



Interpreting High-resolution Spectroscopy of Exoplanets using Cross-correlations and Supervised Machine Learning

Chloe Fisher¹ , H. Jens Hoeijmakers^{1,2}, Daniel Kitzmann¹ , Pablo Márquez-Neila³, Simon L. Grimm¹ ,
Raphael Sznitman³, and Kevin Heng¹ 

¹ University of Bern, Center for Space and Habitability, Gesellschaftsstrasse 6, CH-3012, Bern, Switzerland; chloe.fisher@csh.unibe.ch,
jens.hoeijmakers@space.unibe.ch, kevin.heng@csh.unibe.ch

² Observatoire astronomique de l'Université de Genève, 51 Chemin des Maillettes, 1290 Versoix, Switzerland

³ ARTORG Center for Biomedical Engineering, University of Bern, Bern, Switzerland

Received 2019 October 25; revised 2020 February 14; accepted 2020 February 25; published 2020 April 9

Abstract

We present a new method for performing atmospheric retrieval on ground-based, high-resolution data of exoplanets. Our method combines cross-correlation functions with a random forest, a supervised machine-learning technique, to overcome challenges associated with high-resolution data. A series of cross-correlation functions are concatenated to give a “CCF-sequence” for each model atmosphere, which reduces the dimensionality by a factor of ~ 100 . The random forest, trained on our grid of $\sim 65,000$ models, provides a likelihood-free method of retrieval. The precomputed grid spans 31 values of both temperature and metallicity, and incorporates a realistic noise model. We apply our method to HARPS-N observations of the ultra-hot Jupiter KELT-9b and obtain a metallicity consistent with solar ($\log M = -0.2 \pm 0.2$). Our retrieved transit chord temperature ($T = 6000_{-200}^{+0}$ K) is unreliable as strong ion lines lie outside of the extent of the training set, which we interpret as being indicative of missing physics in our atmospheric model. We compare our method to traditional nested sampling, as well as other machine-learning techniques, such as Bayesian neural networks. We demonstrate that the likelihood-free aspect of the random forest makes it more robust than nested sampling to different error distributions, and that the Bayesian neural network we tested is unable to reproduce complex posteriors. We also address the claim in Cobb et al. 2019 that our random forest retrieval technique can be overconfident but incorrect. We show that this is an artifact of the training set, rather than of the machine-learning method, and that the posteriors agree with those obtained using nested sampling.

Unified Astronomy Thesaurus concepts: [Exoplanet atmospheres \(487\)](#)

1. Introduction

1.1. Observational Motivation I: The Rise of Ground-based High-resolution Spectra

The observational characterization of exoplanetary atmospheres via the measurement of transmission and emission spectra is occurring on two fronts: low-resolution, space-based spectroscopy (mainly with the *Hubble Space Telescope* and the *Spitzer Space Telescope*), and high-resolution spectroscopy using a wide variety of ground-based spectrographs (Table 1). Spectra measured from space have the advantage that the spectral continuum, which encodes information on chemistry and clouds/hazes, may be measured in an absolute sense. Ground-based spectra lose the spectral continuum—and effectively measure *relative* transit depths or fluxes—due to having to correct for the presence of the Earth’s atmosphere, but offer the key advantage that individual spectral lines may be resolved with spectral resolution $\sim 10^5$. A plausible approach is to combine the advantages each has to offer and jointly analyze space- and ground-based spectra (e.g., Brogi et al. 2017).

Following the pioneering work of Snellen et al. (2008, 2010); (see also Wiedemann et al. 2001; Brown et al. 2002; Deming et al. 2005), the use of high-resolution, ground-based spectroscopy to identify the presence of atoms and molecules has become routine (Redfield et al. 2008; Brogi et al. 2012; Birkby et al.

2013, 2017; Brogi et al. 2013, 2014, 2018; de Kok et al. 2013; Lockwood et al. 2014; Wyttenbach et al. 2015, 2017; Piskorz et al. 2016, 2017, 2018; Khalafinejad et al. 2017, 2018; Nugroho et al. 2017; Hoeijmakers et al. 2018, 2019; Cauley et al. 2019; Guilluy et al. 2019; Seidel et al. 2019). These identifications are essentially model independent, relying only on knowledge of the cross sections or opacities of these atoms and molecules as determined by quantum physics (e.g., Rothman et al. 1998; Heng 2017). Line transition databases contain the positions and relative strengths of individual lines, either from experimental measurement or derived from first principles, which are then cross-correlated against the lines detected in the high-resolution spectrum. By matching dozens to hundreds of lines using cross-correlation, robust identifications of atoms and molecules may be obtained (but see Hoeijmakers et al. 2015; Brogi & Line 2019 for examples of detections being dependent on the accuracy of the line database used to compute these opacities). In contrast, the claimed detections of molecules other than water in the Wide Field Camera 3 (WFC3) spectra of exoplanetary atmospheres remains model dependent and an active topic of debate (e.g., Fisher & Heng 2018), because at these resolutions (~ 10) only the shapes of the overall opacities, consisting of a large collection of lines averaged together, are measured.

Interpreting ground-based, high-resolution spectra using the cross-correlation technique has one major shortcoming: cross-correlation is mainly capable of answering the binary question of whether an atom or molecule is absent or present, either in emission or absorption. It does not yield the abundance of that atom or molecule, nor the atmospheric temperature and pressure of the environment in which it lies. It similarly does

Table 1
High-resolution Cross-dispersed Echelle (grating) Spectrographs with Wide Instantaneous Wavelength Coverage

Name	Telescope	Resolving power	Wavelength Range (nm)	Status	Reference(s)
HARPS	ESO 3.6 m	120,000	378–691	Active	Mayor et al. (2003)
HARPS-N	TNG	120,000	378–691	Active	Cosentino et al. (2012)
ESPRESSO	VLT	70,000–190,000	378–691	Active	Pepe et al. (2014)
CARMENES	CAHA 3.5	80,000–100,000	520–1710	Active	Quirrenbach et al. (2010)
GIANO	TNG	50,000	950–2450	Active	Origlia et al. (2014)
CRIRES+	VLT	50,000–100,000	<i>Y, J, H, K, L, M</i> bands	Under development	Follert et al. (2014)
UVES	VLT	40,000–110,000	300–1100	Active	Dekker et al. (2000)
NIRSPEC	Keck	25,000	960–5500	Active	McLean et al. (1998)
PEPSI	LBT	43,000–270,000	383–912	Active	Strassmeier et al. (2015)
HDS	Subaru	90,000–165,000	298–1016	Active	Noguchi et al. (2002)
EXPRES	DCT	150,000	380–844	Active	Fischer et al. (in prep)
HIRES	ELT	100,000	397–2500	Under development	Zerbi et al. (2014)
NIRPS	ESO 3.6 m	80,000	974–1809	Under development	Wildi et al. (2017)
SPIRou	CFHT	70,000	980–2440	Active	Donati et al. (2018)
iShell	IRTF	75,000	<i>J, H, K, L, M</i> bands	Active	Rayner et al. (2016)
IGRINS	HJS	40,000	1450–2450	Active	Park et al. (2014)

not yield cloud or haze properties of the atmosphere. The first study to decisively address this shortcoming was Brogi & Line (2019), who re-analyzed CRILES observations and derived an analytical expression that maps the cross-correlation function to the likelihood function. The ability to compute the likelihood function implies that Bayes’s Theorem may subsequently be invoked to compute posterior distributions of chemical abundances, temperature, etc.

CRILES was an infrared echelle spectrograph mounted on UT1 of ESO’s VLT (Kaeufl et al. 2004). Although the spectrograph achieved high spectral resolution of $\sim 100,000$, the instantaneous wavelength coverage was small because the spectrograph was not cross-dispersed. Consequently, the spectra analyzed by Brogi & Line (2019) contain only 4096 data points (1.9626–2.0045 μm , 2.2875–2.3454 μm in two different modes). As every model being computed in the atmospheric retrieval needs to be cross-correlated against the spectrum, it becomes computationally prohibitive to scale this method up to spectra of cross-dispersed echelle spectrographs that contain $\sim 10^5$ – 10^6 data points, because this increases the computational time by a factor $\sim 10^2$ – 10^3 . However, elucidating such a scalable method is crucial in the era of high-resolution spectrographs with wide *instantaneous* wavelength coverage, an overview of which we list in Table 1. Gibson et al. (2020) also recently performed retrieval on high-resolution data from the blue arm of UVES with $\sim 10^3$ data points using a Markov Chain Monte Carlo (MCMC) method.

A novel method to analyze ground-based, high-resolution spectra with $\sim 10^5$ – 10^6 data points is therefore needed that will allow the computational effort to be reduced at the order-of-magnitude level *and* allow for the computation of posterior distributions of parameters.

1.2. Observational Motivation II: Failure of Direct Retrievals on Noisy Spectra

Another major limitation of ground-based, high-resolution spectra is the observational uncertainty. The level of noise on each individual spectral data point is typically much greater than the signal itself, which causes the direct retrieval to fail (see Section 3.1). While each individual spectral point contains little information, the entire spectrum does encode valuable information on the atmospheric abundances and properties. Any successful interpretation method needs to leverage the

information content of the entire spectrum against the high level of noise present.

This is the rationale behind the cross-correlation technique, which has been adopted by many workers (e.g., Snellen et al. 2010; Brogi et al. 2012; Birkby et al. 2013; de Kok et al. 2013; Lockwood et al. 2014; Wyttenbach et al. 2015; Piskorz et al. 2016; Nugroho et al. 2017; Hoeijmakers et al. 2018; Guilluy et al. 2019; Seidel et al. 2019), including Brogi & Line (2019).

In the current study, we will incorporate the cross-correlation technique into a novel method for performing retrievals on noisy, high-resolution spectra, but in a way that is distinct from Brogi & Line (2019).

1.3. Theoretical Motivation I: Likelihood-free Inference Methods using Machine Learning

In the published exoplanet literature, atmospheric retrievals typically assume the likelihood function to be a Gaussian when implementing the Markov Chain Monte Carlo (MCMC) or nested-sampling routines (e.g., Benneke & Seager 2012; Line et al. 2013; Waldmann et al. 2015; Lavie et al. 2017; MacDonald & Madhusudhan 2017; Fisher & Heng 2018; Brogi & Line 2019),

$$\ln \mathcal{L} = -\frac{1}{2} \sum_i^n \left(\frac{R_i - R_{i,\text{obs}}}{\sigma_i} \right)^2 - \frac{\ln(2\pi\sigma_i^2)}{2}, \quad (1)$$

where the transmission spectrum has n measurements of transit radii ($R_{i,\text{obs}}$) that are compared to the theoretical values of the transit radii (R_i) computed using a model. The standard deviation of the uncertainty on each data point, assumed to follow a Gaussian distribution, is σ_i . It is further assumed that the uncertainties are uncorrelated with one another.

