

Tau-REx I: A next generation retrieval code for exoplanetary atmospheres

I. P. Waldmann, G. Tinetti, M. Rocchetto, E. J. Barton, S. N. Yurchenko, J. Tennyson
Department of Physics & Astronomy, University College London, Gower Street, WC1E 6BT, UK
 ingo@star.ucl.ac.uk

ABSTRACT

Spectroscopy of exoplanetary atmospheres has become a well established method for the characterisation of extrasolar planets. We here present a novel inverse retrieval code for exoplanetary atmospheres. \mathcal{T} -REx (Tau Retrieval for Exoplanets) is a line-by-line radiative transfer fully Bayesian retrieval framework. \mathcal{T} -REx includes the following features: 1) the optimised use of molecular line-lists from the *ExoMol* project; 2) an unbiased atmospheric composition prior selection, through custom built pattern recognition software; 3) the use of two independent algorithms to fully sample the Bayesian likelihood space: nested sampling as well as a more classical Markov Chain Monte Carlo approach; 4) iterative Bayesian parameter and model selection using the full Bayesian Evidence as well as the Savage-Dickey Ratio for nested models, and 5) the ability to fully map very large parameter spaces through optimal code parallelisation and scalability to cluster computing. In this publication we outline the \mathcal{T} -REx framework and demonstrate, using a theoretical hot-Jupiter transmission spectrum, the parameter retrieval and model selection. We investigate the impact of Signal-to-Noise and spectral resolution on the retrievability of individual model parameters, both in terms of error bars on the temperature and molecular mixing ratios as well as its effect on the model's global Bayesian evidence.

Subject headings: methods: data analysis — methods: statistical — techniques: spectroscopic — radiative transfer

1. Introduction

Remote sensing of atmospheres and inverse retrieval methods have a well established and long standing history. Beginning with pioneering work on our own Earth (e.g. Wark & Hilleary 1969; Conrath et al. 1970), we quickly extended our grasp to other planets in our solar system (e.g. Hanel et al. 1972; Conrath et al. 1973; Rodgers 1976; Hanel et al. 1981). With the first detection of exoplanetary atmospheres (Charbonneau et al. 2002) we have taken this work beyond our solar system confines.

In recent years, the field of extrasolar spectroscopy has seen a increased effort in the development of data analysis and de-trending techniques (e.g. Swain et al. 2008; Carter & Winn 2009; Burke et al. 2010; Snellen et al. 2010; Thatte et al. 2010;

Swain et al. 2010; Waldmann et al. 2012, 2013; Waldmann 2012, 2014; Gibson et al. 2012; Crouzet et al. 2012; Berta et al. 2012; Morello et al. 2014; Danielski et al. 2014; Kreidberg et al. 2014). With the maturation of these methodologies, we are obtaining a rapidly increasing number of exoplanetary emission and transmission spectra requiring interpretation. We refer the reader to Seager (2011) and Tinetti et al. (2013) and references within, for reviews of current spectroscopic results.

This ever increasing wealth of spectroscopic data of extrasolar planet atmospheres allows an unprecedented insight into the properties of these foreign worlds.

The interpretation of atmospheric spectra of extrasolar planets through *inverse atmospheric retrieval* modelling (e.g. Fletcher et al. 2007; Ter-rile et al. 2008; Irwin et al. 2008; Madhusudhan

& Seager 2009; Lee et al. 2011; Line et al. 2012; Benneke & Seager 2012; Barstow et al. 2013; Griffith 2014) has become the industry standard. Line et al. (2013b) provides a recent and comprehensive review of currently existing exoplanetary atmospheric retrieval codes.

With greater accuracy in data often comes an increased complexity in its interpretation. In analogy to recent challenges in observational exoplanetary data analysis, one can identify three major objectives for the interpretation of exoplanetary spectra:

Sensitivity: Given the often low resolution and low signal-to-noise (S/N) of currently available exoplanetary spectra, an understanding of the limitations and degeneracies of spectroscopic models is paramount.

Objectivity: Are the results driven by model dependencies, over-constraint inputs or human biases? An idealised atmospheric retrieval should make no prior assumptions about the complex nature of exoplanetary atmospheres. Whilst this is often infeasible, modern retrieval algorithms should be designed to take into account the broadest possible range of atmospheric models. It should then select amongst these models using a consistent and quantifiable metric of parameter and model adequacies.

Big data: With the increasing automation of exoplanet observations, the manual interpretation of atmospheric spectra will become infeasible. A modern retrieval algorithm should bear this in mind and allow for a high degree of intelligent automation and scalability to larger cluster computing.

In this paper, we introduce a new atmospheric retrieval code, \mathcal{T} -REx (Tau Retrieval for Exoplanets), which has been designed with the above objectives in mind. Here we will describe the overall architecture and atmospheric retrieval for transmission spectroscopy and dedicate a subsequent publication (Waldmann et al. in prep.) to the emission/reflection spectroscopy case and the parameterisation of the temperature-pressure (T-P) profile.

1.1. \mathcal{T} -REx

\mathcal{T} -REx is a novel, fully Bayesian, retrieval code for exoplanetary atmospheres. In its current im-

plementation, \mathcal{T} -REx includes the following features:

Line-by-line: \mathcal{T} -REx uses customised molecular and atomic line lists available directly from the *ExoMol*¹ project (Tennyson & Yurchenko 2012). In particular, ExoMol provides computed line lists valid over extended temperature ranges for a variety of molecules including water (Barber et al. 2006), ammonia (Yurchenko et al. 2011), methane Yurchenko & Tennyson (2014) and a variety of diatomic molecules (Yadin et al. 2012; Barton et al. 2013; Barber et al. 2014; Barton et al. 2014). Besides line lists *ExoMol* provides cross sections (Hill et al. 2013) for *ExoMol* and other line lists. In this work cross sections for CO, NO and CO₂ created from *HITEMP* (Rothman et al. 2010) and TiO from Schwenke (Schwenke 1998) are also used. Molecular (and atomic) absorption cross-sections are calculated on an optimal linear or non-linear wavelength grid resulting in a optimally sparse cross-section library with fine gridding of optically thick lines. This guarantees high computational efficiency without loss of accuracy.

Non-parametric prior constraint: The priors to the Bayesian retrieval such as number and type of molecules considered, abundance and temperature ranges are not manually set by individual users but automatically determined by \mathcal{T} -REx based on the probability of individual molecules being present in the exoplanetary spectrum. The `Marple` module is a custom built pattern recognition package capable of rapidly identifying likely absorbers/emitters in the exoplanetary spectra from large line-list archives (e.g. *ExoMol* (Tennyson & Yurchenko 2012), *HITRAN*² (Rothman et al. 2009, 2013) and *HITEMP* (Rothman et al. 2010)). Such an approach is highly efficient as a very large number of molecules/atoms/ions can be considered and minimises the human bias in selecting ‘key atmospheric components’.

Bayesian Model Selection: The code can be run to integrate over the full likelihood space of the Bayesian argument allowing for the Bayesian partition function, also called the Bayesian Evidence, to be calculated. This allows for the posterior distributions of model parameters to be calculated, as well as the adequacy of the model itself, given

¹<http://www.exomol.com>

²<http://www.cfa.harvard.edu/hitran/>

the data, to be assessed and iteratively optimised.

Scaleability: Even the simplest retrieval cases can feature a high dimensional likelihood space. By relying on nested sampling approaches, we can naturally achieve an excellent multi-core processor scalability and full parallelisation of the code, allowing us to fully map possible correlations in the likelihood space.

2. Code overview

The retrieval code discussed here is based on a fully modular, object oriented architecture. Amongst others, one of the main advantages of a modular approach is a more flexible and structured approach to complex programs. Throughout this paper we will follow this modular approach in describing \mathcal{T} -REx’s individual components.

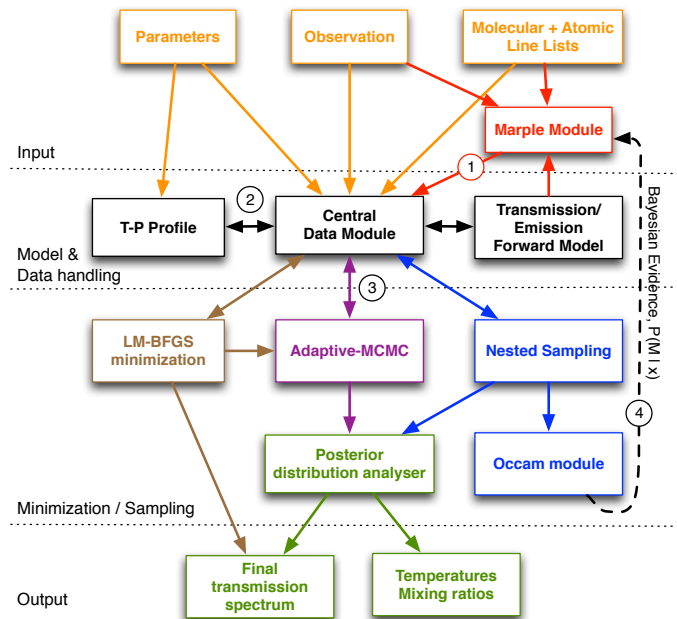


Fig. 1.— Flowchart illustrating the modular design of \mathcal{T} -REx. As described in the text, \mathcal{T} -REx is subdivided into four main parts: Input, Model and data handling, Retrieval/minimization and Output analysis.

The \mathcal{T} -REx design is illustrated in figure 1 and can be broadly subdivided into four main programmatic segments:

1. *Inputs* (sections 3 & 5) - these include global

parameters, the observed exoplanetary spectrum and the absorption cross-sections for the range of molecules to be considered. The **Marple** module (point 1 in figure 1) acts as extra input to the main code, providing a best initial guess of the atmospheric composition of the extrasolar planet to \mathcal{T} -REx.

2. *Model and Data handling* - defines the radiative transfer forward model, the temperature-pressure (TP) profile, overall input data handling. The Central Data Module (point 2 in figure 1) acts as an abstraction layer between model, minimisation and data, allowing for an easy interchangeability of models, minimisation/retrieval techniques and data types.

3. *Retrieval* (sections 6) - contains three minimisation codes: 1) Limited-memory Broyden-Fletcher-Goldfarb-Shannon (LM-BFGS) algorithm, 2) Adaptive Markov Chain Monte Carlo (MCMC) sampling, 3) Bayesian Nested Sampling (NS). For NS runs, the Bayesian Partition function is calculated and model selection is performed. Point 3 in figure 1 illustrates that all these algorithms have a common standardised interface with the Central Data Module. This guarantees exact and comparable results for different retrieval techniques.

The **Occam** module (section 7) performs Bayesian model selection on the outputs of the MCMC and NS. In the case of under or over-complete models, the **Occam** module updates the planetary transmission model in an iterative manner (point 4 in figure 1).

4. *Output* (section 8) - the final exoplanetary spectrum is returned along with all parameter posterior distributions, cross-correlations and Bayesian Evidence.

3. Atomic and Molecular Line-lists

The optimal treatment of atomic and molecular line-lists is key to the accuracy achieved by \mathcal{T} -REx. Throughout the code we perform line-by-line radiative transfer calculations at typically 50-100 times higher spectral resolution than the resolution of the observed spectrum to be analysed. These ‘high-resolution’ spectra are binned

down (at each iteration of the code) to the data to calculate the χ^2 of the model fit. This ensures a correct treatment of optically thick absorption regions. In future versions of the code, we plan to include optimal non-linear binning of line-lists to further increase the computational efficiency without impacting model accuracies (Barton et al. in prep).

\mathcal{T} -REx allows for an easy and seamless inclusion of large numbers of line lists. For the scope of this paper we limit ourselves to line-lists from absorption cross-sections obtained from *ExoMol* but *HITRAN* line-lists (or a combination of both) are equally natively supported. Automated pre-processing steps allow for conversions to a uniform data format with cross sections typically at $\Delta\nu = 1.0 \text{ cm}^{-1}$ resolution. The cross-section library is generated at temperature intervals of 100K (Hill et al. 2013) and upon execution of the main code interpolated to a user-set temperature resolution (typically 10K). Two forms of cross section interpolation are available: 1) linear and 2) optimal. For the optimal case, we follow Hill et al. (2013) where the temperature interpolated cross-section $\varsigma_{m,\lambda}(T)$ is given by

$$\varsigma_{m,\lambda}(T) = a_{m,\lambda} e^{-b_{m,\lambda}/T} \quad (1)$$

where m is the molecular/atomic species index, λ the wavelength, T the final temperature and a and b are scaling factors given by

$$b_{m,\lambda} = \left(\frac{1}{T_2} - \frac{1}{T_1} \right)^{-1} \ln \frac{\varsigma_{m,\lambda}(T_1)}{\varsigma_{m,\lambda}(T_2)} \quad (2)$$

$$a_{m,\lambda} = \varsigma_{m,\lambda}(T_1) e^{b_{m,\lambda}/T_1} \quad (3)$$

where T_1 and T_2 are upper and lower temperatures respectively.

