

# Data analysis pipeline for EChO end-to-end simulations

Ingo P. Waldmann · E. Pascale

Received: 18 February 2014 / Accepted: 7 July 2014

© The Author(s) 2014. This article is published with open access at Springerlink.com

**Abstract** Atmospheric spectroscopy of extrasolar planets is an intricate business. Atmospheric signatures typically require a photometric precision of  $1 \times 10^{-4}$  in flux over several hours. Such precision demands high instrument stability as well as an understanding of stellar variability and an optimal data reduction and removal of systematic noise. In the context of the EChO mission concept, we here discuss the data reduction and analysis pipeline developed for the EChO end-to-end simulator EChOSim. We present and discuss the step by step procedures required in order to obtain the final exoplanetary spectrum from the EChOSim ‘raw data’ using a simulated observation of the secondary eclipse of the hot-Neptune 55 Cnc e.

**Keywords** EChO space-mission · Astronomical data reduction · Time resolved spectroscopy · Atmospheric spectroscopy · EChOSim

## 1 Introduction

Recent successes in characterisation of extrasolar planets are also always tales of characterising the instrument response function to an unprecedented detail. Always being at the edge of technical feasibility means that instrument calibration, observing strategy as well as data analysis and modelling are interdependent. In the light of the EChO ESA-M3 mission concept [1], such interdependence becomes important in the study of engineering decisions and instrument trade-offs. In other words, one

---

I. P. Waldmann (✉)

Department Physics & Astronomy, University College London, Gower Street,  
London WC1E 6BT, UK  
e-mail: ingo@star.ucl.ac.uk

E. Pascale

School of Physics & Astronomy, Cardiff University, Cardiff CF24 3AA, UK

needs to simulate the full observational and data analysis chain in order to gauge the impact the instrument concept has on the achievable error bar of the detection. Such a feat requires an advanced mission end-to-end simulator as well as an advanced data analysis pipeline. In this paper, we discuss the data analysis pipeline which is used in conjunction to the mission simulator, EChOSim [2]. The EChOSim data pipeline (from here on EChOSim-DP) is a stand-alone software custom built for EChOSim but with easy adaptability to other instruments and data sets in mind.

The method by which the EChO mission will characterise the nature of extrasolar planets is by time resolved spectroscopy of their atmospheres, in particular of transiting extrasolar planets. Briefly, when an exoplanet transits in front of its host star (in our line of sight) we observe a diminishing of the stellar flux due to the obscuration of the planet. The depth of the resulting lightcurve allows us to estimate the planetary radius (given the stellar radius is known). This we refer to as transit (or primary eclipse) observation. Should the exoplanet feature an extended atmosphere, we expect some of the stellar light to filter through the terminator region of the planetary atmosphere. Here we are sensitive to molecules absorbing the stellar light at specific wavelengths. We hence perceive a variation of transit depths depending on the wavelength range observed. These variations constitute the signatures of an exoplanetary absorption spectrum. Similarly, we can observe the occultation (or secondary eclipse) where the thermal contribution of the exoplanet's day-side is lost to the observer as the planet passes behind its host star. The study of transmission and emission spectroscopy is now a well established field for both space and ground based observations of exoplanetary atmospheres (e.g. [3–23] also see [24] for a comprehensive review).

## 1.1 EChOSim

EChOSim is the EChO mission end-to-end simulator. EChOSim implements a detailed simulation of the major observational and instrumental effects, and associated systematics. It also allows the influence of individual instrumental and astrophysical parameters to be studied and thus represents a key tool in the optimisation of the instrument design. Observation and calibration strategies, data reduction pipelines and analysis tools can all be designed effectively using the realistic outputs produced by EChOSim [2, 25]. The simulation output closely mimics standard STScI<sup>1</sup> FITS files, allowing for a high degree of compatibility with standard astronomical data reduction routines.

## 1.2 Examples

We illustrate individual steps in EChOSim-DP using diagrams. Unless specified otherwise, we follow a single data processing run of EChOSim simulated data of the hot-Neptune 55 Cnc e. EChOSim was run to simulate the Chemical Census mode of EChO, in which we co-add (in the case of 55 Cnc e) five eclipse observations to obtain a minimal signal-to-noise (S/N) of the final spectrum of  $S/N \sim 5$ .

---

<sup>1</sup><http://archive.stsci.edu/hst/>

For this we assume spectra reconstructed with a resolving powers of 50,50,30,30,30 for the VNIR, SWIR, MWIR-1, MWIR-2, and LWIR channels. The native resolving powers of individual detectors can exceed these requirements, see [1] for a review of the proposed EChO observing modes.