One of the motivations of the current study is to provide an alternative inference approach that is likelihood free, meaning that one does not have to explicitly assume the functional form of the likelihood function. In practice, these likelihood-free inference approaches belong to the class of Approximate Bayesian Computation (ABC) methods (Sisson et al. 2019). Specifically, we use the supervised machine-learning method of the random forest (Ho 1998; Breiman 2001), which was previously adapted by Márquez-Neila et al. (2018) to interpret low-resolution *Hubble*-WFC3 transmission spectra. The

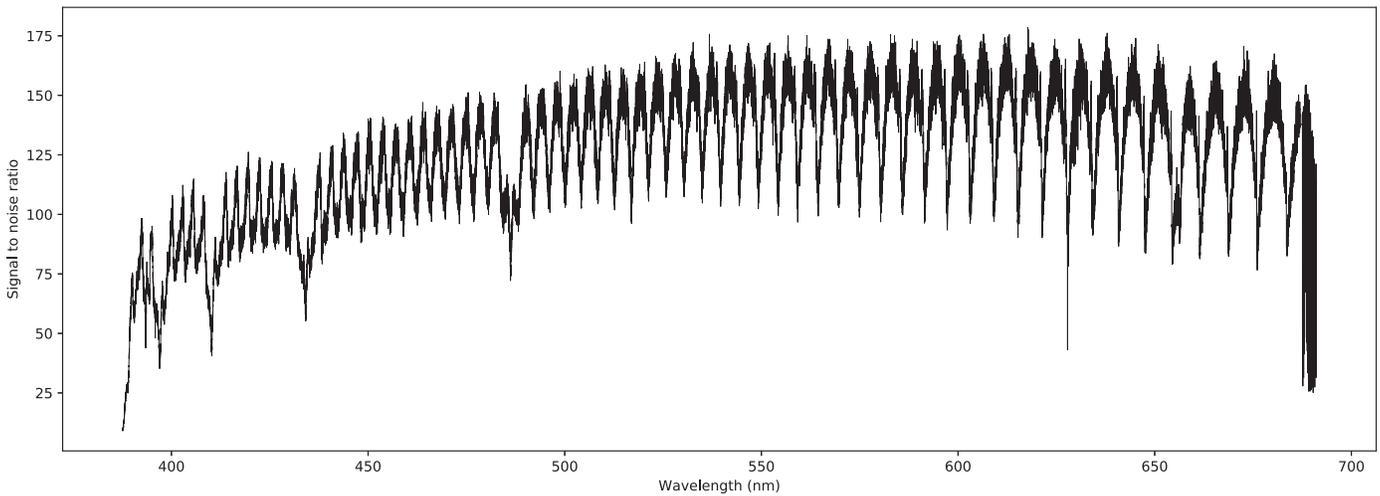


Figure 1. The signal-to-noise level of the spectrum of the host star KELT-9 achieved in a 600 s exposure obtained with the HARPS-N instrument. The signal-to-noise ratio is dominated by the photon (shot) noise, which decreases toward shorter wavelengths due to a reduced efficiency of the instrument, transmission of the Earth’s atmosphere, and lower intrinsic luminosity of the star. The significant narrowband variation is due to the efficiency of the spectrograph falling off at the edges of spectral orders, as well as absorption lines in the star and the Earth’s atmosphere.

method relies on using a grid of precomputed atmospheric models combined with an arbitrary noise model as a training set for the random forest. The uncertainties on each data point in the measured spectrum are incorporated into the noise-free model grid to generate a training set of noisy models. This approach is not unlike that of standard retrieval techniques, which typically compute a grid of atmospheric models on the fly.

The random forest consists of a collection of regression trees. Each regression tree is trained on a subset of the grid of atmospheric models. By identifying regions of the multi-dimensional parameter space that predict similar transmission spectra, each regression tree quantifies the “distance” between the model and measured transmission spectra. This plays the role of the Euclidean distance ($R_i - R_{i,\text{obs}}$) in the Gaussian likelihood function, except that the likelihood is implicitly learned from the training set of noisy models. (See Section 2.6.1 for more information about the random forest.)

Other advantages offered by the random forest retrieval method include the ability to run large suites of mock retrievals to both validate the model grid used and quantify its sensitivity to the parameters, as well as information content analysis to quantify the relative importance of each data point in the spectrum toward determining the value of each parameter (Márquez-Neila et al. 2018).

1.4. Theoretical Motivation II: Feature Engineering

Feature engineering is the process by which the training set used in a machine-learning method is optimized, e.g., a reduction in the dimensionality of the problem. Deep learning methods perform feature engineering in an automated way, but they are significantly more expensive to implement than the random forest. One of the novel aspects of the current study is the use of feature engineering to efficiently interpret noisy, high-resolution spectra. Instead of using the spectra themselves as the training set, we demonstrate that it is sufficient to use a set of cross-correlation functions (CCFs) that sparsely sample the parameter space. The resulting “cross-correlation sequence” serves as the training set for the random forest, resulting in a reduction in the size of the training set by a factor of ~ 100 .

This feature engineering step allows the random forest retrieval method to be scaled up to interpret high-resolution spectra with $\sim 10^5$ – 10^6 data points in a computationally feasible way.

1.5. Layout of Study

In Section 2, we describe our methodology, including the computation of the model grid of transmission spectra (radiative transfer, opacities, and chemistry), the implementation of the random forest method, etc. In Section 3, we show our results from testing the method, and also the retrieval on HARPS-N observations of KELT-9b. In Section 4, we discuss the results and compare our method to nested sampling and other machine-learning techniques. In Section 5, we summarize our conclusions.

2. Methods

2.1. KELT-9b

As a proof of concept and in order to test the method, we have focused the retrieval on the ultra-hot Jupiter, KELT-9b. The brightness of the star combined with the extremely high temperatures allow for a higher signal-to-noise ratio (SNR) than for other exoplanets (see Figure 1), making it a good test subject for a retrieval on ground-based data. Furthermore, this object has been previously studied with high-resolution data in Hoeijmakers et al. (2018, 2019). Kitzmann et al. (2018) demonstrated that chemical equilibrium is a reasonable assumption, significantly reducing the number of parameters required in the atmospheric model, and that it is cloud-free with a continuum dominated by H^- (Arcangeli et al. 2018). However, Hoeijmakers et al. (2019) suggested that there is most likely missing physics in this model, due to the discrepancy between the expected cross-correlation function for Fe^+ and the one obtained from the data. We will discuss this further in Section 3.4.

2.2. Model Grid

To construct the grid of models of KELT-9b, we adopt the system parameters reported by Gaudi et al. (2017) and Hoeijmakers et al. (2019). We generate the models using an

observation simulator, `Helios-o` (Bower et al. 2019), which follows the method described in Gaidos et al. (2017). This algorithm has been validated in Heng & Kitzmann (2017), where it was compared against the models from Fortney et al. (2010), Deming et al. (2013), and Line et al. (2013).

The model atmosphere is one-dimensional, plane-parallel, isothermal, in hydrostatic equilibrium, and in chemical equilibrium. It has 199 layers with 200 pressure levels ranging from 10^{-15} –2 bar. Each one-dimensional model atmosphere may be visualized as an atmospheric column. Ray tracing is performed through a collection of these atmospheric columns to construct the transit chord at each wavelength, taking into account the variation of gravity as different pressure levels are probed. The variation of the effective transit radius with wavelength due to the chemical composition of the atmosphere is the transmission spectrum (Brown 2001).

The volume mixing ratios (relative abundances by number) of atoms, ions, and molecules are computed using the `FastChem` chemical-equilibrium code, which considers gas-phase chemistry for more than 550 molecular species with elements more abundant than germanium (Stock et al. 2018). Additionally, we add most of the first and doubly ionized ions as well as anions for atoms lighter than neptunium (Hoeijmakers et al. 2019). Our volume mixing ratios computed using `FastChem` are pressure dependent, because of our nonisobaric treatment of the transit chord. The opacities are computed using the open-source `HELIOS-K` opacity calculator (Grimm & Heng 2015). The inputs for the Fe, Fe⁺, Ti, and Ti⁺ opacities are sourced from the Kurucz database⁴ (Kurucz 2017). The hydrogen anion (H⁻) cross section is taken from John (1988). For completeness, collision-induced absorption associated with H–He, H₂–H₂, and H₂–He collisions are included (Richard et al. 2012). Pressure broadening is neglected as the spectral continuum in ultra-hot Jupiters is dominated by absorption associated with the hydrogen anion (H⁻), which masks the line wings. The line shape is assumed to be a Voigt profile. The natural line width and thermal broadening are included (Kurucz 2017). Opacities are sampled uniformly across wavenumber with a spectral resolution of 0.01 cm⁻¹, and the transmission spectra are calculated at a resolution of 0.03 cm⁻¹.

The assumption of chemical equilibrium allows us to greatly simplify the theoretical analysis because the abundances of atoms and ions are completely specified by the temperature, pressure, and elemental abundances. By assuming the ratios of elemental abundances follow those of the Sun, we reduce the chemical parameters down to a single number known as the metallicity. Therefore, we have just two parameters in our model—temperature and metallicity. The temperature range of the grid spans from 3000 to 6000 K, in steps of 100 K, and the metallicity ranges from 0.1 to 100 times solar (–1 to 2 for the logarithm of the metallicity, $\log M$, in steps of 0.1). This results in 31 values for each parameter, and thus 961 models in the grid in total.

2.3. Modeling HARPS-N Observations

We use existing observations of KELT-9b produced by the HARPS-N spectrograph (Hoeijmakers et al. 2018) to convert the resulting model grid to models of the observed transmission spectrum. First, the transmission spectrum is convolved with a Gaussian with a full-width-at-half-maximum of 2.7 km s⁻¹

(equivalent to the resolving power of the HARPS-N spectrograph), as well as a rotation-broadening profile that matches the rotation period of KELT-9b. It is subsequently interpolated onto the wavelength grid of the stitched, resampled pipeline-reduced (s1d) observations from HARPS-N. The continuum of the transmission spectrum is removed using a high-pass filter, in the same way as the observations with the HARPS-N spectrograph are filtered to remove broadband spectral variations that are due to the instrument and variable observing conditions (Hoeijmakers et al. 2018).

It would be possible to use this retrieval method for other instruments, such as those listed in Table 1, however these would require different training sets to account for other observational effects. The noise model (see Section 2.5) would also need to be adjusted for different instruments.

2.4. CCF-sequences

We use the cross-correlation operator defined as

$$C(v) = \frac{\sum_i F_i \mathcal{T}_i(v)}{\sum_i \mathcal{T}_i(v)}, \quad (2)$$

where F is the transmission spectrum, \mathcal{T} is the cross-correlation template interpolated onto the same wavelength grid as the spectrum, v is the velocity, and the summation takes place over the spectral data points. The denominator is a normalization factor, and thus the fluxes of the templates do not need to be rescaled when performing the cross-correlation.

Four subsets of cross-correlation templates, consisting of the spectral lines of neutral iron (Fe), singly ionized iron (Fe⁺), neutral titanium (Ti), and singly ionized titanium (Ti⁺), are created. Within each subset, there are 16 templates consisting of 4 values of temperature (3000, 4000, 5000, and 6000 K) and 4 values of metallicity (0.1, 1, 10, and 100×solar). In total, there are 64 cross-correlation templates. These templates are generated in the same way as the models (Section 2.2) with all but the relevant species’ opacities removed from the final model, leaving only the required species’ spectral lines. Broadening is not included as we are not aiming to retrieve dynamic properties. (See Section 4.2 for tests involving velocity parameters.)

Each synthetic transmission spectrum in the model grid is cross-correlated with each of the 64 templates to create a set of 64 CCFs. Additionally, each template is shifted in velocity space from –20 to 20 km s⁻¹ in steps of 1 km s⁻¹, resulting in 40 CCF values per template. These 64 CCFs are concatenated together to give a single sequence containing 2560 points, which we term a “CCF-sequence” (Figure 2). Each of the 64 templates probes different components of the information contained in the spectral lines. In this way, the resulting CCF-sequence encodes the physical properties of the atmosphere over multiple axes. This feature engineering step has essentially reduced the dimensions of the model spectra by a factor of ~100.

2.5. Noise Model

Because KELT-9 is a bright star, the noise is dominated by photon noise, and the SNR mainly varies due to the wavelength-dependent efficiency of the instrument, the stellar spectrum, and Earth’s atmospheric transmission function (see Figure 1). The noise per spectral pixel is empirically measured from the time series of observations used by Hoeijmakers et al. (2018).

⁴ <http://kurucz.harvard.edu/>

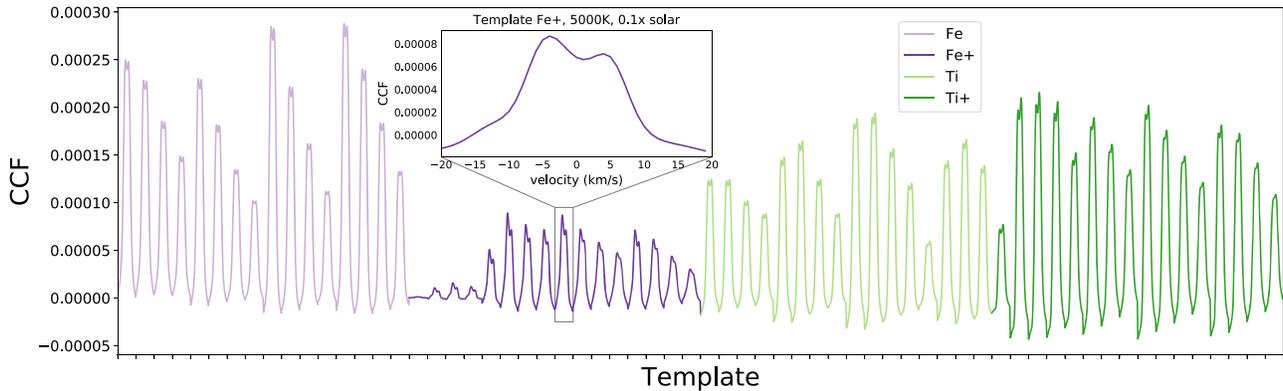


Figure 2. Example of a CCF-sequence constructed by cross-correlating 64 templates with a model transmission spectrum with $T = 3500$ K and $\log M = 0.8$. Each CCF has 40 points across velocity for a total of 2560 points for the entire CCF-sequence. The insert magnifies one of the CCFs (Fe^+ , $T = 5000$ K, and $\log M = 0.1$) for illustration.