4. Forward Model

The transmission forward model is based on the **Tau** code by Hollis et al. (2013) but was optimised for a significantly higher computational efficiency. We will only give a brief summary of the transmission model and refer the interested reader to the relevant literature (e.g. Brown 2001; Liou 2002; Tinetti et al. 2012; Hollis et al. 2013). As previously mentioned, in this paper we will only describe the transmission part of \mathcal{T} -REx and

an isothermal temperature-pressure (T-P) profile. We dedicate a second publication (Waldmann et al. in prep.) to a complete treatment of the emission case and T-P profile parametrisation.

The monochromatic intensity, $I_\lambda(z)$, of radiation passing through a gas is given by the *Beer-Bouguer-Lambert Law* as function of atmospheric altitude, z ,

$$I_\lambda(z) = I_\lambda(0) e^{-\tau_\lambda(z)} \quad (4)$$

where λ is the wavelength of the radiation, $I_\lambda(0)$ the incident radiation intensity at the top of the atmosphere and $\tau_\lambda(z)$ the optical depth of the medium. For a given absorber, m , we can state the optical depth to be the integral of the absorption cross-section, $\varsigma_m(\lambda)$, the column density, $\chi_m(z)$ and the number density, $\rho_N(z)$, over the optical path length $l(z)$

$$\tau_{\lambda,m}(z) = 2 \int_0^{l(z)} \varsigma_m(\lambda) \chi_m(z) \rho_N(z) dl \quad (5)$$

where the path length is dependent on the geometry of the transmission through the planet's terminator (see figure 2 in Hollis et al. 2013). The overall optical depth is now given by the sum of the individual optical depths

$$\tau_\lambda(z) = \sum_{m=1}^{N_m} \tau_{\lambda,m}(z) \quad (6)$$

where N_m is the total number of absorbing species, m . We can now calculate the equivalent atmospheric depth, α_λ , by summing over all atmospheric depth layers, z ,

$$\alpha_\lambda = 2 \int_0^{z_{max}} (R_p + z) (1 - e^{-\tau_\lambda(z)}) dz \quad (7)$$

where z_{max} is the maximum depth of the atmosphere considered. The total transit depth as a function of λ is hence given by

$$\delta_\lambda = \frac{R_p^2 + \alpha_\lambda}{R_*^2} \quad (8)$$

where R_p and R_* are the radii of the planet and star respectively.

\mathcal{T} -REx provides a full implementation of Rayleigh and Mie scattered as well as cloudy atmospheres. We refer the reader to Hollis et al. (2013) for details of implementation. Retrieval degeneracies due to cloud models will be discussed in a separate publication (Rocchetto et al., in prep.).

5. Marple module

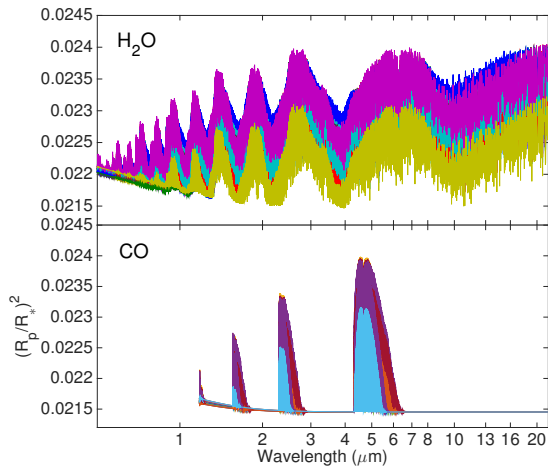


Fig. 2.— Transmission spectra of H₂O (top) and CO (bottom) for temperature and abundance ranges of 600-2000K and 1×10^{-5} - 1×10^{-2} respectively. Planet/star and orbital parameters are taken to be similar to hot-Jupiter HD 209458b.

The purpose of the **Marple** module is to constrain the prior space of the Bayesian retrieval in an unbiased way. Given the unknown, varied and complex nature of exoplanet systems it is difficult to pre-suppose atmospheric compositions from ‘experience’. The most objective approach to atmospheric retrieval of exoplanets would be to assume no prior knowledge at all and to consider all combinations of all atmospheric absorbers known. Whilst desirable, this is computationally infeasible due to the large number of free parameters and often limited spectral resolution of the observed data. The **Marple** module attempts to limit the number of possible absorbers by identifying likely molecules in the observed data using a pattern recognition algorithm. It attempts to identify absorption/emission features that are ‘typical’ for a molecular/atomic species and computes the possibility of such indicative features pertaining to a

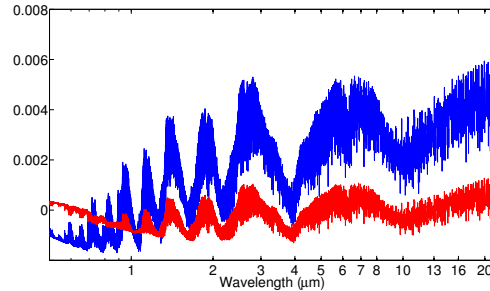


Fig. 3.— First (blue) and second (red) principal components of the water transmission spectrum library shown in figure 2. The first component is used to create the ‘feature mask’, see text, whereas the second component is correlated against the observed data to identify possible matches between the observed spectra and water features.

specific molecule compared to all other options. In this sense, the algorithm is conceptually similar to well established facial recognition algorithms using ‘eigenfaces’ (e.g. Turk & Pentland 1991; Cendrillon & Lovell 2000; Gevaert & de With 2013).

The algorithm is described in the following steps:

1. The **Marple** module generates a library of atmospheric spectra (equation 7) for each atmospheric species using the available absorption cross-sections, $\varsigma_{m,\lambda}(T)$. For each species, m , spectra are produced for a large range of atmospheric temperatures, T , and mixing ratios χ (typically $1 \times 10^{-8} \leq \chi \leq 1 \times 10^{-1}$ and $500K \leq T \leq 2000K$). Figure 2 shows the transmission spectra of water and CO over a range of temperatures and compositions. Where the molecular absorptions are strongest so are the variations. Note that a temperature range can be set for computational efficiency.
2. Characteristic spectral features for an atmospheric species vary significantly over the temperature and mixing ratio ranges computed above. These features are key to the identification of the absorber/emitter.

We can capture these significant variations using a principal component analysis (PCA, Jolliffe 2007) where the first component typically indicates the wavelength range over

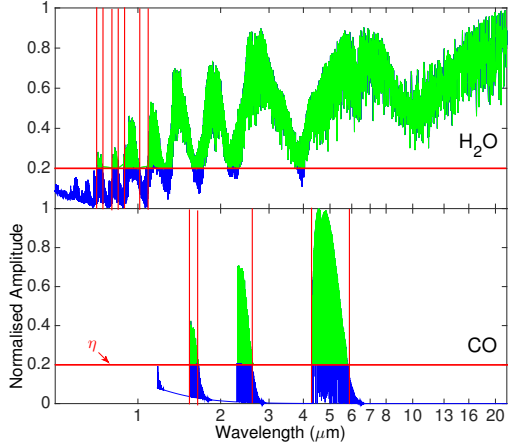


Fig. 4.— showing the creation of the ‘feature mask’, $\psi_m(\lambda)$, for H_2O (top) and CO (bottom). In blue are the first principal components of the molecules whilst in green are the spectral features selected for the mask. The horizontal, red line indicates the threshold parameter η , equation 11. Spectral features above this threshold are taken to be ‘significant features of the absorber’ and included in the feature mask. The vertical, red lines show the major cuts in the feature mask. For the case of H_2O , being an absorber/emitter across a broad wavelength range, most wavelengths are included in the feature mask $\psi_{\text{H}_2\text{O}}(\lambda)$. CO only absorbs/emits in discrete wavelength ranges and only those will be included in the feature mask of CO , $\psi_{\text{CO}}(\lambda)$.

which these features are most prominent (i.e. the amplitude of the variation) and the second component reflects the modulation on the bulk variation, i.e. the features’ morphology. Hence, for each atmospheric species, we compute the first and second principal components (PCs) over the range of spectra produced above using a single-value-decomposition (SVD).

The SVD of the column vector of spectra α is given by

$$\alpha_m(T, \chi) = \mathbf{U}\Sigma\mathbf{V}^T|_m \quad (9)$$

where \mathbf{U} and \mathbf{V} are the left and right unitary matrices respectively and Σ is the diagonal eigenvalue matrix. Due to the large size of

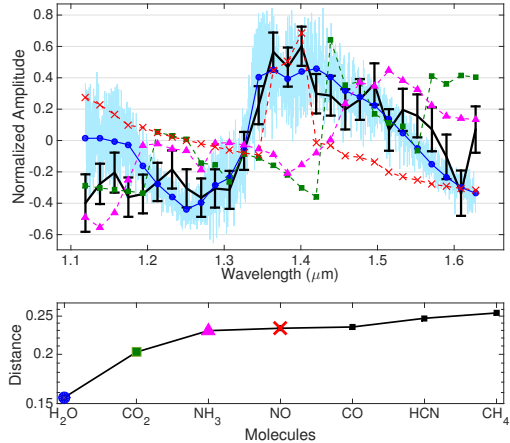


Fig. 5.— TOP: Hubble/WFC3 transmission spectrum of HD 209458b (Deming et al. 2013) (black). Overplotted are the principal components of the four best matching molecules binned to the resolution of the data: H_2O (blue dots), CO_2 (green squares), NH_3 (magenta triangles), NO (red squares). The H_2O principal component is also shown at a resolution of $R = 1000$ (light blue). BOTTOM: The normalised euclidean distance between each individual component and the data. It is clear that H_2O presents the best match to the data with other molecules being significantly worse.

the matrices involved, we approximate equation 9 with a randomised, truncated SVD algorithm (Halko et al. 2011; Martinsson et al. 2011). The individual principal component is then given by

$$\begin{aligned} \mathbf{pc}_{n,m} &= \mathbf{U}_n \Sigma_n |_m \\ &= \alpha_m(T, \chi) \mathbf{V}_{n,m} \end{aligned} \quad (10)$$

where n is the PC index. Figure 3 shows the first (blue) and second (red) principal components of water. In this case spectral features are preserved in both components.

3. We now use the first principal component calculated above to create the ‘feature mask’, $\psi_m(\lambda)$, for a given species. This masking guarantees that only wavelengths regions where a given molecule ab-

sorbs/emits are correlated against the observed spectrum. The feature mask is given by

$$\psi_m(\lambda) = \begin{cases} 1 & \text{if } \frac{\mathbf{pc}_{1,m} - \arg \min(\mathbf{pc}_{1,m})}{\arg \max(\mathbf{pc}_{1,m})} > \eta \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

where ‘arg min’ and ‘arg max’ stand for the minimal and maximal values of the argument or array. This boolean mask identifies at which wavelengths the characteristic spectral features are stronger than the threshold parameter η . In other words, the η parameter sets the threshold between ‘molecule present’ and ‘molecule absent’ over a given wavelength range. We find $\eta = 0.2$ to be a good choice for data with broad wavelength coverage. For spectra consisting of very few (< 20) data points over a narrow wavelength range, the user may not want to exclude (i.e. mask) any data points from the analysis. In such cases setting η to a low value, e.g. $\eta = 0.05$, will effectively prevent any wavelength range masking. A range of **Marple** module sensitivities can be explored by leaving η as free parameter, within user specified limits, over which the **Occam** module (section 7) can iterate.

Figure 4 shows the creation of the ‘feature mask’ for water and CO. Wavelengths with the normalised $\mathbf{pc}_{1,m}$ bigger than η will be included into the feature mask of molecule m .

