming at the smallest scales with the high redshift Lyman-alpha forest”. In: *PRD* 109.4, 043511 (Feb. 2024), p. 043511. DOI: 10.1103/PhysRevD.109.043511. arXiv: 2309.04533 [astro-ph.CO].
  - [40] Sandro D’Odorico, Stefano Cristiani, Hans Dekker, Vanessa Hill, Andreas Kaufer, Taesun Kim, and Francesca Primas. “Performance of UVES, the echelle spectrograph for the ESO VLT and highlights of the first observations of stars and quasars”. In: *Discoveries and Research Prospects from 8- to 10-Meter-Class Telescopes*. Ed. by Jacqueline Bergeron. Vol. 4005. Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series. June 2000, pp. 121–130. DOI: 10.1117/12.390133.
  - [41] George D. Becker, Paul C. Hewett, Gábor Worseck, and J. Xavier Prochaska. “A refined measurement of the mean transmitted flux in the Ly-alpha forest over  $2 < z < 5$  using composite quasar spectra”. In: *Monthly Notices of the Royal Astronomical Society* 430.3 (Feb. 2013), pp. 2067–2081. ISSN: 1365-2966. DOI: 10.1093/mnras/stt031. URL: <http://dx.doi.org/10.1093/mnras/stt031>.

## Bibliography

---

- [42] Prakash Gaikwad, Raghunathan Srianand, Martin G Haehnelt, and Tirthankar Roy Choudhury. “A consistent and robust measurement of the thermal state of the IGM at  $2 < z < 4$  from a large sample of Ly-alpha forest spectra: evidence for late and rapid HeIIreionization”. In: *Monthly Notices of the Royal Astronomical Society* 506.3 (July 2021), pp. 4389–4412. ISSN: 1365-2966. DOI: 10.1093/mnras/stab2017. URL: <http://dx.doi.org/10.1093/mnras/stab2017>.
- [43] Molly Wolfson, Joseph F. Hennawi, Frederick B. Davies, Zarija Lukić, and Jose Oñorbe. *Forecasting constraints on the high-z IGM thermal state from the Lyman- $\alpha$  forest flux auto-correlation function*. 2023. arXiv: 2309.05647 [astro-ph.CO]. URL: <https://arxiv.org/abs/2309.05647>.
- [44] Ravid Shwartz-Ziv and Naftali Tishby. *Opening the Black Box of Deep Neural Networks via Information*. 2017. arXiv: 1703.00810 [cs.LG].
- [45] Vanessa Buhrmester, David Münch, and Michael Arens. *Analysis of Explainers of Black Box Deep Neural Networks for Computer Vision: A Survey*. 2019. arXiv: 1911.12116 [cs.AI].
- [46] Jonas Fischer, Anna Olah, and Jilles Vreeken. “What’s in the Box? Exploring the Inner Life of Neural Networks with Robust Rules”. In: *Proceedings of the 38th International Conference on Machine Learning*. Ed. by Marina Meila and Tong Zhang. Vol. 139. Proceedings of Machine Learning Research. PMLR, 18–24 Jul 2021, pp. 3352–3362. URL: <https://proceedings.mlr.press/v139/fischer21b.html>.
- [47] Intekhab Hossain, Jonas Fischer, Rebekka Burkholz, and John Quackenbush. *Not all tickets are equal and we know it: Guiding pruning with domain-specific knowledge*. 2024. arXiv: 2403.04805 [cs.LG]. URL: <https://arxiv.org/abs/2403.04805>.
- [48] George D. Becker, James S. Bolton, Martin G. Haehnelt, and Wallace L. W. Sargent. “Detection of extended HeII reionization in the temperature evolution of the intergalactic medium: IGM temperatures over  $2 < z < 5$ ”. In: *Monthly Notices of the Royal Astronomical Society* 410.2 (Nov. 2010), pp. 1096–1112. ISSN: 0035-8711. DOI: 10.1111/j.1365-2966.2010.17507.x. URL: <http://dx.doi.org/10.1111/j.1365-2966.2010.17507.x>.
- [49] William H. Press, Brian P Flannery, Saul A Teukolsky, and William T Vetterling. *Numerical Recipes in C book set: Numerical Recipes in C: The Art of Scientific Computing*. 2nd ed. Cambridge, England: Cambridge University Press, Oct. 1992.

- [50] Michael T Murphy, Glenn G Kacprzak, Giulia A D Savorgnan, and Robert F Carswell. “The UVES Spectral Quasar Absorption Database (SQUAD) data release 1: the first 10 million seconds”. In: *Monthly Notices of the Royal Astronomical Society* 482.3 (Oct. 2018), pp. 3458–3479. ISSN: 1365-2966. DOI: 10.1093/mnras/sty2834. URL: <http://dx.doi.org/10.1093/mnras/sty2834>.
- [51] Ann S. Almgren, John B. Bell, Mike J. Lijewski, Zarija Lukić, and Ethan Van Andel. “Nyx: A MASSIVELY PARALLEL AMR CODE FOR COMPUTATIONAL COSMOLOGY”. In: *The Astrophysical Journal* 765.1 (Feb. 2013), p. 39. ISSN: 1538-4357. DOI: 10.1088/0004-637x/765/1/39. URL: <http://dx.doi.org/10.1088/0004-637X/765/1/39>.
- [52] Feige Wang, Xue-Bing Wu, Xiaohui Fan, Jinyi Yang, Zheng Cai, et al. “AN ULTRA-LUMINOUS QUASAR AT  $z = 5.363$  WITH A TEN BILLION SOLAR MASS BLACK HOLE”. In: *The Astrophysical Journal* 807.1 (June 2015), p. L9. ISSN: 2041-8213. DOI: 10.1088/2041-8205/807/1/19. URL: <http://dx.doi.org/10.1088/2041-8205/807/1/L9>.
- [53] Tom Charnock, Guilhem Lavaux, and Benjamin D. Wandelt. “Automatic physical inference with information maximizing neural networks”. In: *Physical Review D* 97.8 (Apr. 2018). ISSN: 2470-0029. DOI: 10.1103/physrevd.97.083004. URL: <http://dx.doi.org/10.1103/PhysRevD.97.083004>.
- [54] Alexander Ly, Maarten Marsman, Josine Verhagen, Raoul Grasman, and Eric-Jan Wagenmakers. *A Tutorial on Fisher Information*. 2017. arXiv: 1705.01064 [math.ST]. URL: <https://arxiv.org/abs/1705.01064>.
- [55] Soumik Maitra, Stefano Cristiani, Matteo Viel, Roberto Trotta, and Guido Cunradi. *Parameter estimation from Ly $\alpha$  forest in Fourier space using Information Maximising Neural Network*. 2024. arXiv: 2404.04327 [astro-ph.CO]. URL: <https://arxiv.org/abs/2404.04327>.
- [56] Clara Grazian and Yanan Fan. *A review of Approximate Bayesian Computation methods via density estimation: inference for simulator-models*. 2019. arXiv: 1909.02736 [stat.CO]. URL: <https://arxiv.org/abs/1909.02736>.