For each spectral pixel a value may be drawn randomly from an assumed Gaussian distribution, creating a model of the noise of the entire spectrum that can be propagated through the cross-correlation function.

We assume each point in the spectrum F has a Gaussian error bar with standard deviation σ_{F_i} . The noise model for the CCF then becomes a linear combination of Gaussians, therefore also a Gaussian, with a variance of

$$\sigma_C^2 = \frac{\sum_i \sigma_{F_i}^2 \mathcal{T}_i(v)^2}{\sum_i \mathcal{T}_i(v)^2}. \quad (3)$$

We can then add the noise to the model grid of CCF-sequences. Since we require many instances of noise for the random forest, and the cross-correlation is computationally quite expensive, this provides a great advantage over applying the CCF to the noisy spectra.

2.6. Random Forest

2.6.1. Theory

The random forest consists of a collection of regression trees—decision trees for interpreting continuous data. Each regression tree is trained on a subset of the grid of atmospheric models. During training, a tree is constructed by locating divisions in each wavelength dimension that sort the training spectra into groups with similar parameter values, known as leaves. Each leaf then has an assigned set of parameter values given by the training spectra in its group. When predicting on a real data set, the spectrum is passed down each tree until it lands in a leaf, and the predicted parameter values are given by the corresponding set. The sets for every tree in the forest are then combined to give a distribution for each parameter.

The random forest falls into a class of inference methods known as “ABC” (Sisson et al. 2019). ABC methods were invented to treat problems where it was either infeasible or impossible to explicitly specify the functional form of the likelihood (e.g., in the study of human populations). Instead of seeking the maximum likelihood in a multidimensional parameter space, ABC methods seek to minimize some abstract distance (with the Euclidean distance being one specific example) between a set of simulated models and data to below some stated tolerance (Chapter 1.3 of Sisson et al. 2019). If the tolerance is formally zero, then ABC methods become exact Bayesian methods, which have been shown to produce accurate posteriors (Chapter 1.6 of Sisson et al. 2019). In practice,

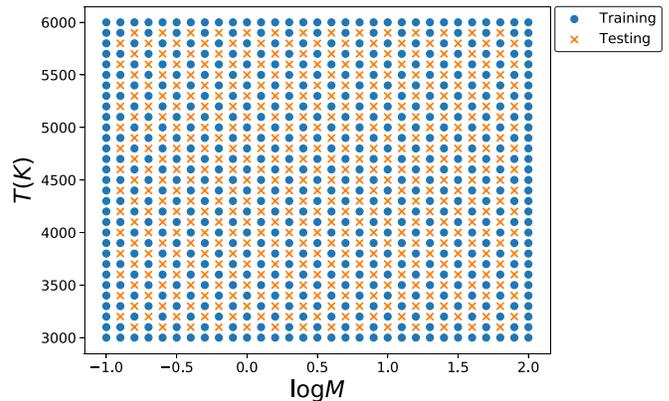


Figure 3. Separation of the 961 members of the model grid into training and testing sets for the random forest. The edges of this parameter space are intentionally included in the training set as the forest is unable to extrapolate.

nonzero tolerances generally imply that the computed posterior distributions are approximate (hence the “A” in “ABC”), where the degree of accuracy depends on the tolerance specified (Chapter 1.5 of Sisson et al. 2019). ABC methods often employ “summary statistics” as a dimensionality reduction step (Chapter 1.7 of Sisson et al. 2019). In the current study, the use of the CCF-sequence qualifies as a use of summary statistics.

2.6.2. Setup

Starting from our grid of CCF-sequences, we divide the parameter space into training and testing sets, as shown in Figure 3. This is to ensure the two sets are sufficiently distinct such that we can accurately test the performance of the forest. Next, we sample each point of the CCF-sequence within its respective uncertainty to generate 120 noisy instances of each CCF-sequence. We do this by drawing from Gaussian distributions with variance defined by Equation (3). The entire set therefore amounts to 115,320 noisy CCF-sequences, with 64,920 in training and 50,400 in testing.

Our random forests consists of 1000 trees. Tree splitting is performed using a threshold variance of 0.01. Each time a tree is split, a random subset of 50 (approximately the square root) of the 2560 sequence points is used. Tree pruning methods are not used (see Breiman et al. 1984; Hastie et al. 2001 for clarification of the terminology). For the predictions, the data is passed down through each tree until it reaches an end point,

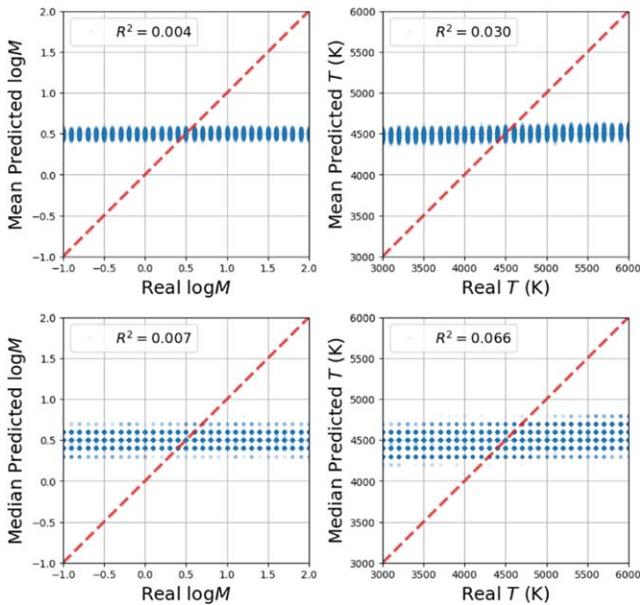


Figure 4. Predicted vs. real values of the logarithm of metallicity ($\log M$) and temperature (T) for the random forest trained using a section of the high-resolution spectrum containing 10^4 points from 400 to 410 nm. The top and bottom sets of plots correspond to the mean and median predictions, respectively. The coefficient of determination (R^2) varies from -1 to 1 , where values near unity indicate strong anticorrelations or correlations between the real and predicted values of a given parameter, based on the variance of outcomes. See Figure 5 for a mock retrieval.

known as a leaf. The set of all training parameters that lie in this leaf are then given as the prediction for that tree. We call this the “full-leaf” prediction. These training parameters come from the bootstrapped training data set—built using random sampling with replacement from the original training data set—that was used to train each tree. The final posterior is constructed by combining these predictions for all of the 1000 trees. This full-leaf prediction is an improvement on the previous method in Márquez-Neila et al. (2018), in which only the mean parameter values corresponding to the predicted leaf were used, as it gives a more accurate approximation of the posterior. The implementation of the random forest method and R^2 metric are adopted from the open-source `scikit.learn` library (Pedregosa et al. 2011) in the Python programming language.

3. Results

3.1. Failure of Direct Retrieval

Initially we attempted to perform the random forest retrieval directly on the transmission spectra, set up in the same way as described in Section 2.6.2 but with the model spectra instead of the CCF-sequences. Since the random forest method has been demonstrated to work for a dimensionality of at most $\sim 10^4$ (Hastie et al. 2001; Sznitman et al. 2013; Zikic et al. 2014; Rieke et al. 2015; Zhang et al. 2017), we consider only a section of 10^4 wavelength points from 400 to 410 nm in each synthetic spectrum. Other sampling strategies (e.g., selecting line peaks only) produce similar outcomes⁵ (not shown). Figure 4 shows the results of testing this forest, using both the mean (top panels) and median (bottom panels) predictions. The

⁵ While selecting line peaks is conceptually similar to a cross-correlation, by not averaging the spectral lines the noise remains high and hence the retrieval still fails.

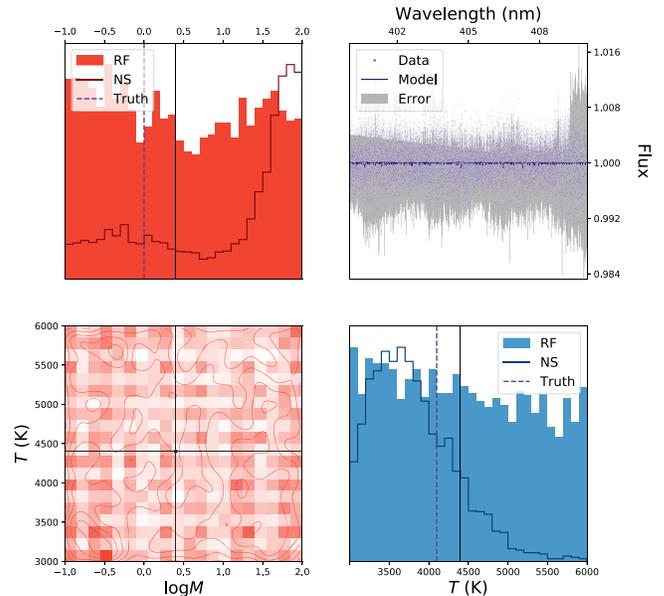


Figure 5. A mock retrieval using a section of the high-resolution spectrum containing 10^4 points from 400 to 410 nm, from the test set shown in Figure 4. The mock spectrum has solar metallicity and a temperature of 4100 K. In the top left and bottom right panels, the solid posteriors show the results of the retrieval using the random forest (RF), and the empty line posteriors show the results from nested sampling (NS). The purple, dashed lines show the true values. The top right panel shows the data points (lilac) with the error region (gray), along with the model (dark purple) corresponding to the medians from the $\log M$ and T posteriors.

coefficient of determination, R^2 , which measures the degree of agreement between the real versus predicted parameter values, is essentially zero for temperature and metallicity for both mean and median predictions, implying that the random forest has no predictive power when applied to the synthetic spectra themselves. Figure 5 includes an example of the posterior distributions of temperature and metallicity for a mock retrieval, which are unconstrained and consistent with their prior distributions. In addition, we tested a traditional retrieval algorithm using nested sampling (Skilling 2006; Feroz & Hobson 2008; Feroz et al. 2009, 2019) with the open-source `PyMultinest` package (Buchner et al. 2014). Due to the high number of spectral points and complex forward model, we are unable to compute models on the fly as in a regular nested-sampling retrieval (see Section 4.1). Instead, we take the same grid of models as the forest, but without the added noise, and interpolate on it to produce forward models. Figure 5 also shows the results from the nested-sampling mock retrieval. These posteriors span essentially the whole prior, with peaks offset from the correct values.

In summary, ground-based high-resolution spectra of exoplanets reside in a qualitatively different regime than the same measurements of stars or space-based low- to medium-resolution spectra of exoplanets. Individual data points hold little information as they are overwhelmed by noise, but the entire spectrum does encode useful information. This motivates our use of the CCFs, which effectively select the most informative lines in the spectrum.