4. For each species we now convolve the second PC and the observed data \mathbf{x} with the ‘feature mask’ to obtain the masked PC and observed data vectors $\widehat{\mathbf{pc}}_{2,m}$ and $\hat{\mathbf{x}}_m$ respectively

$$\widehat{\mathbf{pc}}_{2,m} = \psi_m \otimes \mathbf{pc}_{2,m} \quad (12)$$

$$\hat{\mathbf{x}}_m = \psi_m \otimes \bar{\mathbf{x}} \quad (13)$$

where \otimes denotes the convolution operator and $\bar{\mathbf{x}}$ is the normalised observed data vector given by

$$\bar{\mathbf{x}} = \left(\frac{\mathbf{x} - \arg \min(\mathbf{x})}{\arg \max(\mathbf{x})} \right) \quad (14)$$

5. In order to select the set of best matching molecular/atomic species to the observed spectrum, we implement a variant of the K-nearest-neighbour (K-NN, Cover & Hart 1967; Altman 1992) algorithm based on the euclidian distance between the spectral library principal component vectors and the observed data. For this we calculate the L^2 norm (also known as Euclidean norm or Euclidean distance), \mathfrak{d}_m , between masked data and the masked second PC for each molecule, m

$$\mathfrak{d}_m = \frac{1}{N} \|\hat{\mathbf{x}}_m - \widehat{\mathbf{pc}}_{2,m}\|_2 \quad (15)$$

where N is the total number of data points in \mathbf{x} . We now sort the euclidian distances in ascending order to form the monotonically increasing sequence of \mathfrak{d}_m

$$f(\mathfrak{d}_m) = \{\mathfrak{d}_{1,m}, \mathfrak{d}_{2,m}, \dots, \mathfrak{d}_{\phi,m}\} \quad (16)$$

where ϕ is the sequence index and $|\mathfrak{d}_{\phi-1}| < |\mathfrak{d}_{\phi}|$. The total distance is given by $\mathfrak{d}_{total} = \sum_{\phi}^M \{\mathfrak{d}_{\phi}\}$, where M is the total number of molecules considered. The algorithm distinguishes between the cluster of best matching (low \mathfrak{d}) and worst matching species (high \mathfrak{d}) by finding the series index associated with the highest second derivative of $f(\mathfrak{d}_m)$

$$\varphi = \phi \text{ where } \left[\arg \max \left(\frac{d^2 f}{d\mathfrak{d}^2} \right) \right] \quad (17)$$

the number of molecules selected by the pre-processor, N_m , is then given by

$$N_m = \begin{cases} \varphi & \text{if } \varphi > N_{m,min} \\ N_{m,min} & \text{otherwise} \end{cases} \quad (18)$$

where $N_{m,min}$ is a minimal number of molecules to be selected which can be set by the user. The selected molecules are given by $m_{select} = m_{\phi < \varphi}$.

Once determined the **Marple** module passes its list of selected atmospheric species to the central

data module (see figure 1) which will prepare all inputs for further analysis. The efficiency of the **Marple** module is a function of spectral resolution and signal to noise (S/N) of the data. This is self evident as any identification of features is impaired by a too coarse wavelength grid ($R < 10$) or high noise levels ($S/N < 5$). For those extreme cases, the user may specify a list of ‘must include’ molecules in the parameter files to be considered by the \mathcal{T} -REx.

5.1. Example: HD 209458b

We demonstrate the **Marple** module using a transmission spectrum of the hot-Jupiter HD 209458b obtained by the *Hubble*/WFC3 camera (Deming et al. 2013). The top of figure 5 shows the transmission spectrum in black and the principal components of the four best matching molecules. Note that all amplitudes are normalised and we only compare morphologies. The bottom panel summarises the normalised Euclidean distances (equation 15) for individual molecules. Here the black continuous line represents the sequence $f(\mathfrak{d}_m)$ in equation 16. The **Marple** module returns water as the most likely molecule present with CO₂ a more distant second. The presence of water as main absorbing species is in good agreement with the results of previous analyses (Deming et al. 2013; Madhusudhan et al. 2014).

6. Retrieval

\mathcal{T} -REx features three independent retrieval methods: 1) least-square minimisation using a quasi-Newtonian Limited-Memory Broyden-Fletcher-Goldfarb-Shannon (LM-BFGS) algorithm, section 6.2, 2) an Adaptive, multi-chain Markov Chain Monte Carlo algorithm, section 6.3.1 and 3) a nested-sampling algorithm using MultiNest, section 6.3.2.

Programatically, individual minimisation routines submit standardised requests to the central data module which in turn handles all calls to the forward module, the T-P profile and required inputs (figure 1). This modular approach guarantees consistency between model, data and retrieval codes as well as a high degree of flexibility in the analysis of the observed data.

6.1. Prior bounds

\mathcal{T} -REx by default uses uniform priors for all free parameters. As default, the isothermal temperature bounds are $T_{equ} \pm 200$ K, where T_{equ} is the planetary equilibrium temperature (this can either be derived by \mathcal{T} -REx given planetary/orbital parameters or set by the user). The molecular mixing ratios are bounded between 0.0 - 1.0×10^{-1} by default. All prior bounds can be manually specified by the user. The planet-star ratio, $(R_p/R_*)^2$, is treated as free parameter by default, with its upper/lower bounds derived from the reported observational uncertainty on this ratio. Griffith (2014) and Benneke & Seager (2013), amongst others, have noted strong degeneracies between the value of $(R_p/R_*)^2$ and various retrieval parameters (e.g. H₂O abundance and cloud opacities). We will further explore these degeneracies in a subsequent publication (Rocchetto et al. in prep.).

6.2. LM-BFGS minimization

The least-square minimisation allows us to obtain a quick look at the optimal model fit for the data without using the computationally expensive MCMC or Nested sampling routines. In this respect it is key to the pre-burning of the MCMC chain, as at least one chain can be started at the optimal solution and hence does not require a burn-in time (Brooks et al. 2011), as well as providing a valuable consistency check between model fits produced by the MCMC and MultiNest routines.

Large numbers of free parameters are often a limiting factor for simplex-downhill algorithms (e.g. Nelder & Mead 1965) commonly used. We find such amoeba algorithms to be insufficient and to often get stuck in local minima. \mathcal{T} -REx uses the LM-BFGS (Zhu et al. 1997; Morales & Nocedal 2011) algorithm instead, which being quasi-Newtonian uses the inverse Hessian matrix of the χ^2 surface to efficiently and robustly converge to the global maximum. We furthermore find the LM-BFGS to be more robust in the presence of observational noise than comparable methods.

6.3. Bayesian analysis

The Bayesian argument is given by

$$P(\theta|\mathbf{x}, \mathcal{M}) = \frac{P(\mathbf{x}|\theta, \mathcal{M})P(\theta, \mathcal{M})}{P(\mathbf{x}|\mathcal{M})} \quad (19)$$

where $P(\theta, \mathcal{M})$ is the Bayesian prior which we take to be uniform throughout this paper. The number and type of absorbing species, as well as the equilibrium temperature of the planet defining the forward model, \mathcal{M} , are set by the `Marple` module (section 5). $P(\theta|\mathbf{x}, \mathcal{M})$ is the posterior probability of the model parameters θ given the data, \mathbf{x} assuming the forward model \mathcal{M} . The likelihood, $P(\mathbf{x}|\theta, \mathcal{M})$ is given by the Gaussian

$$P(\mathbf{x}|\theta, \mathcal{M}) = \frac{1}{\epsilon\sqrt{2\pi}} \exp \left[-\frac{1}{2} \sum_{\lambda}^N \left(\frac{x_{\lambda} - \mathcal{M}_{\lambda}}{\epsilon_{\lambda}} \right)^2 \right] \quad (20)$$

where ϵ is the error on the observed spectral point. As opposed to the nested sampling described in the next section, an MCMC does not sample the Bayesian partition function (also known as Bayesian Evidence) and equation 19 reduces to

$$P(\theta|\mathbf{x}, \mathcal{M}) \propto P(\mathbf{x}|\theta, \mathcal{M})P(\theta, \mathcal{M}). \quad (21)$$

6.3.1. MCMC

MCMC routines are commonly used in the field of extrasolar planets (e.g. Ford 2006; Burke et al. 2010; Bakos et al. 2007; Knutson et al. 2007; Cameron et al. 2007; Charbonneau et al. 2009; Bean et al. 2010; Kipping & Bakos 2011; Gregory 2011; Crouzet et al. 2012; Kreidberg et al. 2014; Braak 2006; Line et al. 2013b; Madhusudhan et al. 2014; Benneke & Seager 2012; Ter Braak & Vrugt 2008; Foreman-Mackey et al. 2013; Goodman & Weare 2010; Madhusudhan et al. 2014). `T-REx` provides an implementation of the Delayed-Rejection Adaptive-MCMC (DRAM, Haario et al. 2006). We refer the interested reader to the cited literature and here only provide a brief overview. The DRAM algorithm differs from a more classical Metropolis-Hastings sampler (Metropolis & Rosenbluth 1953; Hastings 1970; Brooks et al. 2011) in two aspects: 1) It implements a delayed rejection algorithm and 2) an adaptive proposal distribution calibrated using the covariance of the

sample path of the MCMC chain. For additional information on DRAM, we refer the reader to Appendix A and the relevant literature.

`T-REx` runs several MCMC chains in parallel to check convergence and increase the sampling of the likelihood space. The number of chains is user defined but usually set to 4-5 and limited by the number of available CPUs. The first primary chain is started at the optimal values determined by the LM-BFGS, avoiding significant burn-in time (Brooks et al. 2011). All secondary chains' starting positions are offset from the optimum by a random distance and direction of at least 10% of the prior width. These secondary chains are run with a burn-in period of typically 10% of the total chain length. Burn-in and chain lengths are user defined.

6.3.2. Nested Sampling

Nested sampling (NS) algorithms (Skilling 2004, 2006; Mukherjee et al. 2006; Chopin & Robert 2010; Keeton 2011; Jasia & Xiang 2005) are becoming increasingly popular in extrasolar planets (e.g. Kipping et al. 2012; Placek et al. 2013) as well as Benneke & Seager (2013) for atmospheric retrieval. Here we include an implementation of `MultiNest` (Feroz & Hobson 2008; Feroz et al. 2009, 2013). MCMC algorithms are commonly used for parameter estimation by solving equation 21. Whereas MCMC explores the likelihood space by means of a Markovian chain, NS performs a general Monte Carlo (MC) analysis which is periodically constrained by ellipsoids encompassing spaces of highest likelihoods. Note that unlike MCMC, NS does not depend on a pre-determined proposal density and can hence better explore highly degenerate and non-Gaussian regimes. Using NS, we can compute the Bayesian evidence (or simply evidence), which is given by the integral required to normalising equation 21

$$E = \int P(\theta|\mathcal{M})P(\mathbf{x}|\theta, \mathcal{M})d\theta \quad (22)$$

where $E = P(\mathbf{x}|\mathcal{M})$ is the evidence. The evidence allows us to test the adequacy of the model itself and to perform model selection as described in the following section. Posterior distributions for parameter estimations are returned as by-product of `MultiNest` and should be similar to posteri-

ors obtained by the MCMC. Note that through the very different sampling techniques and fewer constraints on the proposal density for NS, we expect MCMC posteriors to be a ‘smoothed’ version of the NS’s. \mathcal{T} -REx allows the choice between importance nested sampling (INS) and the more classical NS. Through the sampling process, INS retrains all accepted as well as rejected proposal points which allows for a more accurate integration of the evidence (Feroz et al. 2013). Nested Sampling (in its `MultiNest` implementation) is highly efficient and easily parallelisable, allowing an easy scaling to cluster computing. We here use the NS approach as our main means of retrieval with the MCMC implementation providing a valuable cross check on the final results.

7. Model selection

For an inverse retrieval problem, such as the one discussed here, the idea of model selection is highly relevant but rarely discussed due to the computational expense and complexities involved. Notable examples of Bayesian model selection in atmospheric retrieval are Benneke & Seager (2013); Line et al. (2013a); Swain et al. (2014). Here we explicitly make the distinction between optimal estimation of parameters and the adequacy of the parameter and/or model itself.

We perform two tests after each \mathcal{T} -REx run:

1. Parameter adequacy: is a parameter (e.g. a given molecular species) required to describe the underlying physics? If not, is the model considered *over-complete*? In the case of *over-completeness* the forward model may be too complex (not obeying Occam’s razor). This can lead to overfitting in the worst case or in the best case a reduction in retrieval efficiency.
2. Model adequacy: are parameters missing in the model considered, i.e. is the model *under-complete*? In the case of model *under-completeness* the data is better accounted for by a more complex model. For example, a cloudy exoplanetary atmosphere cannot be modelled adequately by a cloud-free atmospheric model. Here the presence of clouds could force a cloud-free model to compensate for the extra absorption using molecu-

lar/atomic absorbers. This introduced bias often cannot be discerned from parameter estimating algorithms such as MCMC, maximum likelihood and similar methods.

Determining a model that is adequate to the data’s complexity is hence paramount.

\mathcal{T} -REx tries to perform model selection in an intelligent way through the `Occam` module. The `Occam` module will perform model selection until a complete model is determined. It will iterate through models, appropriately increasing or decreasing the model complexity through interaction with the `MarpLe` module (see figure 1, point 4).