3.2. Random Forest Mock Retrievals

Figure 6 shows the results of testing the random forest trained on the CCF-sequences. The predictive power of the

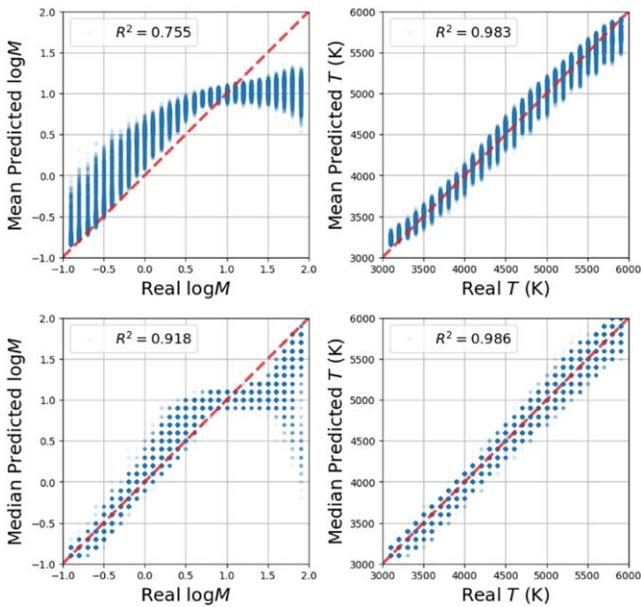


Figure 6. Predicted vs. real values of the logarithm of metallicity ($\log M$) and temperature (T) for the random forest trained on the CCF-sequences. The top and bottom panels show the results using the mean and median predictions, respectively. The coefficient of determination (R^2) varies from -1 to 1 , where values near unity indicate strong anticorrelations or correlations between the real and predicted values of a given parameter, based on the variance of the outcomes. See Figure 7 for a mock retrieval.

random forest has increased significantly. The difference in the predictability of the two parameters, metallicity and temperature, follows our intuition. The strength of spectral features are proportional to the logarithm of the opacity multiplied by the abundance of an atom. Because opacities have an exponential dependence on temperature (Rothman et al. 1998; Heng 2017), the line strengths are highly sensitive to temperature, and the ability of the random forest to predict temperature is strong. The ability to predict metallicity is somewhat weaker, because the metallicity linearly controls the atomic abundances, the logarithm of which determines the line depths (e.g., Heng & Kitzmann 2017). At high metallicities, the predictive power of the random forest tapers off, because the pressure scale height of the atmosphere decreases and the size of spectral features starts to decrease (see Section 3.3). The top and bottom panels of Figure 6 correspond to the mean and median predictions of the trees, respectively. Traditionally, random forests produce mean predictions, but given the focus of atmospheric retrieval on posteriors and confidence intervals, we are more interested in the medians, which are more robust against asymmetric posteriors. The increase in R^2 scores when using the median comes particularly from these more complex posteriors. Figure 7 also shows an example of the posterior distributions obtained from the hybrid CCF retrieval, which recovers the injected values of temperature and metallicity accurately.

A useful, natural outcome of the random forest is the information content analysis known as the “feature importance.” This determines which data points hold the most importance for retrieving each parameter. Figure 8 shows the feature importance when predicting metallicity and temperature. As suggested by the bottom panel of Figure 8, the ion species control the temperature prediction. Rising temperatures cause the neutral species to collisionally ionize, initially increasing the abundances of Fe^+ and Ti^+ by orders of

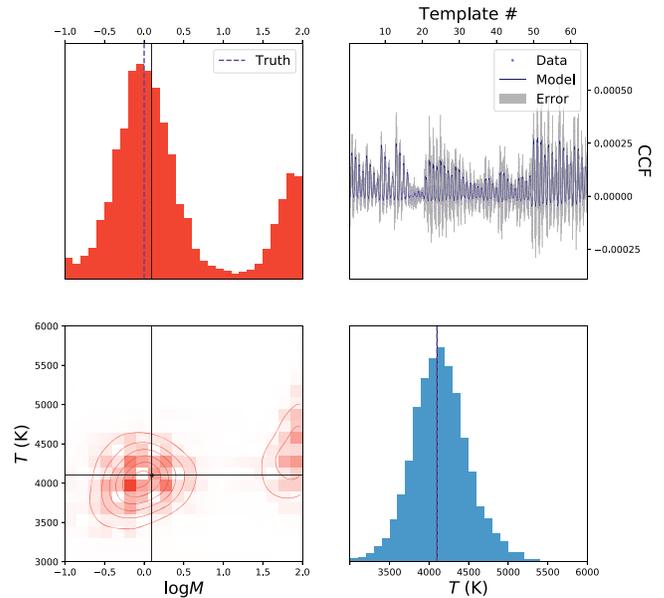


Figure 7. A mock retrieval performed on a model with solar metallicity and $T = 4100$ K, using the random forest trained on the CCF-sequences (see Figure 6). The black lines show the median values. The purple, dashed lines show the true values. The top right panel shows the data points (lilac) with the error region (gray), along with the model (dark purple) corresponding to the medians from the $\log M$ and T posteriors.

magnitude while the corresponding decrease in neutral abundance is relatively small.

As the metallicity increases, the depths of all metal absorption lines will tend to increase. However, in Figure 8 there appears to be a greater feature importance for the neutrals when predicting metallicity. A possible explanation for this is that as metallicity increases, the atmosphere will be more laden by free electrons from easily ionized species. Following the Saha equation (Saha 1920), this will lead to a decrease in the ionization fraction, partially negating the enhancement to the ion mixing ratios that stems from the higher metal abundance. Therefore, the neutral species are expected to be more sensitive to metallicity.

3.3. Metallicity Degeneracy

From our tests on the random forest in Figure 6, we can see that some of the high metallicity spectra yield much lower metallicity predictions. This is demonstrated further in Figure 9, which shows a retrieval on one of these high metallicity spectra. The double-peaked posterior leads to a mean prediction that is heavily offset from the true value. This multimodal structure is due to a degeneracy between line depth and metal abundance for high metallicity values. As discussed in Section 3.2, as the metallicity increases to very high levels, the atmosphere is no longer hydrogen dominated, causing the mean molecular weight to increase significantly. This in turn decreases the scale height and absorption line depths, reminiscent of lower metallicity values. We tested all the spectra with the highest metallicity value in the testing set ($\log M = 1.9$), and plotted the median predictions in Figure 10. This plot shows that the degeneracy is stronger at lower temperatures. This follows our physical intuition because at lower temperatures the pressure scale height is smaller, thus compressing the features and reducing the spectrum’s

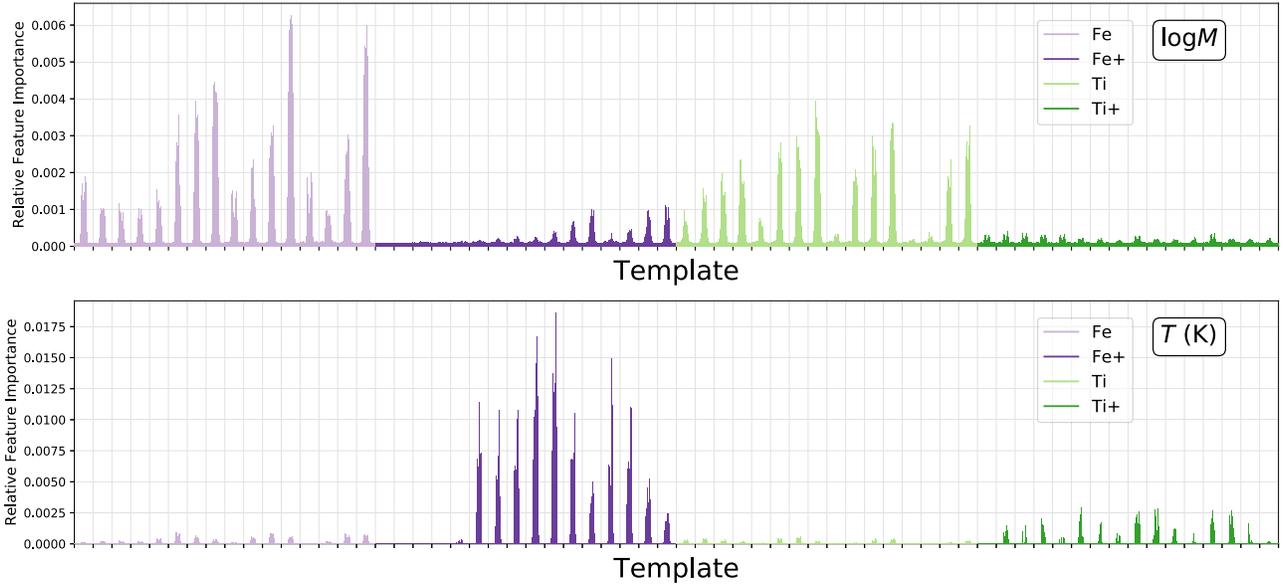


Figure 8. Feature importance plots describing the relative importance of each CCF in the sequence toward constraining metallicity and temperature.

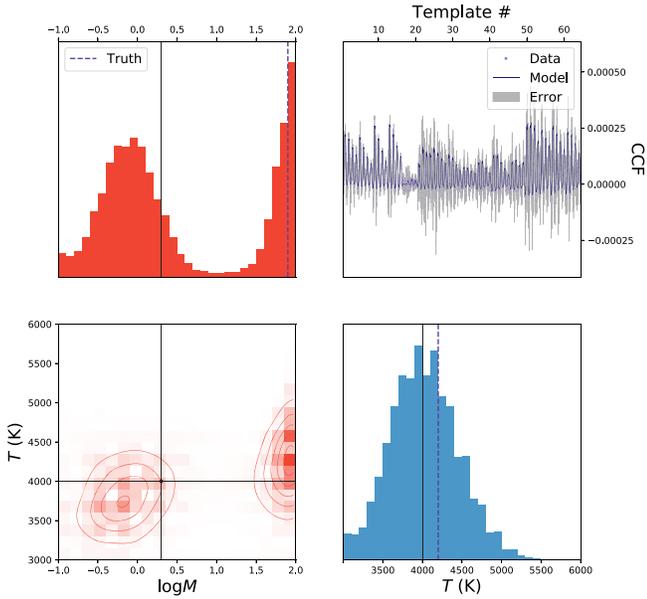


Figure 9. A mock retrieval performed on a model with $\log M = 1.9$ and $T = 4200$ K, using the random forest trained on the CCF-sequences (see Figure 6). The black lines show the median values. The purple, dashed lines show the true values. The top right panel shows the data points (lilac) with the error region (gray), along with the model (dark purple) corresponding to the medians from the $\log M$ and T posteriors.

sensitivity to metallicity. This makes these spectra harder to distinguish from one another for a given SNR.

This degeneracy is also visible in Figure 11, which shows noise-free spectra with $T = 3000$ K and varying metallicities, and a cross-correlation with those spectra. As the metallicity increases, the troughs in the left-hand plot deepen up to a point, after which they become shallower again. Similarly, the height of the CCFs in the right-hand plot increase with metallicity until $\log M \gtrsim 1.0$, after which the peaks decrease again. While the shape of the high and low metallicity noise-free spectra do differ slightly from each other, these variations are within the error bars of the data, making the noisy spectra indistinguishable.

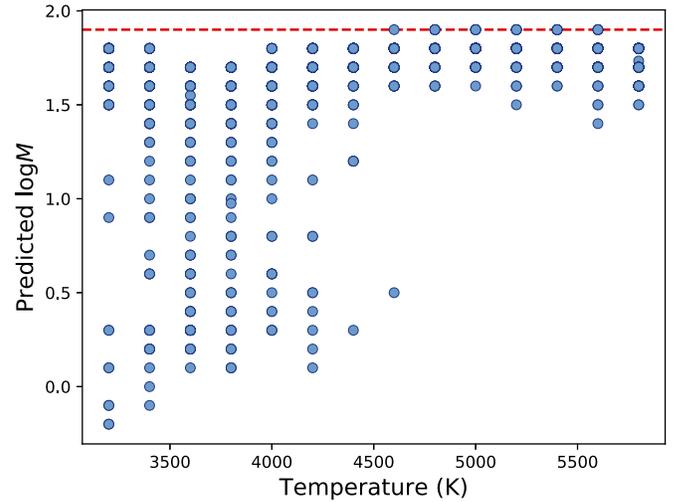


Figure 10. Median predictions for metallicity vs. the true temperature value for the test spectra with $\log M = 1.9$, from Figure 6. The red, dashed line shows the true metallicity value, 1.9.

3.4. KELT-9b Retrieval

Finally, we performed the hybrid CCF retrieval on the real HARPS-N data set for the ultra-hot Jupiter KELT-9b. Figure 12 shows our results for several different retrievals. As described in Hoeijmakers et al. (2019), the ionized iron lines in the spectrum of KELT-9b appear to be much larger than predicted, possibly resulting from an outflowing envelope not present in the model. This leads to a CCF-sequence for the real KELT-9b data that features significantly higher peaks in the Fe^+ CCFs when compared to the training set, as shown in Figure 13. With the intent of comparing the effects of the different species, we performed three independent retrievals on the KELT-9b data set—one containing the full CCF-sequence, as described in Section 2.4, a second containing only the neutral elements, and a third containing only the ions. Each retrieval uses a separate random forest trained on the corresponding sections of the model CCF-sequences. The three retrievals are compared in Figure 12, where the empty

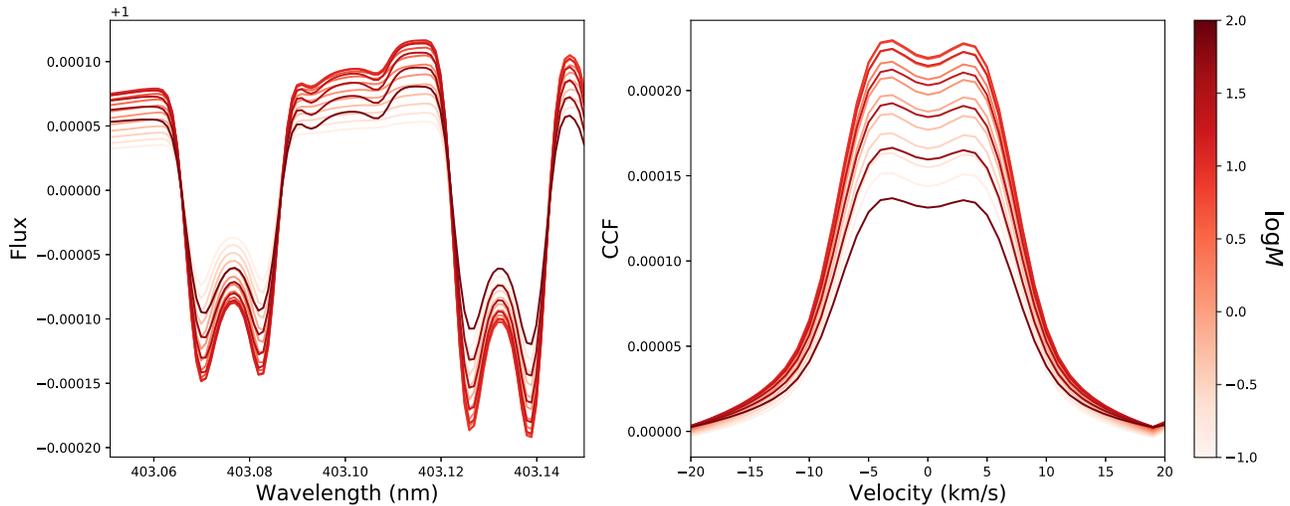


Figure 11. Noise-free synthetic spectra with $T = 3100$ K and varying metallicity values. The left-hand plot shows a zoomed-in section of the transmission spectra themselves, while the right-hand plot shows a single cross-correlation with each spectrum and the template for Fe at $T = 3000$ K and $\log M = -1.0$. The darker color corresponds to higher metallicity values.

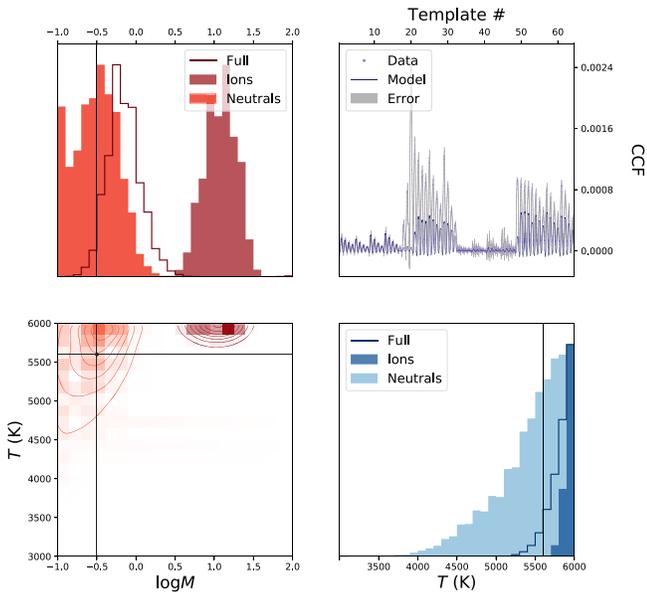


Figure 12. Retrieval performed on the CCF-sequence of the transmission spectrum of KELT-9b measured by the HARPS-N spectrograph. The retrieval is performed in three different ways: using only neutrals (Fe, Ti); (see Figure 14), using only ions (Fe^+ , Ti^+) or using all four species (“Full”). The vertical and horizontal lines indicate the median values of the posterior distributions corresponding to the neutrals-only retrieval. The top right panel shows the data points (lilac) with the error region (gray) for the CCF-sequence produced by the KELT-9b HARPS-N data, along with the model (dark purple) corresponding to the medians from the $\log M$ and T posteriors.

lined, darker colored, and lighter colored posteriors show the results from the full, ionized, and neutral retrievals, respectively.

The metallicity prediction greatly varies between the different retrievals, which is not unexpected here. The extremely high temperatures cause most of the neutral species to be ionized, leading to low abundances for Fe and Ti. Thus, in the neutral retrieval we predict a low $\log M$ value of $-0.5^{+0.2}_{-0.4}$, while the ion retrieval predicts 1.0 ± 0.2 . The full retrieval lies further toward the neutral prediction, with $\log M = -0.2 \pm 0.2$, which is unsurprising due to the stronger feature importance in the neutral CCFs for metallicity.

When the Fe^+ CCFs are included, i.e., in the full and ion retrievals, the temperature prediction is forced to its upper limit in an attempt to match the strong Fe^+ lines ($T = 6000^{+0}_{-200}$ K and $T = 6000^{+0}_{-100}$ K for the full and ion retrievals, respectively). However, in the neutral retrieval we still obtain a very high temperature value of 5600^{+400}_{-600} K, suggesting it is not only the excess Fe^+ that escalates the temperature prediction. Figure 14 shows the “predicted versus real” graphs for the forest trained only on the neutrals. As the temperature increases, this forest’s predictive ability decreases, as expected due to ionization. This suggests that the neutral posterior for temperature in Figure 12 may not be reliable. A positive conclusion is that this method is able to identify when a model is flawed.

Using *TESS* photometry, Wong et al. (2019) constrain the dayside and nightside temperatures of KELT-9b to be 4570 ± 90 K and 3020 ± 90 K, respectively. However, this is not inconsistent with a higher retrieved temperature from transmission spectroscopy. The dayside spectrum traces higher pressures than the transmission spectrum, which probes tenuous layers of the upper atmosphere. The present retrieval would be consistent with the scenario of an inversion layer, as is predicted in highly irradiated exoplanets (Hubeny et al. 2003; Fortney et al. 2008).

4. Discussion

4.1. Comparison to Nested Sampling

One of the most common techniques for performing atmospheric retrieval is nested sampling (Skilling 2006; Feroz & Hobson 2008; Feroz et al. 2009, 2019). In a traditional retrieval, a relatively computationally inexpensive forward model is used to generate spectra on the fly, while the sampling method searches the parameter space for the optimal solution. Brogi & Line (2019) demonstrate a method for performing retrieval on high-resolution data with nested sampling, but are restricted to ~ 4000 spectral data points of the CRRES instrument. As the number of spectral points increases, so does the time required to compute the models, making this method infeasible for a full HARPS-N spectrum with $\sim 300,000$ points and multiple free parameters.

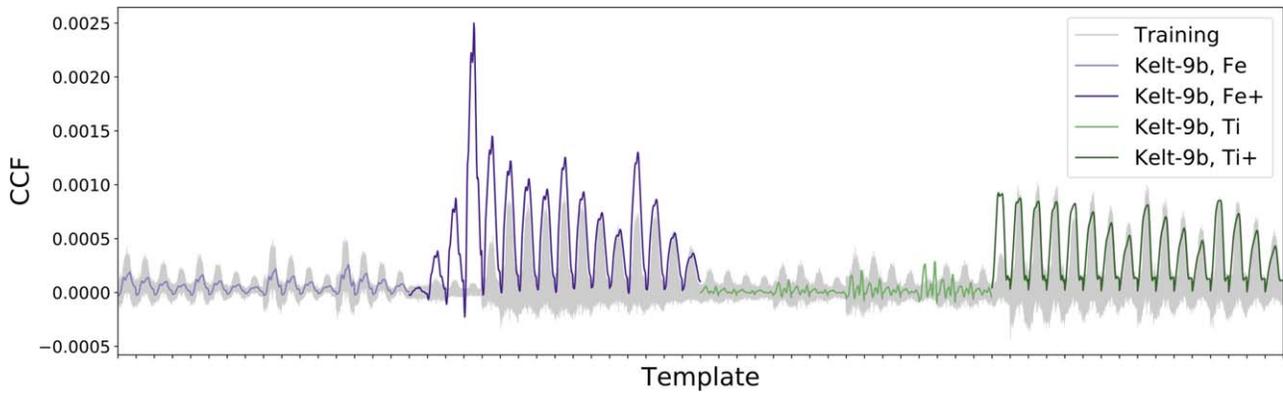


Figure 13. Training vs. measured KELT-9b CCF-sequences. The measured CCF-sequence for Fe^+ lies outside of the range of the model CCF-sequences, thus flagging missing physics in the model grid.

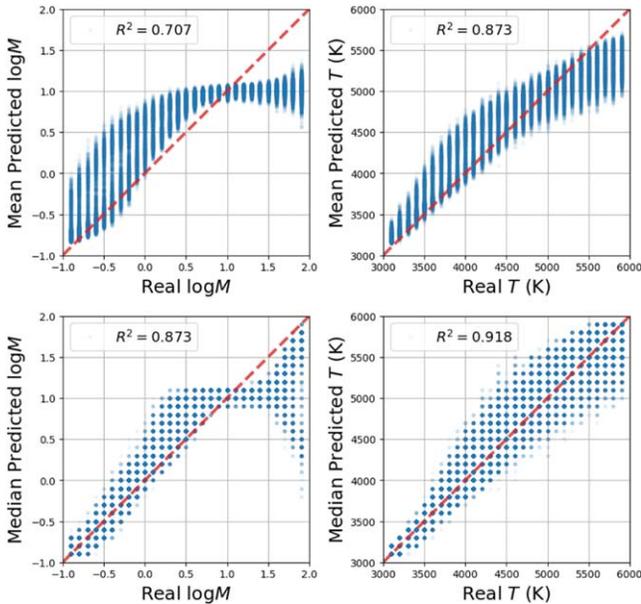


Figure 14. Predicted vs. real values for the forest trained on the CCF-sequences with only neutral species, Fe and Ti. The top and bottom panels show the predictions using the means and medians, respectively. The coefficient of determination (R^2) varies from -1 to 1 , where values near unity indicate strong anticorrelations or correlations between the real and predicted values of a given parameter, based on the variance of the outcomes. The retrieval using this forest on the KELT-9b data is shown as the lighter colored posteriors in Figure 12.

Our method of constructing CCF-sequences allows us to reduce the dimensionality down from $\sim 300,000$ to ~ 2500 . However, now the computational time for each model is much greater as it involves first generating the spectrum and then cross-correlating 64 times with the different templates. Therefore, it remains infeasible to use a standard nested-sampling retrieval for this technique. The random forest requires a grid of precomputed models to train on, allowing the computational burden to be shifted offline. An alternative method using nested sampling could be employed by interpolating on the same grid of models, but without the added noise. There are a few disadvantages involved with this when compared with the forest.

First, the prediction time on a single spectrum is still orders of magnitudes slower than the pretrained forest (~ 20 s versus ~ 0.05 s). This increased computational speed allows the forest to produce “predicted versus real” graphs, as shown in Figure 6 for $\sim 50,000$ models. These graphs give crucial information

about the ability to predict each parameter and the performance of one’s retrieval over a vast range of models. We also gain additional information from the random forest, such as the feature importance plots shown in Figure 8. This quantifies the information content in each spectral point with respect to each parameter being retrieved and can be used to infer which areas of the spectrum are most affected by each parameter. It gives us a deeper insight into how the retrieval works, and even indicates which spectral regions might be most informative when considering future observations.

Second, the use of the likelihood function in nested sampling assumes that the error bars on each spectral point are independent. While this is usually a good assumption, in the process of generating the CCF-sequences we repeatedly cross-correlate a single spectrum with multiple templates and then concatenate these into a sequence. This implies that the noise samples corresponding to each individual cross-correlation cannot be independent as they propagate from the same spectrum. With this assumption broken, it becomes unclear how to proceed with a nested-sampling retrieval on the CCF-sequences.