7.1. Over-complete models

The over-complete model features an unnecessary complexity. Here the desired model is a subset of the more complex model initially run. We referred to these models as being ‘nested’. Complexity in parametric models (such as the forward models of atmospheric retrieval) is usually synonymous with number of free-parameters. Hence we can define the nested model $\mathcal{M}_{\theta-\theta_\gamma}$ as sub-set of the more complex one \mathcal{M}_θ ,

$$\mathcal{M}_{\theta-\theta_\gamma} = \mathcal{M}_\theta|_{\theta_\gamma=0} \quad (23)$$

where θ is a column-vector of all model parameters and θ_γ is an individual parameter. The Bayes factor (see section 7.2) allows us to perform this model selection by marginalising out individual parameters (Benneke & Seager 2013; Swain et al. 2014). For ‘nested’ models we can derive the simpler to compute Savage-Dickey density ratio (SDR) (Dickey 1971; Verdinelli & Wasserman 2012; Marin & Robert 2010; Verde et al. 2013)

$$SDR = \frac{P(\theta_\gamma|\mathbf{x}, \mathcal{M}_\theta)}{P(\theta_\gamma|\mathcal{M}_\theta)} \Big|_{\theta_\gamma=0} \quad (24)$$

where $P(\theta_\gamma|\mathbf{x}, \mathcal{M}_\theta)$ is the marginalised posterior of θ_γ and $P(\theta_\gamma|\mathcal{M}_\theta)$ its respective prior distribution. A comprehensive derivation and discussion of equation 24 can be found in Verde et al. (2013). This ratio of densities at $\theta_\gamma = 0$ is indicative of whether a simpler model not containing θ_γ is sufficient or whether a more complex model is preferred by the data. To assess the significance of the evidence towards a complex rather than a simpler

model, we compare the outcome of equation 24 to the Jeffrey’s scale (Jeffreys 1961). We adopt a slightly modified version of Kass & Raftery (1995) in table 1

Table 1: Jeffrey’s scale for model selection

$2\ln(SDR)$	Preference for simplified model $\mathcal{M}_{\theta-\theta_\gamma}$
> 10	Very strong preference for excluding θ_γ
10 to 6	Strong preference for excluding θ_γ
6 to 2	Substantial preference for excluding θ_γ
2 to 0	Insignificant preference for excluding θ_γ
Preference for complex model \mathcal{M}_θ	
0 to -2	Insignificant preference for including θ_γ
-2 to -6	Substantial preference for including θ_γ
-6 to -10	Strong preference for including θ_γ
< -10	Very strong preference for including θ_γ

The `Occam` module calculates the Savage-Dickey ratio for each model parameter and adjust the model complexity accordingly in case of a strong preference for a simpler forward model.

7.2. Under-complete models

Should the model at hand not be over-complete, the `Occam` module tests for model under-completeness, i.e. is the model complex enough. For this we iteratively re-run the retrieval process allowing the `MarpLe` module to add the two next most likely molecular opacities to the current selection of opacities. We then compute the global model evidence, E , and compute the Bayes factor (Kass & Raftery 1995; Weinberg 2012). The Bayes factor is given by the ratio of model probabilities $P(\mathcal{M}_1|\theta)$

$$\frac{P(\mathcal{M}_2|\theta)}{P(\mathcal{M}_1|\theta)} = \frac{P(\mathcal{M}_2) P(\mathbf{x}|\mathcal{M}_2)}{P(\mathcal{M}_1) P(\mathbf{x}|\mathcal{M}_1)} = \frac{P(\mathcal{M}_2) E_2}{P(\mathcal{M}_1) E_1} \quad (25)$$

which can be expressed as fraction of the Evidences and the prior distribution of the models. Most times we can assume the model priors to be identical $P(\mathcal{M}_1) = P(\mathcal{M}_2)$, reducing equation 25 to

$$E_2/E_1 = \frac{P(\mathbf{x}|\mathcal{M}_2)}{P(\mathbf{x}|\mathcal{M}_1)}. \quad (26)$$

Using the Jeffrey’s scale the `Occam` module determines whether an improvement to the fit is

achieved using a more complex model.

8. Outputs

The output module generates the best fit transmission model, plots of all marginalised and conditional posteriors as well as statistics on individual parameters and model adequacy. Examples of these outputs can be found in the following section.

9. Example

In this section we demonstrate the output of \mathcal{T} -REx using a simulated hot-Jupiter. We base the simulation on a HD209458b like planet/star system (Charbonneau et al. 2000; Southworth 2010) with temperature and bulk composition taken from Venot et al. (2014). We choose a wavelength range of $1 - 20\mu\text{m}$ at a constant resolution of $R = 300$ and constant error bars of 50ppm. Table 2 summarises the inputs and figure 6 shows the input spectrum to \mathcal{T} -REx. Whilst such an example may be optimistic given currently available data we would like to note the following: 1) In order to demonstrate the retrieval accuracy of \mathcal{T} -REx one needs a precise data set, 2) Future observatories and missions (e.g. JWST, E-ELT and dedicated missions) will yield data of comparable or better quality over broad wavelength ranges.

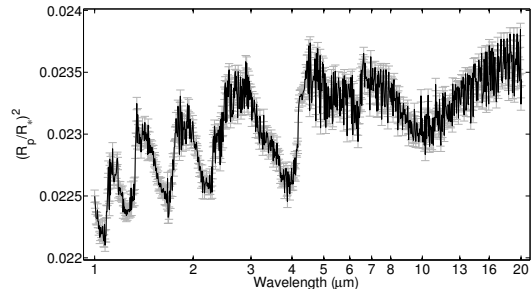


Fig. 6.— Simulated example spectrum of a carbon-rich hot-Jupiter used in section 9. Temperature and abundances of the main absorbers, H_2O , CO , CO_2 , NH_3 and CH_4 are given in table 2. The bulk planet/star and orbital properties are based on the hot-Jupiter HD209458b.

The data is passed through \mathcal{T} -REx as described in the previous sections. The `MarpLe` module suggested the correct molecules as potentially im-

Table 2: Model input and retrieval results for LM-BFGS, MCMC and NS. All values (but temperature) are in units of fractional column density.

Parameters	Model	LM-BFGS	MCMC	Nested Sampling
Temp. (K)	1400	1419.14	1403.30 \pm 9.88	1403.87 \pm 9.26
H ₂ O	2 \times 10 ⁻³	1.94 \times 10 ⁻³	1.90 \times 10 ⁻³ \pm 3.27 \times 10 ⁻⁵	1.90 \times 10 ⁻³ \pm 3.11 \times 10 ⁻⁵
CH ₄	2 \times 10 ⁻⁶	2.04 \times 10 ⁻⁶	2.25 \times 10 ⁻⁶ \pm 1.45 \times 10 ⁻⁶	2.17 \times 10 ⁻⁶ \pm 1.42 \times 10 ⁻⁶
CO	2 \times 10 ⁻³	1.95 \times 10 ⁻³	1.97 \times 10 ⁻³ \pm 1.26 \times 10 ⁻⁴	1.97 \times 10 ⁻³ \pm 1.22 \times 10 ⁻⁴
CO ₂	2 \times 10 ⁻⁵	2.41 \times 10 ⁻⁵	2.48 \times 10 ⁻⁵ \pm 2.59 \times 10 ⁻⁶	2.48 \times 10 ⁻⁵ \pm 2.60 \times 10 ⁻⁶
NH ₃	2 \times 10 ⁻⁷	1.48 \times 10 ⁻⁶	1.18 \times 10 ⁻⁶ \pm 9.69 \times 10 ⁻⁷	1.18 \times 10 ⁻⁶ \pm 9.69 \times 10 ⁻⁷

portant absorbers given the data and their wavelength ranges. In addition to the molecules listed in table 2, it also identified H₂C₂ as possible absorber which was subsequently rejected by the *Occam* module and the transmission module was updated to reflect the true model of the data.

The retrieved temperature and abundance values for the LM-BFGS, MCMC and Nested Sampling algorithms are summarised in table 2. *TREx* does not compute a formal error for the LM-BFGS result as only the MCMC and Nested Sampling results are considered to be final data products. Figure 7 (top spectrum) shows the best-fit transmission model for the MCMC (green) and the Nested Sampling (red) algorithms. Figures 8 & 9 show the marginalised and conditional posterior distributions for the MCMC and Nested Sampling results respectively. The MCMC results consist of 8 independent chains (note the different colours in the marginalised posteriors representing the results of individual chains) and 2.5 \times 10⁴ samples each, including a 10% burn-in period. The Nested sampling results used 4000 initial live-points and 8.3 \times 10⁴ replacements. Note the Nested Sampling posteriors to be more kurtotic than the MCMC results. This is due to a finer sampling of the maximum likelihood space by the NS.

Table 2 summarises the retrieved abundances. All major species in the simulated transmission spectrum as well as the isothermal temperature were retrieved with great fidelity by all retrieval methods. Mixing ratios for NH₃ and CH₄ were set purposefully low to test the retrievability of molecular abundances at the limits of data uncertainties. As shown in section 9.1 and table 3, these detections were identified as ‘insignificant’ by the *Occam* module.

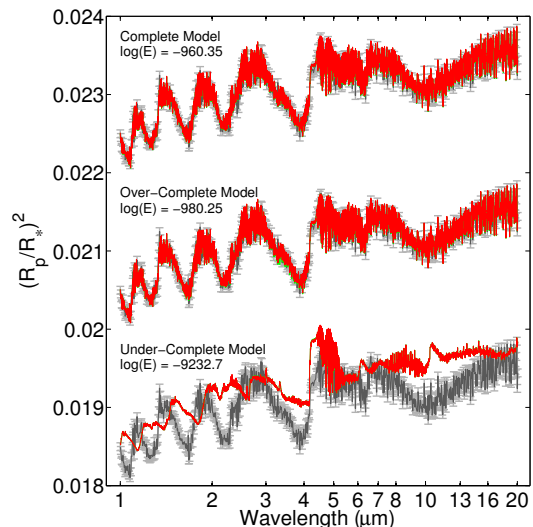


Fig. 7.— showing the best fitting models, in red, for the complete (top), over-complete (middle) and under-complete (bottom) model cases as described in section 9. The best fitting models are offset along the ordinate for clarity and over-plotted on the ‘observed’ spectrum in grey. The global Bayesian Evidences, $\log(E)$, are given for each case, quantifying the adequacy of each model given the data. As expected the evidence strongly favours the correct, complete model.

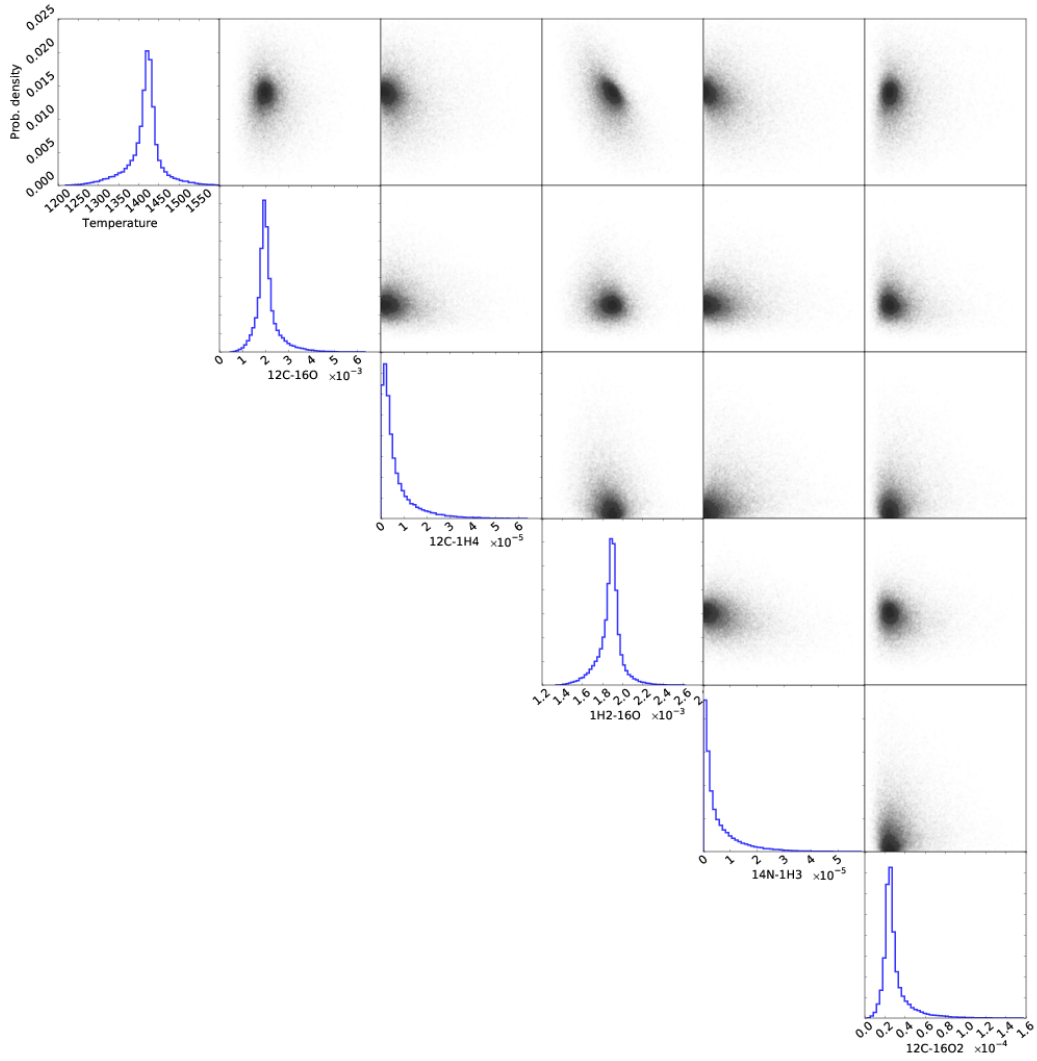


Fig. 8.— showing the marginalised and conditional posterior distributions for the Nested Sampling for the complete model in section 9. We find the highest correlation between atmospheric temperature and water absorption. With H₂O being the strongest absorber across the broadest wavelength range, this correlation between abundance and thermal broadening due to temperature changes is to be expected.