Third, as discussed in Section 1.3, another assumption one needs to make with nested sampling is a form for the likelihood function, and thus the error bars. For example, it is commonly assumed that the error bars are Gaussians, leading to a likelihood function as shown in Equation (1). The forest also requires an assumption of a random distribution when adding noise to the training set, however it does not depend on a likelihood function. As a test, we generated a model CCF-sequence for a mock retrieval, but this time we added noise by drawing from a Cauchy distribution as opposed to a Gaussian. The motivation behind using a Cauchy distribution is that it does not obey the central limit theorem, and thus the likelihood across many points in a spectrum does not behave as a Gaussian. This provides a challenging test for retrieval methods that assume normally distributed error bars. We performed these retrievals using a forest trained on models with Gaussian errors and a nested-sampling algorithm assuming a Gaussian likelihood function, as shown in Equation (1). Note that this likelihood does not use the cross-correlation function, unlike in Brogi & Line (2019). The results are shown in Figure 15. We can see that while the posteriors are wide for the forest, they still encapsulate the true values, whereas the nested-sampling retrieval produces tightly constricted, incorrect posteriors. This suggests that the forest is more robust to differences in error distributions.

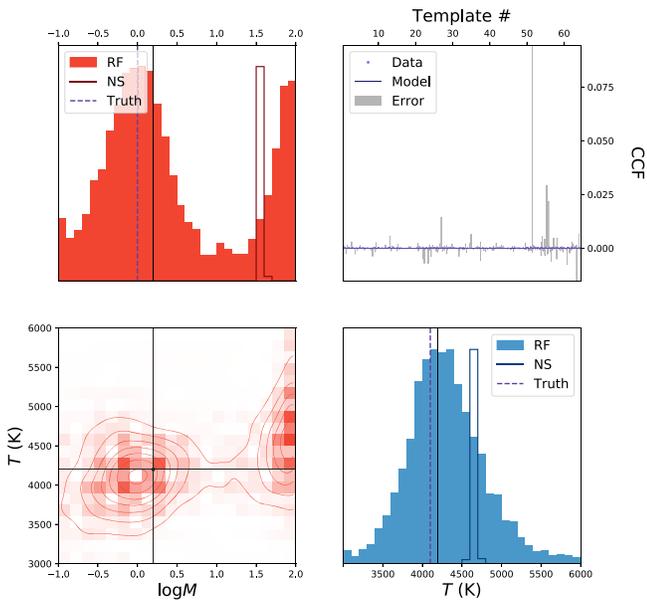


Figure 15. A mock retrieval performed on a model with $\log M = 0$ and $T = 4100$ K, using the CCF-sequence where the noise has been drawn from a Cauchy distribution. The solid posteriors show the random forest (RF) retrieval results, trained on the CCF-sequences with Gaussian noise models (see Figure 6). The empty line posteriors show the nested-sampling (NS) retrieval results using a model that interpolates on the grid of noise-free CCF-sequences and has a Gaussian likelihood. The black lines show the median values of the random forest. The purple, dashed lines show the true values. The top right panel shows the data points (lilac) with the error region (gray), along with the model (dark purple) corresponding to the medians from the $\log M$ and T posteriors.

4.2. Velocity–Velocity Space Performance

So far we have only explored the effects of temperature and metallicity and assumed the velocity parameters, V_{sys} and K_p , are fixed to previously determined values. It is possible that neglecting velocities could lead to severe biases in retrievals. To investigate this, we added the systemic velocity (V_{sys}) and the error in the semiamplitude of the planet radial velocity (ΔK_p) to the method. We took V_{sys} from -10 km s^{-1} to $+10 \text{ km s}^{-1}$ in steps of 2 km s^{-1} , and ΔK_p from 0 km s^{-1} to 60 km s^{-1} in steps of 6 km s^{-1} . Figure A1 shows the results of testing the random forest trained on the CCF-sequences, including the velocity parameters.

This test shows that the addition of the velocity parameters does somewhat reduce the predictive ability of the other parameters, however this reduction is extremely minor for the temperature and not too problematic for the metallicity. The method is able to perfectly retrieve the systemic velocity (V_{sys}), but struggles with the error in semiamplitude of the planet radial velocity (ΔK_p). An error in the assumed value of K_p leads to a misalignment of the planet absorption line when summing in the planet rest frame, effectively resulting in a broadening of the CCF. This makes it more challenging to distinguish between sequences of different metallicities. This explains the greater uncertainty in metallicity that we see in Figure A1 when compared with the results from the method without the velocity parameters (Figure 6).

4.3. Comparison to Other Machine-learning Techniques

There are several other machine-learning methods that can be used to perform atmospheric retrieval (Waldmann 2016; Zingales & Waldmann 2018; Cobb et al. 2019), each with their

own advantages. We tested the same CCF-sequence retrieval as before, but now using a standard neural network and a standard Bayesian neural network (BNN); (Gal 2016). In both cases we used a standard multilayer perceptron architecture with three layers. Each layer consists of a linear transformation with bias followed by a ReLU activation, except the last layer, which does not apply an activation function. The first layer transforms spectra from the input space \mathbb{R}^{2560} to an intermediate representation \mathbb{R}^{512} . Similarly, layer 2 maps elements to \mathbb{R}^{32} , and layer 3 maps elements to the space of parameters \mathbb{R}^2 . The BNN also applies dropout (Srivastava et al. 2014) with probability 0.15 on the output of layers 1 and 2, as explained in Gal (2016). We implemented both networks using the PyTorch library for automatic differentiation (Paszke et al. 2017) and used Adam (Kingma & Ba 2014) as the optimization method.

The results of the test predictions are shown in Figure A2. Compared to the random forest, they both perform with slightly improved R^2 scores. However, this is only a measure of the average prediction. In atmospheric retrieval, we are predominantly interested in the range of possible parameter values given by a retrieval, and therefore the posteriors of each parameter. A traditional neural network does not produce posteriors, so it cannot be meaningfully applied to this retrieval problem. The BNN does provide posteriors, so we are able to compare these to the forest. Figure A3 shows the comparison for two mock retrievals, one with $\log M = 1.0$ and $T = 5100$ K (top panel), and one with $\log M = 1.9$ and $T = 4200$ K (bottom panel). For the first retrieval, the forest and the BNN produce very similar results, with the BNN posteriors slightly tighter and more centered on the true values. However, in the second retrieval the BNN does not perform well for the metallicity prediction. This mock spectrum was selected as one of the retrievals with a strong metallicity degeneracy, as discussed in Section 3.3, in order to test how the two methods deal with these issues. The results for the metallicity prediction are $\log M = 0.3^{+1.7}_{-0.7}$ for the forest, and $\log M = 0.7^{+0.2}_{-0.2}$ for the BNN. Both the average predictions are heavily offset from the correct value, however the posterior from the forest captures the degenerate behavior in metallicity, and therefore encompasses the correct value inside the 1σ interval. In contrast, the BNN posterior sits in the middle of the degenerate peaks and remains tightly constrained around the offset value. It is worth noting that this implementation of the BNN is not equivalent to the one used in Cobb et al. (2019), as they use a different form of the likelihood which has not been tested on such high-resolution data.

4.4. Clarification with Respect to Cobb et al. (2019)

In Cobb et al. (2019), it was suggested that the random forest in Márquez-Neila et al. (2018) has the potential to produce overconfident, incorrect posteriors based on a mock retrieval from a test data set. This forest was trained on WFC3 spectra with 13 data points and predicted 5 parameters—temperature, free chemical abundances of H_2O , HCN , and NH_3 , and a gray cloud opacity, κ_0 . The opacities were calculated with HELIOS-K (Grimm & Heng 2015), using the EXOMOL⁶ (Tennyson et al. 2016) spectroscopic line lists for H_2O (Polyansky et al. 2018), HCN (Barber et al. 2014), and NH_3 (Yurchenko et al. 2011).

⁶ <http://exomol.com>

The mock spectrum tested on by Cobb et al. (2019) has $T = 1479.6$ K, $\log X_{\text{H}_2\text{O}} = -9.79$, $\log X_{\text{HCN}} = -9.04$, $\log X_{\text{NH}_3} = -5.91$, and $\log \kappa_0 = 1.87$. The retrieved posterior for NH_3 was tightly constrained and offset from the correct value, which was used to infer that the forest could produce spurious results. However, we ran the same retrieval with nested sampling, using the same model with the same assumptions. Figure A4 shows the results from the random forest retrieval (left panel) and the nested-sampling retrieval (right panel). The posteriors appear very comparable, with the same behavior in the ammonia abundance.

At the time of publishing Márquez-Neila et al. (2018), there were no opacity linelists available for NH_3 for temperatures above 1500 K. To deal with this, as stated in Márquez-Neila et al. (2018), the opacity for NH_3 was artificially set to zero, and the abundance to the minimum in the range, 10^{-13} . Notable in this particular mock spectrum is the high cloud opacity, equivalent to a cloud top pressure of $\sim 1 \mu\text{bar}$. This results in an essentially flat spectrum. When retrieving on a flat line, the only two parameters in this model having an effect are the temperature and the cloud opacity, which are perfectly degenerate with each other (i.e., an increase in either results in an upwards shift of the line, so by decreasing the other, one obtains the same spectrum). This degeneracy means one can only obtain lower bounds for the temperature and cloud opacity, corresponding to the upper bound of the other parameter's prior. A consequence of this is a collection of posterior samples in the region $T > 1500$ K, which, as forced by the model, have $\log X_{\text{NH}_3} = -13$, resulting in the peaked posterior for NH_3 . Therefore, this offset posterior is actually an artifact of the training set rather than of the random forest. This is shown conclusively in Figure A4, as the forest's posteriors agree with the true Bayesian posteriors from nested sampling.

5. Conclusion

This paper presents a new method for performing atmospheric retrieval on ground-based, high-resolution spectroscopic observations of exoplanets. By using a combination of CCFs we are able to reduce the dimensionality of the problem and decrease the high levels of uncertainty on each spectral data point. Using our previously demonstrated random forest retrieval technique (Márquez-Neila et al. 2018), we can execute the retrieval quickly and run a multitude of tests of the method. These show that the method performs well on mock observations, with a high predictive power for metallicity and temperature ($R^2 = 0.918$ and 0.986 , respectively). The random forest also provides feature importance plots, which show that the neutral cross-correlations are most important for determining the metallicity, while the temperature prediction relies predominantly on the ions. Our method also highlights the metallicity degeneracy in the model, which accounts for the reduced predictability at high metallicity values.

We performed the retrieval on HARPS-N observations for the ultra-hot Jupiter KELT-9b. The metallicity appears to be consistent with solar, with the retrieval seemingly driven by the neutral species. The prediction for temperature is forced up to exceptionally high values, due to excess Fe^+ absorption that appears in the high-resolution transmission spectrum, suggesting the need for more complex physics in the model. This can be seen when comparing the data to the training set, which also implies that this method is able to recognize when the model is incomplete.

We also compared the use of our random forest to other approaches, such as the traditional nested-sampling technique and other machine-learning methods. We showed that the forest is more robust to the use of different error distributions than nested sampling, due to it being likelihood free. When compared with a BNN, although the BNN obtains marginally improved R^2 scores, only the forest was able to produce complex posteriors, e.g., in the case of degenerate metallicity values. We also demonstrated that the claim in Cobb et al. (2019), that the forest can be overconfident but incorrect, is actually an outcome of the atmospheric model itself and that the forest's posteriors agree with the results from nested sampling.

We thank Heather Knutson for suggesting the Cauchy distribution test and Ewan Cameron for stimulating discussions on the random forest. We acknowledge financial support from the Swiss National Science Foundation, the European Research Council (via a Consolidator Grant to K.H.; grant No. 771620), the PlanetS National Center of Competence in Research (NCCR), the Center for Space and Habitability (CSH), the Swiss-based MERAC Foundation, and the University of Bern International 2021 PhD Fellowship.

Software: FastChem (Stock et al. 2018), HELIOS-K (Grimm & Heng 2015), Helios-o (Bower et al. 2019), scikit.learn (Pedregosa et al. 2011), PyTorch (Paszke et al. 2017), Adam (Kingma & Ba 2014), PyMultinest (Buchner et al. 2014), Astropy (The Astropy Collaboration et al. 2018), numpy (Van Der Walt et al. 2011), scipy (Virtanen et al. 2020), matplotlib (Hunter 2007).