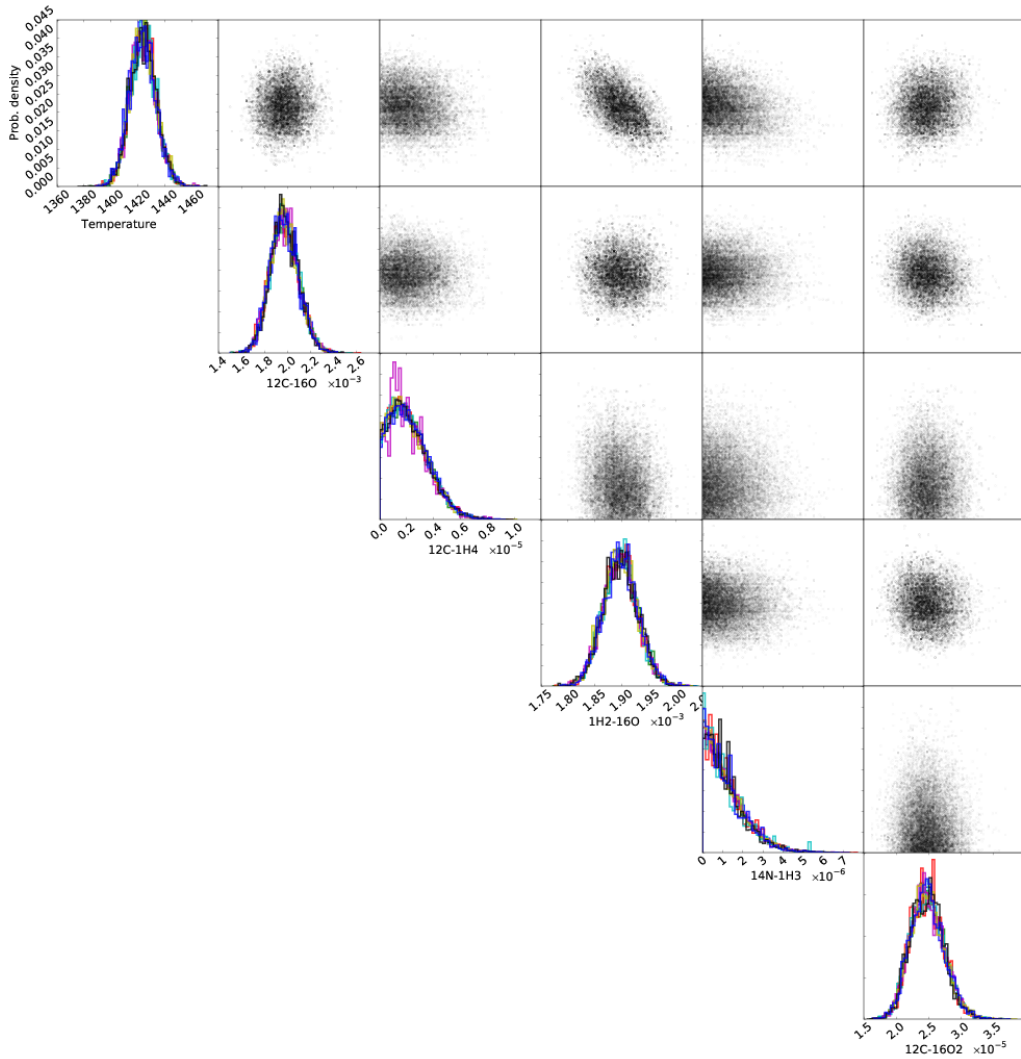


Fig. 9.— showing the marginalised and conditional posterior distributions for the MCMC run of the complete model. Different colours in the marginalised posterior plots represent individual MCMC chains. The very good overlap of these independent chains indicates a good convergence of the code.

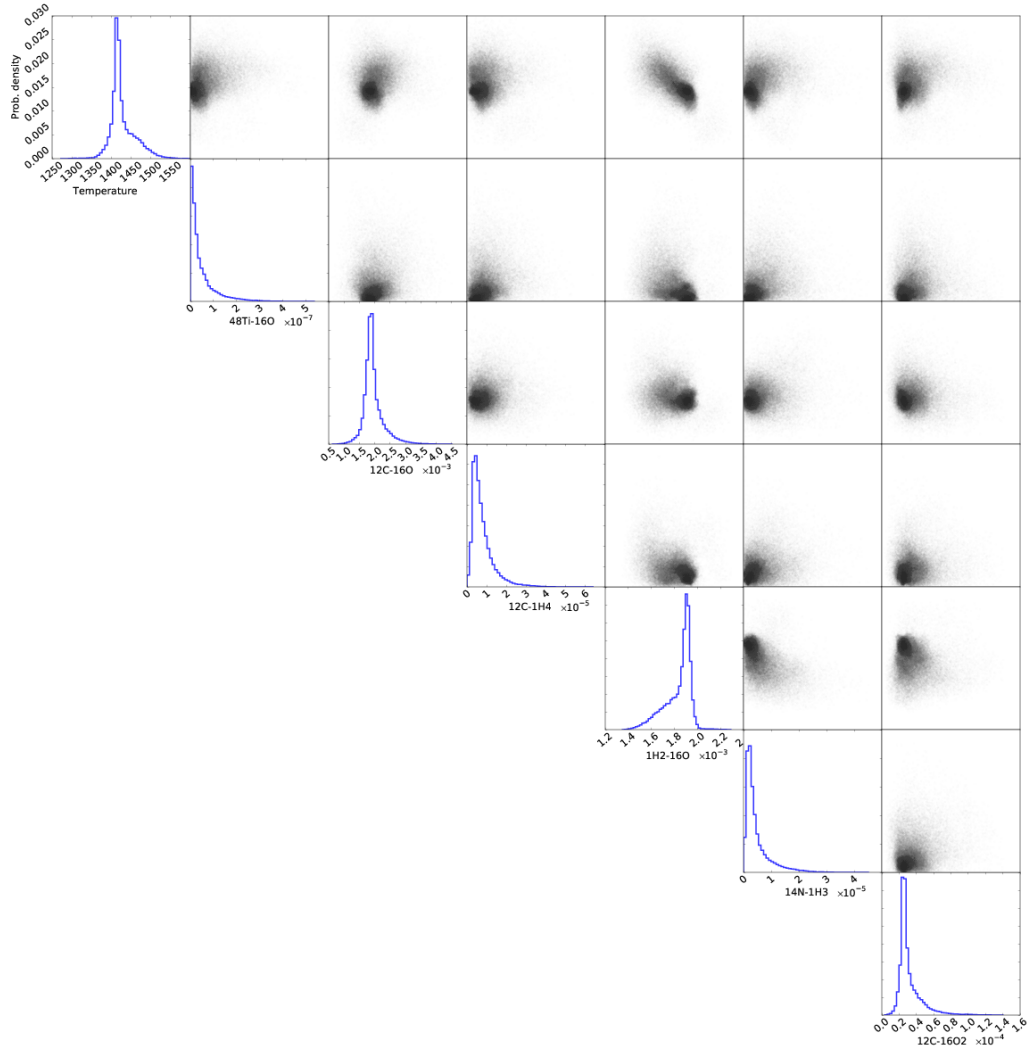


Fig. 10.— showing marginalised and conditionals posterior distributions of the Nested Sampling run for the over-complete model in section 9.1. Otherwise identical to figure 8.

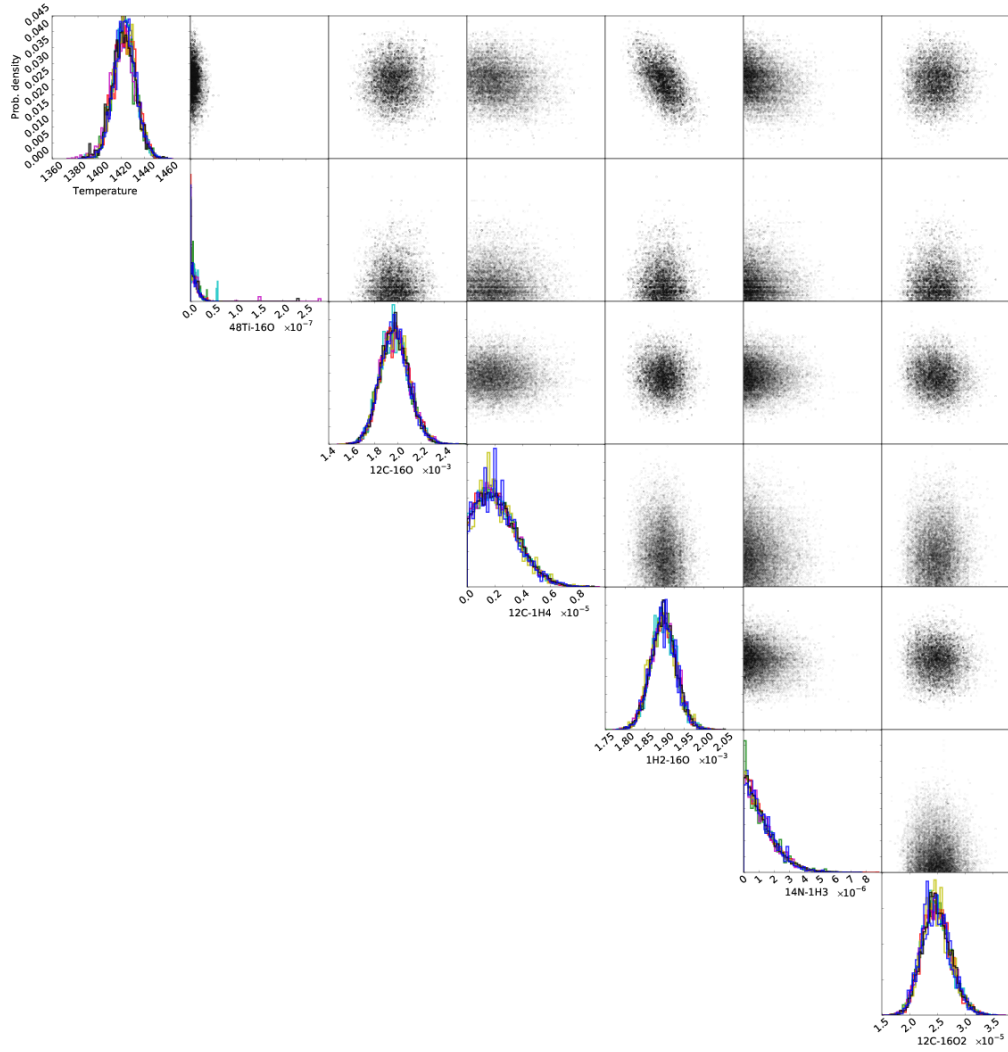


Fig. 11.— showing marginalised and conditional posterior distributions of the MCMC run for the over-complete model in section 9.1. Otherwise identical to figure 9.

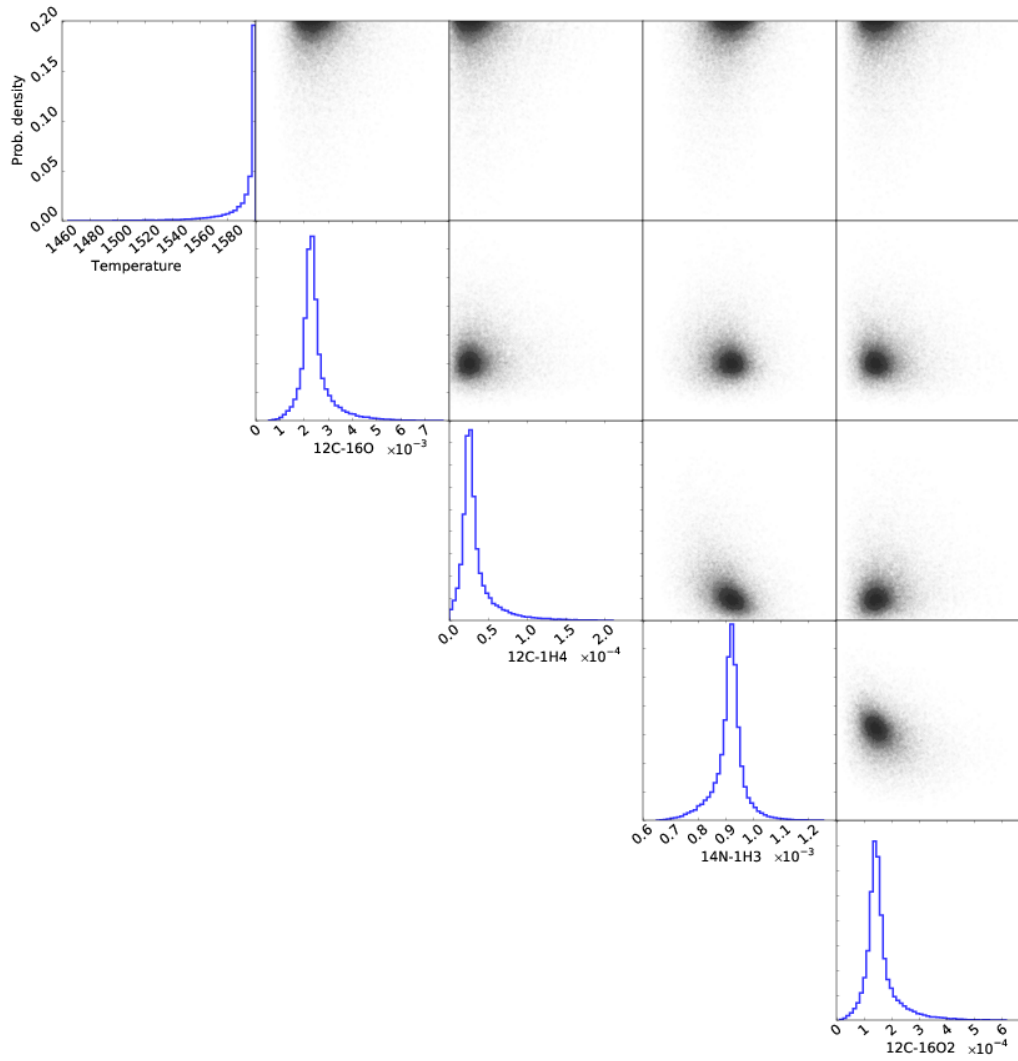


Fig. 12.— showing marginalised and conditionals posterior distributions of the Nested Sampling run for the under-complete model in section 9.1. In the absence of the main absorbing species, H_2O , \mathcal{T} -REx tries to compensate for lacking opacity by increasing thermal broadening. This results in the planetary temperature converging to the upper end of the prior. Otherwise identical to figure 8.

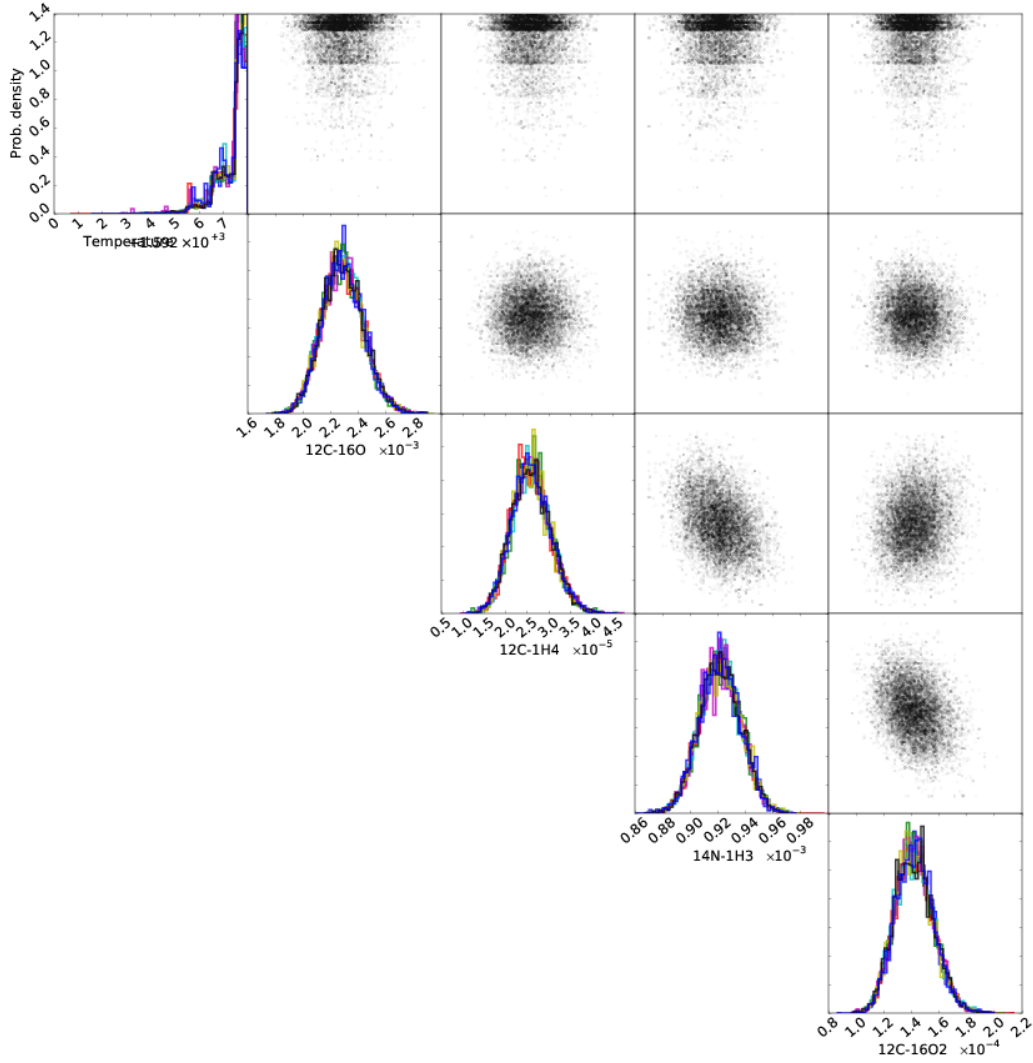


Fig. 13.— showing marginalised and conditionals posterior distributions of the MCMC run for the under-complete model in section 9.1. Similar to figure 12 the MCMC is converging to the upper temperature prior rapidly. This rapid convergence leaves ‘step’ artefacts in the temperature posterior due to the discrete temperature resolution of 1K. Step sizes can be set to an arbitrarily small value but in this case convergence to the upper prior bound is fast enough to only ever sample the top bins after the burn in period is completed. Otherwise identical to figure 9.

9.1. Model Selection

In addition to the complete model shown above, we have simulated an over-complete and under-complete model to test the model selection abilities of \mathcal{T} -REx. The over-complete model contains TiO as additional absorber and the under-complete model lacks H₂O. Whereas in terms of χ^2 statistics we would expect a similarly valid fit for the over-complete model compared to the complete case (as the excessive parameters should converge to small values), we would expect a decrease in the overall model evidence as well as a clear discrimination of unnecessary complexity in the SDR. This behaviour is indeed demonstrated by \mathcal{T} -REx. Figure 7 (middle) shows the model fit of the over-complete model and figures 10 and 11 the posterior distributions of the NS and MCMC fits respectively. As figure 7 shows, the fit is maintained but at a lower global evidence, $\log(\mathbf{E}) = -980$ compared to $\log(\mathbf{E}) = -960$ for the correct model. On the Jeffrey’s scale this results in a very strong preference for the overall simpler (i.e. complete) model. Table 3 shows the SDRs calculated from the NS and MCMC posteriors for $\mathcal{M}_{\theta-\theta_\gamma}/\mathcal{M}_\theta$, where θ_γ is the molecule in question. The SDRs show a substantial to strong preference for the exclusion of TiO from the model and strongly confirm the inclusion of CO, H₂O and CO₂. For the two low column density species, CH₄ and NH₃, the SDRs neither include nor exclude either species but do not support a significant detection of the molecule in the data, as expected. Differences in the SDR derived between NS and MCMC are due to the NS providing a tighter constraint on the marginalised posterior distributions than the MCMC, see figures 10 & 11.

Figure 7 (bottom) shows the under-complete model excluding water. For under-complete models the χ^2 increases significantly as well as a very low global evidence of $\log(\mathbf{E}) = -9232$. Figures 12 and 13 show the posterior distributions of the NS and MCMC runs respectively. Both show a

Table 3: Savage-Dickey density ratio (SDR) for overcomplete model

$2\ln(\text{SDR})$	TiO	CO	CH ₄	H ₂ O	NH ₃	CO ₂
NS	5.5	-31.9	-0.5	-31.9	1.9	-31.9
MCMC	10.3	-29.3	1.6	-29.3	3.6	-29.3

strong over dependence on high atmospheric temperatures, trying to fill in the missing opacities with and increased absorption due to an increased planetary scale height and an increased spectral broadening through the emergence of molecular hot-bands at higher temperatures.

9.2. Resolution and Signal-to-Noise

We now take the complete model from the previous section and reduce the signal to noise (S/N) and resolution to explore the impact on the retrieval of exoplanetary spectra. Given the potentially large scope of such an exercise we here limit ourselves to a S/N-Resolution grid most representative of current data: R = 300, 200, 100, 50, 30 and spectral error bars of $\sigma = 10\text{ppm}$, 50ppm, 100ppm and 500ppm. Figure 14 shows the input data in red and the best fitting transmission model at a resolution of R = 500. Whereas visually all spectra fit equally well, degeneracies between mixing ratios and temperature increase as resolution and S/N decrease. Whereas this result is intuitive, we find that the effect of degrading resolution and S/N is not uniform amongst parameters. Figure 15 shows the increase in error-bar (in percent) for the temperature posterior distribution derived. Here the reduction in S/N (i.e. increase in σ) has a much more pronounced effect than the reduction in resolution, meaning that the planetary temperature can be derived with high confidence for low resolution data but not vice versa. Figure 16 shows the same plot for the retrieval of the water mixing ratio. Water being very broad absorber in the NIR to mid-IR its abundance retrieval depends on resolution and S/N in approximately equal measures. Figure 17 plots the retrieval error bars of carbon-monoxide. With the spectral feature of CO being less broad than water we see that the dependence on S/N exceeds that of resolution (assuming that each CO feature is captured by at least one data point).

The loss of information through a reduced S/N is also demonstrated in figure 18 showing the posterior distribution for all σ considered at R = 200. Here we can see the transition of the Bayesian argument going from data to prior dominated as the S/N decreases. Such a transition is gradual and we find a more regular occurrence of small local likelihood maxima as the likelihood surface “flattens out” at increasing σ . This is demonstrated

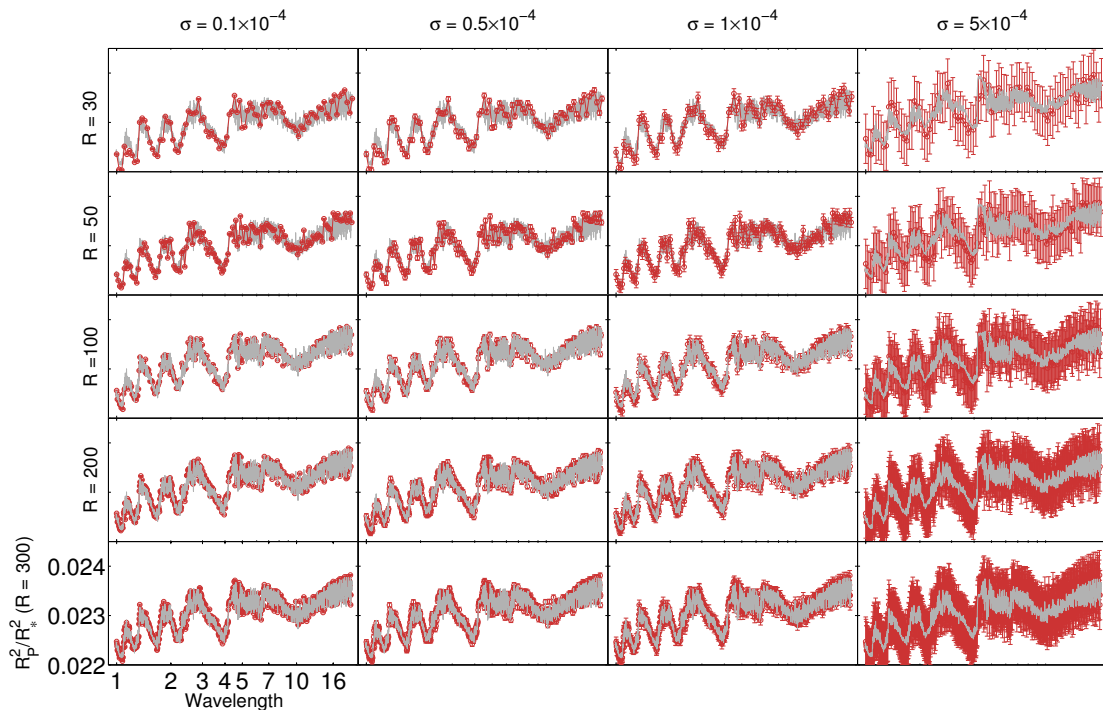


Fig. 14.— Best fitting transmission model (grey) at $R = 500$ superimposed on simulated input data (red) at resolutions $R = 30, 50, 100, 200, 300$ and data-error bars $\sigma = 10\text{ppm}, 50\text{ppm}, 100\text{ppm}, 500\text{ppm}$.