Appendix Additional Figures

Figure A1 shows the results for the forest including the velocity parameters, as discussed in Section 4.2. Figures A2 and A3 show the results for the neural network and Bayesian neural network trained on the CCF-sequences, as discussed in Section 4.3. Figure A4 shows the comparison between a random forest and a nested sampling retrieval for the model and spectrum considered in Cobb et al. (2019), as discussed in Section 4.4.

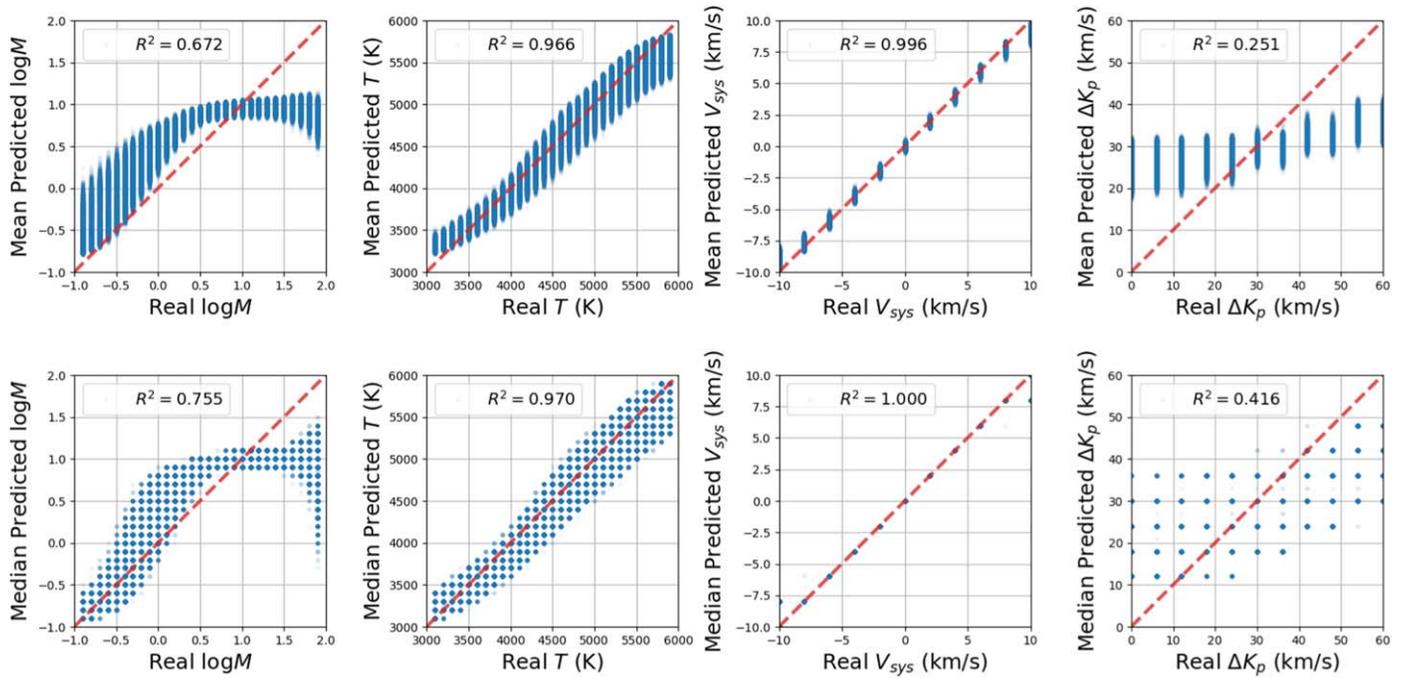


Figure A1. Predicted vs. real values of the logarithm of metallicity ($\log M$), temperature (T), systemic velocity (V_{sys}), and error in semiamplitude of the planet radial velocity (ΔK_p) for the random forest trained on the CCF-sequences. The top and bottom rows show the predictions using the means and medians, respectively. The coefficient of determination (R^2) varies from -1 to 1 , where values near unity indicate strong anticorrelations or correlations between the real and predicted values of a given parameter, based on the variance of the outcomes.

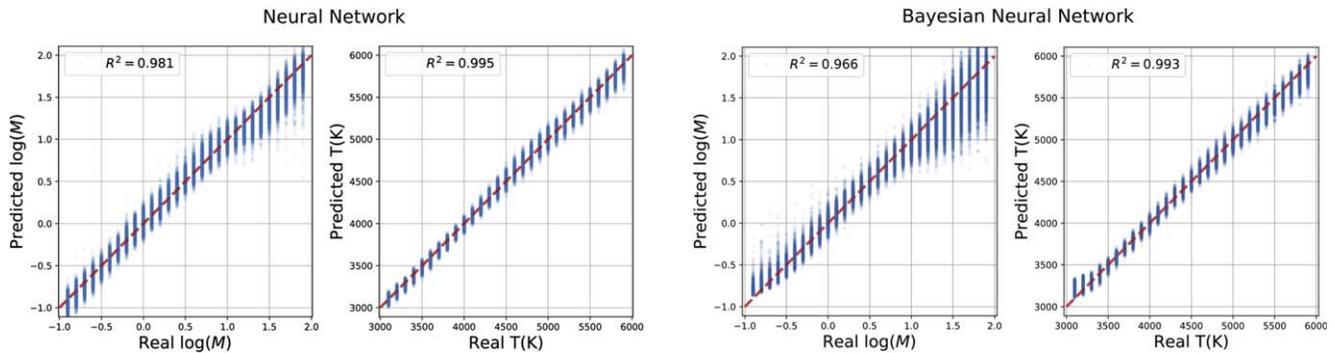


Figure A2. Predicted vs. real values for neural networks trained on the CCF-sequences. The left and right pairs show the results for a standard neural network and a Bayesian neural network (BNN), respectively. The coefficient of determination (R^2) varies from -1 to 1 , where values near unity indicate strong anticorrelations or correlations between the real and predicted values of a given parameter, based on the variance of the outcomes. Mock retrievals for the BNN are shown as the empty line posteriors in Figure A3.

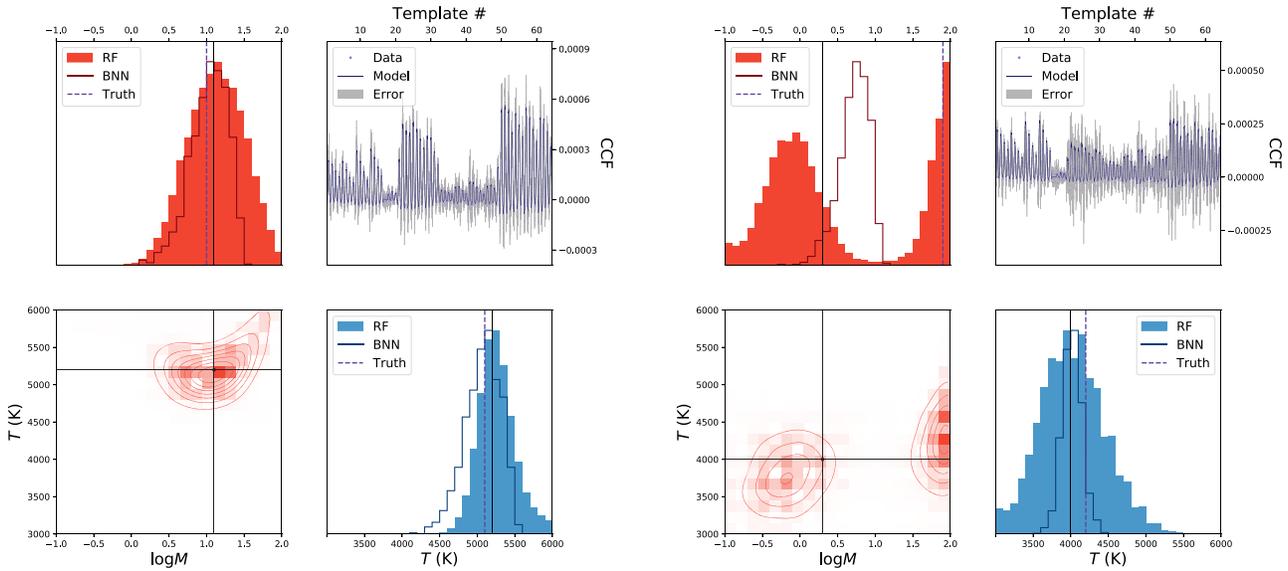


Figure A3. Two mock retrievals performed using the CCF-sequences. The solid posteriors show the random forest retrieval results (see Figure 6). The empty line posteriors show the Bayesian neural network retrieval results (see bottom panels of Figure A2). The black lines show the median values of the random forest, and the purple, dashed lines show the true values. The left figure corresponds to a retrieval on a model with $\log M = 1.0$ and $T = 5100$ K. The right figure corresponds to a retrieval on a model with $\log M = 1.9$ and $T = 4200$ K.

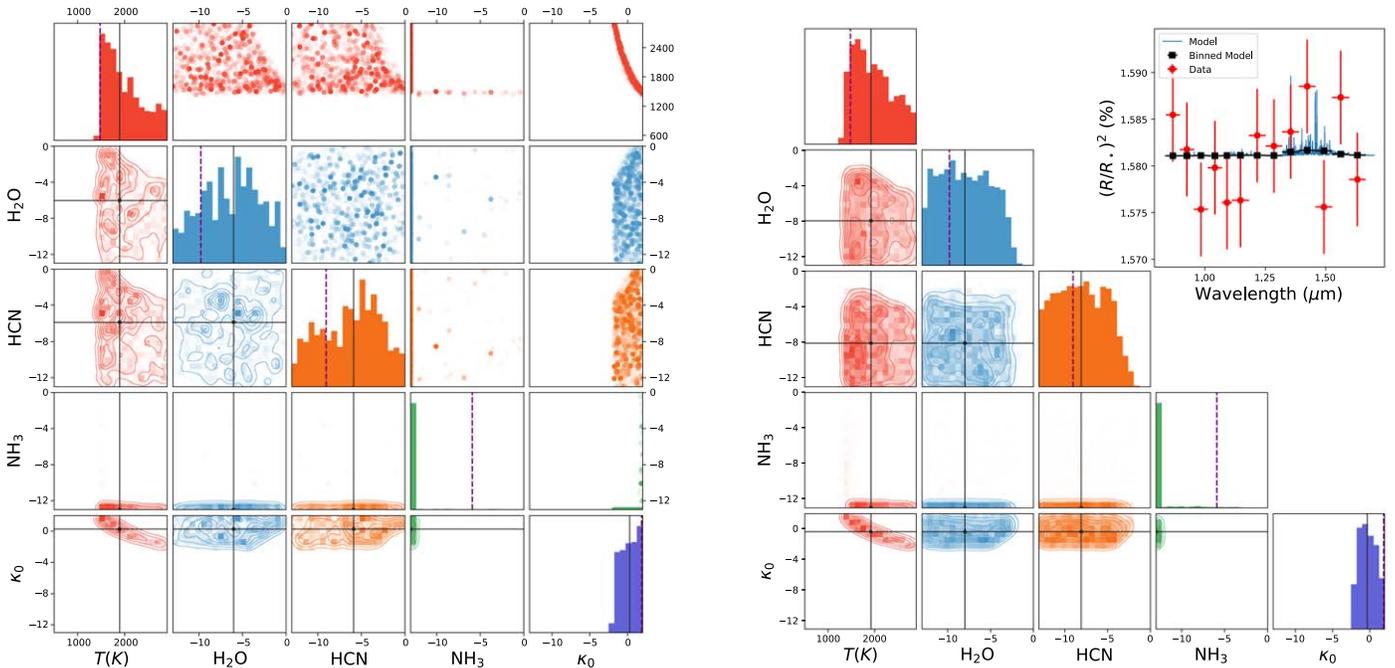


Figure A4. Retrieval results for a mock spectrum with $T = 1479.6$ K, $\log X_{\text{H}_2\text{O}} = -9.79$, $\log X_{\text{HCN}} = -9.04$, $\log X_{\text{NH}_3} = -5.91$, and $\log \kappa_0 = 1.87$. The left- and right-hand plots show the results using a random forest and nested sampling, respectively. The black lines show the median predictions. The purple, dashed lines show the true values.