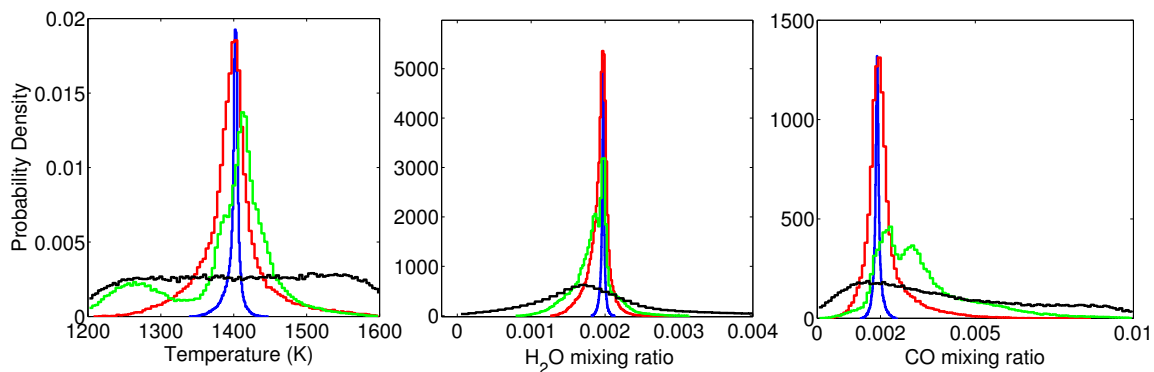


Fig. 18.— Three figures showing the posterior distributions for the planetary temperature (left), H_2O (middle) and CO (right) mixing ratios at resolution $R = 200$. Colours represent data error bars. Blue: $\sigma = 10\text{ppm}$ (scaled by a factor of $1/5$ to improve comparability to other curves); Red: 50ppm ; Green: 100ppm ; Black: 500ppm .

Table 4: Savage-Dickey Ratios (SDRs) for H₂O, CO and NH₃ for data error-bars of $\sigma = 10$ ppm, 50ppm, 100ppm, 500ppm and resolutions of R = 300, 200, 100, 50, 30. Negative values signify a detection of the molecule with values < -10 being a very strong detection. Similarly, positive values > 6 strongly indicate a non-detection. Values between -2 and +2 are inconclusive.

Resolution	10ppm	50ppm	100ppm	500ppm	
H ₂ O	300	-28.3	-28.3	-28.5	-1.6
	200	-28.3	-28.6	-28.5	-2.6
	100	-28.3	-28.4	-28.5	-1.9
	50	-28.2	-28.2	-28.6	-0.8
	30	-28.7	-28.7	-28.7	-1.2
CO	300	-28.3	-28.3	-28.5	-1.4
	200	-28.3	-28.6	-28.5	-1.0
	100	-28.3	-28.4	-2.7	-0.1
	50	-28.2	-28.2	-2.3	-0.1
	30	-28.7	-28.7	-2.1	0.2
NH ₃	300	8.0	5.7	4.8	9.1
	200	6.6	2.3	4.4	7.9
	100	0.5	3.0	3.6	6.6
	50	-1.6	4.8	6.2	5.6
	30	7.7	7.7	8.4	5.5

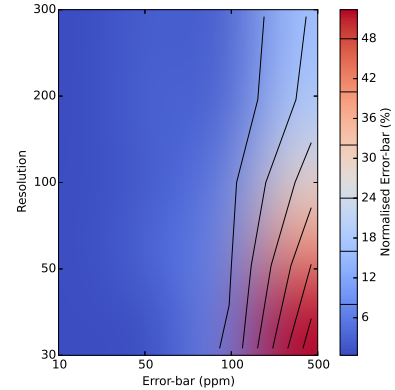


Fig. 16.— H₂O posterior standard deviation normalised by the ground-truth abundance ($\chi_{(H_2O)} = 2 \times 10^{-3}$) as function of spectral resolution (R) and the data-error bar. The ability to retrieve water abundances remain relatively stable for high-resolution and low S/N data but significantly decreases as both S/N and R drop. Here posterior error-bars can reach the size of prior space.

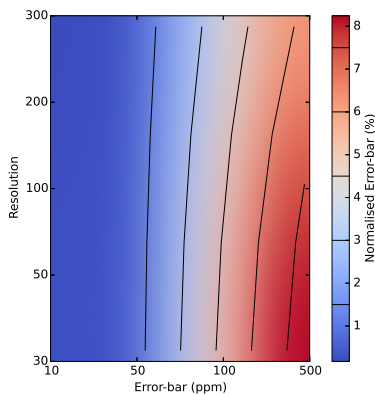


Fig. 15.— Temperature posterior standard deviation normalised by the ground-truth temperature (1400K) as function of spectral resolution (R) and the data-error bar. We find that the retrieval of the planetary temperature is dominated by the signal-to-noise of the data and less dominated by the resolution of the spectrum.

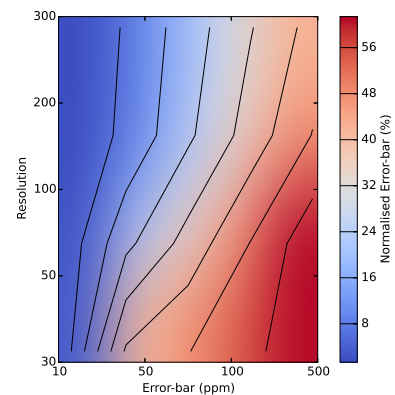


Fig. 17.— CO posterior standard deviation normalised by the ground-truth abundance ($\chi_{(CO)} = 2 \times 10^{-3}$) as function of spectral resolution (R) and the data-error bar. The CO retrieval more strongly depends on S/N other than for the resolution grid considered here.

by the the local maxima in the green curve ($\sigma = 100\text{ppm}$) compared to the best determined solution (blue curve, $\sigma = 10\text{ppm}$) and the prior driven solution (black curve, $\sigma = 500\text{ppm}$). The theoretical behaviour of a Bayesian retrieval at low R and low S/N is an important result which will be address in detail in future work.

Finally, we explore the effect of S/N and R on model selection. Similar to table 3 we calculated the SDR for H_2O , CO and NH_3 for the above S/N and R grid (table 4). As with previous examples negative values below < -6 indicate a strong detection with < -10 being a decisive detection whilst equally positive values rule out a detection in the data. As expected, we see the detection evidence for H_2O and CO decrease at high σ and low R. NH_3 (and CH_4 not shown) remain undetected throughout as expected.

Line et al. (2012) present a complementary analysis for the emission spectroscopy case where the overall informational content, and the resulting possible number of retrievable free parameters, is calculated. Their analysis points at the inverse relationship between S/N and R. The lower the S/N of an observation, the higher the resolution must be to obtain the same degree of retrievability, and vice versa. This relationship we also find for the retrievability of individual parameters in the transmission case, e.g. CO in figure 17.

10. Summary & Conclusion

In this publication we have introduced the \mathcal{T} -REx retrieval code for exoplanetary atmospheres. As described in the introduction and shown throughout the text, we have based the design of \mathcal{T} -REx on three guiding principles: 1) Sensitivity, 2) Objectivity, and 3) Big data.

\mathcal{T} -REx incorporates a line-by-line radiative transfer code using state-of-the-art molecular opacity line lists by the *ExoMol* project. Atmospheric transmission models are run at ~ 50 - 100 times higher resolution than the observed data to ensure a correct treatment of optically thick absorption lines as well as allowing for a precise treatment of thermal line broadening through an arbitrarily finely sampled temperature grid.

Given the large number of potential absorbing/emitting species of an extrasolar planet, we have developed custom build pattern recognition

software (the **Marple** module) to rapidly scan large molecular and atomic line-list archives for possible absorbing/emitting signatures in the observed spectrum. By not manually specifying a list of molecules ‘expected’ to be present in the atmospheres of exoplanets we break potential human biases in the selection of the atmospheric model. In other words by not assuming anything about the atmospheres composition and structure we maximise the objectivity of the analysis from the start.

Whereas the **Marple** module sets the potential prior space of the fully Bayesian retrieval, the **Occam** module performs iterative Bayesian model selection and iteratively verifies the adequacy of individual parameters as well as the overall evidence of the atmospheric model itself.

By using efficient MCMC and Nested Sampling techniques throughout, we are able to parallelise the sampling of the likelihood space making \mathcal{T} -REx natively scaleable to cluster computing. This allows \mathcal{T} -REx to explore very large parameter spaces and accurately map correlation manifolds.

We demonstrated individual properties of \mathcal{T} -REx using a simulated hot-Jupiter and explored the model selection process of over-complete and under-complete models. The quality of the retrieval was investigated for varying resolutions and signal-to-noise ratios of the input data and found to be consistent with expectations.

Future work will see a detailed treatment of emission spectroscopy in the framework of \mathcal{T} -REx, explore modelling degeneracies over large and short wavelength ranges and see the application to individual data sets.

With the maturation of data reduction techniques for exoplanetary spectroscopy we obtain higher and higher precision spectroscopy of these exotic atmospheres. With higher precision of the data often comes higher complexity in the interpretation. The goal of an ideal retrieval of atmospheric properties is to be able to capture said complexity whilst maintaining the highest possible degree of objectivity in the analysis. \mathcal{T} -REx presents a significant step towards this goal.

Acknowledgements

We thank the referee for providing useful comments. This work was supported by the ERC

project numbers 617119 (ExoLights) and 267219 (ExoMol).

A. DRAM

In standard Metropolis-Hastings samplers each proposal step can either be accepted or rejected based on a fitness criterion and often a probability of acceptance when the fitness criterion is not met. Should the proposal be rejected the MCMC chain remains in the same position on the likelihood space. The delayed rejection (DR) mechanism allows for a second (and third) proposal attempt to be made which is dependent on the previous chain as well as the previously rejected proposals. The adaptive proposal distribution based on its past history furthermore increases the efficiency and accuracy of the chain's exploration as the proposal distribution is iteratively adapted to the target distribution. These features can be shown to significantly improve the efficiency of the MCMC chain in high dimensional likelihood space.

B. Glossary

Variable	Description	Equation example
N	Number of spectral points in data	
N_m	Number of molecules selected for retrieval	
m	Molecular species index	
λ	Wavelength index	
\mathbf{x}	Data column vector	
$\bar{\mathbf{x}}$	Normalised data column vector	14
ς	Absorption cross section	1
T	Temperature (K)	
a, b	Absorption cross section temperature interpolation coefficients	1, 2, 3
χ	Atmospheric mixing ratio	
$\boldsymbol{\tau}$	Optical depth column vector over λ	
$\mathbf{I}(z)$	Intensity column vector over λ as function of z	4
z	Height in atmosphere	
$\boldsymbol{\alpha}$	Total atmospheric absorption column vector	7
$\boldsymbol{\alpha}_m(T, \chi)$	Total atmospheric absorption column vector as function of molecule, temperature, mixing ratio	9
R_p	Planetary radius	
R_*	Stellar radius	
$\boldsymbol{\delta}$	Transit-depth column vector	8
\mathbf{U}	Left unitary matrix of single value decomposition	9, 10
$\boldsymbol{\Sigma}$	Diagonal matrix of single value decomposition	9, 10
\mathbf{V}	Right unitary matrix of single value decomposition	9, 10
\mathbf{pc}	Principal component vector	10
n	Principal component index	
$\boldsymbol{\psi}_m$	Boolean data masking vector as function of molecule and wavelength	11
η	Molecule detection threshold coefficient	11
$\hat{\mathbf{x}}$	Masked data vector	13
$\hat{\mathbf{pc}}$	Masked PCA vector	12
\mathfrak{d}_m	l_2 -norm between normalised data and 2^{nd} principal component	15
$f(\mathfrak{d}_m)$	Monotonically increasing function of \mathfrak{d}_m	16
ϕ	Marple cluster index	16
φ	Marple cluster index for highest second derivative of $f(\mathfrak{d})$	17
m_{select}	Marple determined molecular/atmoic species	
\mathcal{M}	Exoplanet model	
θ_γ	Generic parameter of model \mathcal{M}	
$\boldsymbol{\theta}$	Column vector of parameters of model \mathcal{M}	
σ_λ	One sigma error at wavelength λ	
$P(\boldsymbol{\theta} \mathbf{x}, \mathcal{M})$	Posterior probability distribution of $\boldsymbol{\theta}$ given \mathbf{x} and \mathcal{M}	19
$P(\mathbf{x} \boldsymbol{\theta}, \mathcal{M})$	Likelihood distribution of $\boldsymbol{\theta}$	19
$P(\boldsymbol{\theta}, \mathcal{M})$	Prior distribution of $\boldsymbol{\theta}$	19
$P(\mathbf{x} \mathcal{M})$	Bayesian partition function	19
E	Bayesian Evidence	22

REFERENCES

- Altman, N. S. 1992, *The American Statistician*, 46, 175
- Bakos, G. Á., Noyes, R. W., Kovacs, G., et al. 2007, *ApJ*, 656, 552
- Barber, R. J., Strange, J. K., Hill, C., et al. 2014, *Monthly Notices of the Royal Astronomical Society*, 437, 1828
- Barber, R. J., Tennyson, J., Harris, G. J., & Tolchenov, R. N. 2006, *VizieR On-line Data Catalog*, 6119, 0
- Barstow, J. K., Aigrain, S., Irwin, P. G. J., Fletcher, L. N., & Lee, J. M. 2013, *Monthly Notices of the Royal Astronomical Society*, 434, 2616
- Barton, E. J., Chiu, C., Golpayegani, S., et al. 2014, *arXiv.org*, 7952
- Barton, E. J., Yurchenko, S. N., & Tennyson, J. 2013, *Monthly Notices of the Royal Astronomical Society*, 434, 1469
- Bean, J. L., Kempton, E. M.-R., & Homeier, D. 2010, *Nature*, 468, 669
- Benneke, B., & Seager, S. 2012, *APJ*, 753, 100
- . 2013, *APJ*, 778, 153
- Berta, Z. K., Charbonneau, D., Desert, J.-M., et al. 2012, *APJ*, 747, 35
- Braak, C. J. F. T. 2006, *Statistics and Computing*, 16, 239
- Brooks, S., Gelman, A., Jones, G., & Meng, X. L. 2011, *Handbook of Markov Chain Monte Carlo* (Chapman & Hall)
- Brown, T. M. 2001, *APJ*, 553, 1006
- Burke, C. J., McCullough, P. R., Bergeron, L. E., et al. 2010, *ApJ*, 719, 1796
- Cameron, A. C., Wilson, D. M., West, R. G., et al. 2007, *ApJ*, 380, 1230
- Carter, J. A., & Winn, J. N. 2009, *APJ*, 704, 51
- Cendrillon, R., & Lovell, B. 2000, *Proc. SPIE Vol. 4067*, 4067, 269
- Charbonneau, D., Brown, T. M., Latham, D. W., & Mayor, M. 2000, *ApJL*, 529, L45
- Charbonneau, D., Brown, T. M., Noyes, R. W., & Gilliland, R. L. 2002, *APJ*, 568, 377
- Charbonneau, D., Berta, Z. K., Irwin, J., et al. 2009, *Nature*, 462, 891
- Chopin, N., & Robert, C. P. 2010, *Biometrika*
- Conrath, B., Curran, R., Hanel, R., et al. 1973, *Journal of Geophysical Research*, 78, 4267
- Conrath, B. J., Hanel, R. A., & Kunde, V. G. 1970, *Journal of Geophysical . . .*, 75, 5831
- Cover, T., & Hart, P. 1967, *Information Theory, IEEE Transactions on*, 13, 21
- Crouzet, N., McCullough, P. R., Burke, C., & Long, D. 2012, *The Astrophysical Journal*, 761, 7
- Danielski, C., Deroo, P., Waldmann, I. P., et al. 2014, *The Astrophysical Journal*, 785, 35
- Deming, D., Wilkins, A., McCullough, P., et al. 2013, *The Astrophysical Journal*, 774, 95
- Dickey, J. M. 1971, *The Annals of Mathematical Statistics*, 42, 204
- Feroz, F., Gair, J. R., Hobson, M. P., & Porter, E. K. 2009, *Classical and Quantum Gravity*, 26, 215003
- Feroz, F., & Hobson, M. P. 2008, *Monthly Notices of the Royal Astronomical Society*, 384, 449
- Feroz, F., Hobson, M. P., Cameron, E., & Pettitt, A. N. 2013, *ArXiv e-prints*, 1
- Fletcher, L. N., Irwin, P. G. J., Teanby, N. A., et al. 2007, *ICARUS*, 189, 457
- Ford, E. B. 2006, *APJ*, 642, 505
- Foreman-Mackey, D., Hogg, D. W., Lang, D., & Goodman, J. 2013, *Publications of the Astronomical Society of the Pacific*, 125, 306
- Gevaert, W. J. R., & de With, P. H. N. 2013, in *Image Processing: Algorithms and Systems XI. Proceedings of the SPIE*, ed. K. O. Egiazarian, S. S. Agaian, & A. P. Gotchev, Technische Univ. Eindhoven (Netherlands) (SPIE), 03–865503–11

- Gibson, N. P., Aigrain, S., Roberts, S., et al. 2012, *Monthly Notices of the Royal Astronomical Society*, 419, 2683
- Goodman, J., & Weare, J. 2010, *Communications in Applied Mathematics and Computational Science*, 5, 65
- Gregory, P. C. 2011, *MNRAS*, 410, 94
- Griffith, C. A. 2014, *Philosophical Transactions of the Royal Society A: Mathematical*, 372, 30086
- Haario, H., Laine, M., Mira, A., & Saksman, E. 2006, *Statistics and Computing*, 16, 339
- Halko, N., Martinsson, P.-G., Shkolnisky, Y., & Tygert, M. 2011, [dx.doi.org](https://doi.org/10.1007/s11464-011-9258-0), 33, 2580
- Hanel, R., Conrath, B., Flasar, F. M., et al. 1981, *Science*, 212, 192
- Hanel, R. A., Conrath, B. J., Hovis, W. A., et al. 1972, *Science*, 175, 305
- Hastings, W. K. 1970, *Biometrika*, 57, 97
- Hill, C., Yurchenko, S. N., & Tennyson, J. 2013, *ICARUS*, 226, 1673
- Hollis, M. D. J., Tessenyi, M., & Tinetti, G. 2013, *Computer Physics Communications*, 184, 2351
- Irwin, P. G. J., Teanby, N. A., De Kok, R., et al. 2008, *Journal of Quantitative Spectroscopy & Radiative Transfer*, 109, 1136
- Jasa, T., & Xiang, N. 2005, in *25th AIP Conf. Proc. (AIP)*, 189–196
- Jeffreys, H. 1961, *The theory of probability* (Oxford University Press)
- Jolliffe, I. T. 2007, *Principal Component Analysis, Second Edition* (Springer Verlag, New York)
- Kass, R. E., & Raftery, A. E. 1995, *Journal of the American Statistical Association*, 90, 773
- Keeton, C. R. 2011, *Monthly Notices of the Royal Astronomical Society*, 414, 1418
- Kipping, D., & Bakos, G. 2011, *APJ*, 733, 36
- Kipping, D. M., Bakos, G. Á., Buchhave, L., Nesvorný, D., & Schmitt, A. 2012, *APJ*, 750, 115
- Knutson, H. A., Charbonneau, D., Allen, L. E., et al. 2007, *APJ*, 447, 183
- Kreidberg, L., Bean, J. L., Desert, J.-M., et al. 2014, *Nature*, 505, 69
- Lee, J. M., Fletcher, L. N., & Irwin, P. G. J. 2011, *Monthly Notices of the Royal Astronomical Society*, 420, 170
- Line, M. R., Knutson, H., Deming, D., Wilkins, A., & Desert, J.-M. 2013a, *The Astrophysical Journal*, 778, 183
- Line, M. R., Zhang, X., Vasisht, G., et al. 2012, *The Astrophysical Journal*, 749, 93
- Line, M. R., Wolf, A. S., Zhang, X., et al. 2013b, *The Astrophysical Journal*, 775, 137
- Liou, K. N. 2002, *An introduction to atmospheric radiation* (London: Academic Press)
- Madhusudhan, N., Crouzet, N., McCullough, P. R., Deming, D., & Hedges, C. 2014, *The Astrophysical Journal Letters*, 791, L9
- Madhusudhan, N., & Seager, S. 2009, *The Astrophysical Journal*, 707, 24
- Marin, J. M., & Robert, C. P. 2010, *Electronic Journal of Statistics*, 4, 643
- Martinsson, P.-G., Rokhlin, V., & Tygert, M. 2011, *Applied and Computational Harmonic Analysis*, 30, 47
- Metropolis, N., & Rosenbluth, A. W. 1953, *The journal of . . .*, 21, 1087
- Morales, J. L., & Nocedal, J. 2011, *ACM Transactions on Mathematical Software (TOMS)*, 38, 7
- Morello, G., Waldmann, I. P., Tinetti, G., et al. 2014, *The Astrophysical Journal*, 786, 22
- Mukherjee, P., Parkinson, D., & Liddle, A. R. 2006, *The Astrophysical Journal*, 638, L51
- Nelder, J. A., & Mead, R. 1965, *The computer journal*, 7, 308
- Placek, B., Knuth, K. H., & Angerhausen, D. 2013, [arXiv.org](https://arxiv.org/abs/1310.6764), 1310.6764

- Rodgers, C. D. 1976, *Reviews of Geophysics and Space Physics*, 14, 609
- Rothman, L. S., Gordon, I. E., Barbe, A., et al. 2009, *Journal of Quantitative Spectroscopy & Radiative Transfer*, 110, 533
- Rothman, L. S., Gordon, I. E., Barber, R. J., et al. 2010, *Journal of Quantitative Spectroscopy & Radiative Transfer*, 111, 2139
- Rothman, L. S., Gordon, I. E., Babikov, Y., et al. 2013, *Journal of Quantitative Spectroscopy & Radiative Transfer*, 130, 4
- Schwenke, D. W. 1998, *Chemistry and Physics of Molecules and Grains in Space. Faraday Discussions No. 109. The Faraday Division of the Royal Society of Chemistry*, 109, 321
- Seager, S. 2011, *Contemporary Physics*, 52, 602
- Skilling, J. 2004, in *24th AIP Conf. Proc., Killaha East, Kenmare, Kerry, Ireland (AIP)*, 395–405
- Skilling, J. 2006, *Bayesian Analysis*, 1, 833
- Snellen, I. A. G., de Kok, R. J., de Mooij, E. J. W., & Albrecht, S. 2010, *Nature*, 465, 1049
- Southworth, J. 2010, *Monthly Notices of the Royal Astronomical Society*, 408, 1689
- Swain, M. R., Line, M. R., & Deroo, P. 2014, *The Astrophysical Journal*, 784, 133
- Swain, M. R., Vasisht, G., & Tinetti, G. 2008, *Nature*, 452, 329
- Swain, M. R., Deroo, P., Griffith, C. A., et al. 2010, *Nature*, 463, 637
- Tennyson, J., & Yurchenko, S. N. 2012, *Monthly Notices of the Royal Astronomical Society*, 425, 21
- Ter Braak, C. J. F., & Vrugt, J. A. 2008, *Statistics and Computing*, 18, 435
- Terrile, R. J., Lee, S., Tinetti, G., et al. 2008, in *Aerospace Conference, 2008 IEEE (IEEE)*, 1–9
- Thatte, A., Deroo, P., & Swain, M. R. 2010, *Astronomy and Astrophysics*, 523, 35
- Tinetti, G., Encrenaz, T., & Coustenis, A. 2013, *The Astronomy and Astrophysics Review*, 21, 63
- Tinetti, G., Tennyson, J., Griffith, C. A., & Waldmann, I. P. 2012, *Philosophical Transactions A*, 370, 2749
- Turk, M. A., & Pentland, A. P. 1991, in *Computer Vision and Pattern Recognition, 1991. Proceedings CVPR '91., IEEE Computer Society Conference on (IEEE Comput. Soc. Press)*, 586–591
- Venot, O., Agúndez, M., Selsis, F., Tessenyi, M., & Iro, N. 2014, *Astronomy and Astrophysics*, 562, 51
- Verde, L., Feeney, S. M., Mortlock, D. J., & Peiris, H. V. 2013, *JCAP*, 09, 013
- Verdinelli, I., & Wasserman, L. 2012, *Journal of the American Statistical Association*, 90, 614
- Waldmann, I. P. 2012, *The Astrophysical Journal*, 747, 12
- . 2014, *The Astrophysical Journal*, 780, 23
- Waldmann, I. P., Tinetti, G., Deroo, P., et al. 2013, *The Astrophysical Journal*, 766, 7
- Waldmann, I. P., Tinetti, G., Drossart, P., et al. 2012, *APJ*, 744, 35
- Wark, D. Q., & Hilleary, D. T. 1969, *Science*, 165, 1256
- Weinberg, M. D. 2012, *Bayesian Analysis*, 7, 737
- Yadin, B., Veness, T., Conti, P., et al. 2012, *Monthly Notices of the Royal Astronomical Society*, 425, 34
- Yurchenko, S. N., Barber, R. J., & Tennyson, J. 2011, *Monthly Notices of the Royal Astronomical Society*, 413, 1828
- Yurchenko, S. N., & Tennyson, J. 2014, *Monthly Notices of the Royal Astronomical Society*, 440, 1649
- Zhu, C., Byrd, R. H., Lu, P., & Nocedal, J. 1997, *ACM Transactions on Mathematical ...*, 23, 550

This 2-column preprint was prepared with the AAS L^AT_EX macros v5.2.