ORCID iDs

Chloe Fisher <https://orcid.org/0000-0003-0652-2902>
 Daniel Kitzmann <https://orcid.org/0000-0003-4269-3311>
 Simon L. Grimm <https://orcid.org/0000-0002-0632-4407>
 Kevin Heng <https://orcid.org/0000-0003-1907-5910>

References

Arcangeli, J., Désert, J.-M., Line, M. R., et al. 2018, *ApJL*, 855, L30
 Barber, R. J., Strange, J. K., Hill, C., et al. 2014, *MNRAS*, 437, 1828
 Benneke, B., & Seager, S. 2012, *ApJ*, 753, 100
 Birkby, J. L., de Kok, R. J., Brogi, M., et al. 2013, *MNRASL*, 436, L35

Birkby, J. L., de Kok, R. J., Brogi, M., et al. 2017, *ApJ*, 153, 138
 Bower, D. J., Kitzmann, D., Wolf, A. S., et al. 2019, *A&A*, 631, A103
 Breiman, L. 2001, *Machine Learning*, 45, 5
 Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. 1984. Classification and Regression Trees (Boca Raton, FL: Chapman & Hall/CRC Press)
 Brogi, M., de Kok, R. J., Birkby, J. L., et al. 2014, *A&A*, 565, A124
 Brogi, M., Giacobbe, P., Guilluy, G., et al. 2018, *A&A*, 615, A16
 Brogi, M., Line, M., Bean, J., et al. 2017, *ApJL*, 839, L2
 Brogi, M., & Line, M. R. 2019, *AJ*, 157, 114
 Brogi, M., Snellen, I. A. G., de Kok, R. J., et al. 2012, *Natur*, 486, 502
 Brogi, M., Snellen, I. A. G., de Kok, R. J., et al. 2013, *ApJ*, 767, 27
 Brown, T. M. 2001, *ApJ*, 553, 1006
 Brown, T. M., Libbrecht, K. G., & Charbonneau, D. 2002, *PASP*, 114, 826
 Buchner, J., Georgakakis, A., Nandra, K., et al. 2014, *A&A*, 564, A125

- Cauley, W., Shkolnik, E. L., Ilyin, I., et al. 2019, *AJ*, **157**, 69
- Cobb, A. D., Himes, M. D., Soboczenski, F., et al. 2019, *AJ*, **158**, 33
- Cosentino, R., Lovis, C., Pepe, F., et al. 2012, *Proc. SPIE*, **8446**, 84461V
- Criminisi, A., Shotton, J., & Konukoglu, E. 2011, Decision Forests for Classification, Regression, Density Estimation, Manifold Learning and Semi-Supervised Learning, Microsoft Research Technical Report, TR-2011-114
- de Kok, R. J., Brogi, M., Snellen, I. A. G., et al. 2013, *A&A*, **554**, A82
- Dekker, H., D'Odorico, S., Kaufer, A., et al. 2000, *Proc. SPIE*, **4008**, 534
- Deming, D., Brown, T. M., Charbonneau, D., et al. 2005, *ApJ*, **622**, 1149
- Deming, D., Wilkins, A., McCullough, P., et al. 2013, *ApJ*, **774**, 95
- Donati, J.-F., Kouach, D., Lacombe, M., et al. 2018, in Handbook of Exoplanets, ed. H. Deeg & J. Belmonte (Cham: Springer), 107
- Draine, B. T. 2011, Physics of the Interstellar and Intergalactic Medium (Princeton, NJ: Princeton Univ. Press)
- Feroz, F., & Hobson, M. P. 2008, *MNRAS*, **384**, 449
- Feroz, F., Hobson, M. P., & Bridges, M. 2009, *MNRAS*, **398**, 1601
- Feroz, F., Hobson, M. P., Cameron, E., & Pettitt, A. N. 2019, *OJAp*, **2**, 19
- Fisher, C., & Heng, K. 2018, *MNRAS*, **481**, 4698
- Follert, R., Dorn, R. J., Oliva, E., et al. 2014, *Proc. SPIE*, **9147**, 19
- Fortney, J., Lodders, K., Marley, M., & Freedman, R. 2008, *ApJ*, **678**, 1419
- Fortney, J., Shabram, M., Showman, A. P., et al. 2010, *ApJ*, **709**, 1396
- Gaidos, E., Kitzmann, D., & Heng, K. 2017, *MNRAS*, **468**, 3418
- Gal, Y. 2016, PhD thesis, Univ. Cambridge
- Gaudi, B. S., Stassun, K. G., Collins, K. A., et al. 2017, *Natur*, **546**, 514
- Gibson, N. P., Merritt, S., Nugroho, S. K., et al. 2020, *MNRAS*, **493**, 2215
- Grimm, S. L., & Heng, K. 2015, *ApJ*, **808**, 182
- Guilluy, G., Sozzetti, A., Brogi, M., et al. 2019, *A&A*, **625**, A107
- Hastie, T., Tibshirani, R., & Friedman, J. 2001, The Elements of Statistical Learning (New York: Springer)
- Heng, K. 2017, Exoplanetary Atmospheres: Theoretical Concepts and Foundations (Princeton, NJ: Princeton Univ. Press)
- Heng, K., & Kitzmann, D. 2017, *MNRAS*, **470**, 2972
- Ho, T. K. 1998, *ITPAM*, **20**, 832
- Hoeijmakers, H. J., de Kok, R. J., Snellen, I. A. G., et al. 2015, *A&A*, **575**, A20
- Hoeijmakers, H. J., Ehrenreich, D., Heng, K., et al. 2018, *Natur*, **560**, 453
- Hoeijmakers, H. J., Ehrenreich, D., Kitzmann, D., et al. 2019, *A&A*, **627**, A165
- Hubeny, I., Burrows, A., & Sudarsky, D. 2003, *ApJ*, **594**, 1011
- Hunter, J. D. 2007, *CSE*, **9**, 90
- John, T. L. 1988, *A&A*, **193**, 189
- Kaeufl, H.-U., Ballester, P., Biereichel, P., et al. 2004, *SPIE*, **5492**, 1218
- Khalafinejad, S., Salz, M., & Cubillos, P. E. 2018, *A&A*, **618**, A98
- Khalafinejad, S., von Essen, C., & Hoeijmakers, H. J. 2017, *A&A*, **598**, A131
- Kingma, D. P., & Ba, J. 2014, arXiv:1412.6980
- Kitzmann, D., Heng, K., Rimmer, P. B., et al. 2018, *ApJ*, **863**, 183
- Kurucz, R. L. 2017, *CajPh*, **95**, 825
- Lavie, B., Mendonça, J. M., Mordasini, C., et al. 2017, *AJ*, **154**, 91
- Line, M. R., Knutson, H., Deming, D., et al. 2013, *ApJ*, **778**, 183
- Line, M. R., Wolf, A. S., Zhang, X., et al. 2013, *ApJ*, **775**, 137
- Lockwood, A. C., Johnson, J. A., Bender, C. F., et al. 2014, *ApJL*, **783**, L29
- MacDonald, R. J., & Madhusudhan, N. 2017, *MNRAS*, **469**, 1979
- Márquez-Neila, P., Fisher, C., Sznitman, R., & Heng, K. 2018, *NatAs*, **2**, 719
- Mayor, M., Pepe, F., Queloz, D., et al. 2003, *Msngr*, **114**, 20
- McLean, I. S., Becklin, E. E., Bendiksen, O., et al. 1998, *Proc. SPIE*, **3354**, 566
- Noguchi, K., Aoki, W., Kawanomoto, S., et al. 2002, *PASJ*, **54**, 855
- Nugroho, S. K., Kawahara, H., Masuda, K., et al. 2017, *AJ*, **154**, 221
- Origlia, L., Oliva, E., Baffa, C., et al. 2014, *Proc. SPIE*, **9147**, 91471E
- Park, C., Jaffe, D. T., Yuk, I.-S., et al. 2014, *Proc. SPIE*, **9147**, 91471D
- Paszke, A., Gross, S., Chintala, S., et al. 2017, Automatic differentiation in PyTorch, <https://openreview.net/forum?id=BJJsmfCZ>
- Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, *JMLR*, **12**, 2825
- Pepe, F., Molaro, P., Cristiani, S., et al. 2014, *AN*, **335**, 8
- Piskorz, D., Benneke, B., Crockett, N. R., et al. 2016, *ApJ*, **832**, 131
- Piskorz, D., Benneke, B., Crockett, N. R., et al. 2017, *AJ*, **154**, 78
- Piskorz, D., Buzard, C., Line, M. R., et al. 2018, *AJ*, **156**, 133
- Polyansky, O. L., Kyuberis, A. A., Zobov, N. F., et al. 2018, *MNRAS*, **480**, 2597
- Quirrenbach, A., Amado, P. J., Mandel, H., et al. 2010, *Proc. SPIE*, **7735**, 773513
- Rayner, J., Tokunaga, A., Jaffe, D., et al. 2016, *Proc. SPIE*, **9908**, 990884
- Redfield, S., Endl, M., Cochran, W. D., & Koesterke, L. 2008, *ApJL*, **673**, L87
- Richard, C., Gordon, I. E., Rothman, L. S., et al. 2012, *JQSRT*, **113**, 1276
- Rieke, N., Tan, D. J., Alshekhali, M., et al. 2015, in Medical Image Computing and Computer-assisted Intervention—MICCAI 2015, ed. N. Navab et al. (Cham: Springer), 266
- Rothman, L. S., Rinsland, C. P., Goldman, A., et al. 1998, *JQSRT*, **60**, 665
- Saha, M. N. 1920, *PMag*, **40**, 472
- Seidel, J. V., Ehrenreich, D., Wyttenbach, A., et al. 2019, *A&A*, **623**, A166
- Sisson, S. A., Fan, Y., & Beaumont, M. A. 2019, Handbook of Approximate Bayesian Computation (Boca Raton, FL: CRC Press)
- Skilling, J. 2006, *BayAn*, **1**, 833
- Snellen, I. A. G., Albrecht, S., de Mooij, E. J. W., & Le Poole, R. S. 2008, *A&A*, **487**, 357
- Snellen, I. A. G., de Kok, R. J., de Mooij, E. J. W., & Albrecht, S. 2010, *Natur*, **465**, 1049
- Srivastava, N., Hinton, G., Krizhevsky, A., et al. 2014, *JMLR*, **15**, 1929
- Stock, J. W., Kitzmann, D., Patzer, A. B. C., & Sedlmayr, E. 2018, *MNRAS*, **479**, 865
- Strassmeier, K. G., Ilyin, I., Järvinen, A., et al. 2015, *AN*, **336**, 324
- Sznitman, R., Becker, C., Fleuret, F., & Fua, P. 2013, in 2013 IEEE Conf. on Computer Vision and Pattern Recognition (Piscataway, NJ: IEEE), 3270
- Tennyson, J., Yurchenko, S. N., Al-Refaie, A. F., et al. 2016, *JMoSp*, **327**, 73
- The Astropy Collaboration, Price-Whelan, A. M., Sipöcz, B. M., et al. 2018, *AJ*, **156**, 123
- Van Der Walt, S., Colbert, S. C., & Varoquaux, G. 2011, *CSE*, **13**, 22
- Virtanen, P., Gommers, R., & Oliphant, T. E. 2020, *Nature Methods*, **17**, 261
- Waldmann, I. P. 2016, *ApJ*, **820**, 107
- Waldmann, I. P., Tinetti, G., Rocchetto, M., et al. 2015, *ApJ*, **802**, 107
- Wiedemann, G., Deming, D., & Bjoraker, G. 2001, *ApJ*, **546**, 1068
- Wildi, F., Blind, N., Reshetov, V., et al. 2017, *Proc SPIE*, **10400**, 1040018
- Wong, I., Shporer, A., Morris, B. M., et al. 2019, arXiv:1910.01607
- Wyttenbach, A., Ehrenreich, D., Lovis, C., et al. 2015, *A&A*, **577**, A62
- Wyttenbach, A., Lovis, C., Ehrenreich, D., et al. 2017, *A&A*, **602**, A36
- Yurchenko, S. N., Barber, R. J., & Tennyson, J. 2011, *MNRAS*, **413**, 1828
- Zerbi, F. M., Bouchy, F., Fynbo, J., et al. 2014, *Proc. SPIE*, **9147**, 914723
- Zhang, L., Varadarajan, J., Suganthan, P., et al. 2017, in 2017 IEEE Conf. on Computer Vision and Pattern Recognition (Piscataway, NJ: IEEE), 5825
- Zikic, D., Glocker, B., & Criminisi, A. 2014, *Medical Image Analysis*, **18**, 1262
- Zingales, T., & Waldmann, I. P. 2018, *AJ*, **156**, 268