## 2 Data reduction

The EChOSim-DP is a stand-alone package delivered with the EChOSim code but can easily be adapted to observations produced by any spectrograph. It is written in fully object orientated Python allowing for a cross platform compatibility and an easy adaptability through its modular design. EChOSim-DP is subdivided into five main modules: 1) The data and parameter read-in and object initialisation, 2) data reduction, going from two dimensional focal plane illuminations to 1D time series data, 3) time series de-trending using non-parameteric de-trending algorithms, 4) lightcurve fitting using simplex-downhill minimisations as well as Markov Chain Monte Carlo (MCMC) techniques, 5) collection of results and computation of the final spectrum. We summarise this flow in Fig. 1.

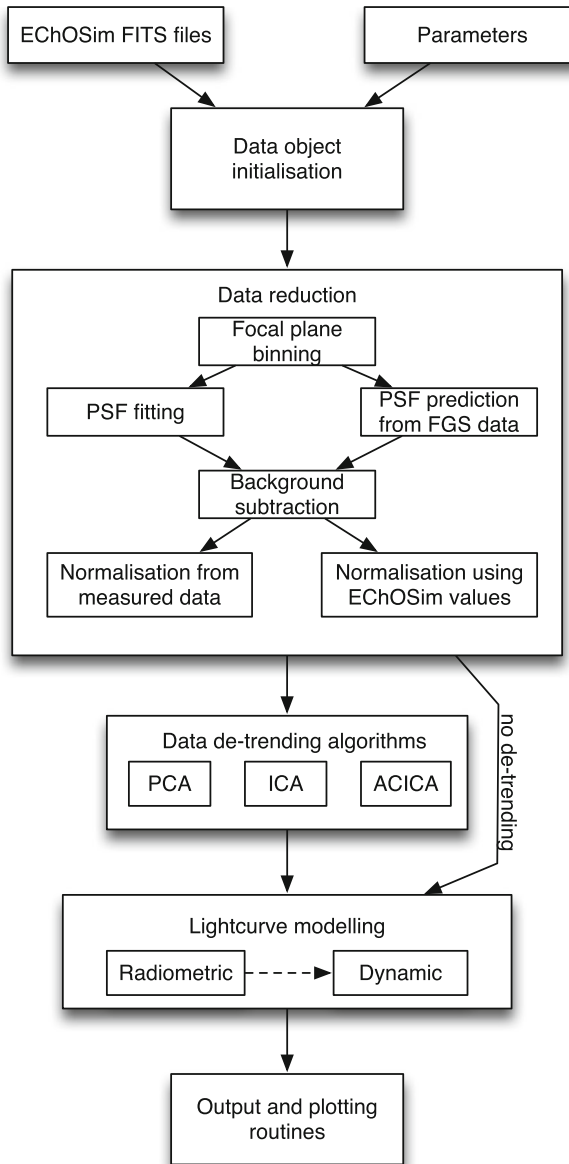
### 2.1 Configuration and data formats

The output of EChOSim follows the standard FITS file conventions with the aim to make the raw data generated by EChOSim as universally readable as possible. The payload of EChO is subdivided into individual channels defined as: VNIR (0.4 - 2.5 $\mu\text{m}$ ), SWIR (2.5 - 5.0 $\mu\text{m}$ ), MWIR-1 (5.0 - 8.5 $\mu\text{m}$ ), MWIR-2 (8.5 - 11.0 $\mu\text{m}$ ) and LWIR (11.0 - 16.0 $\mu\text{m}$ ). For a detailed description of the individual channels we refer the reader to [1] and publications in this special issue. Due to varying detector array sizes, it is not possible to combine all focal plane read-outs (for an individual frame) in one conventional FITS data-cube. EChOSim hence utilises extensions to the Primary FITS Header Data Unit (PrimaryHDU). This allows the inclusion of meta data on each detector as well as additional auxiliary information carried in binary tables (BinaryHDUs). EChOSim produces one FITS file per integration interval resulting in 10s to 100s of files per simulated observation run (Fig. 2). Whilst the high number of output files produced seems cumbersome, it reflects the data handling strategies of current space and ground based instruments. EChOSim-DP is designed to be fully compatible to this customised FITS convention using a custom build read-in routine based on the PyFITS<sup>2</sup> package. EChOSim-DP can also natively read single HDU FITS files generated by other instruments.

Auxiliary information contained in BinaryHDUs contains: the EChOSim generated stellar limb-darkening grid, EChOSim generated noiseless stellar, zodi and thermal fluxes from the instrument and its optical elements, EChOSim generated exoplanetary eclipse/transit depths, EChOSim generated Keplerian solutions. If specified

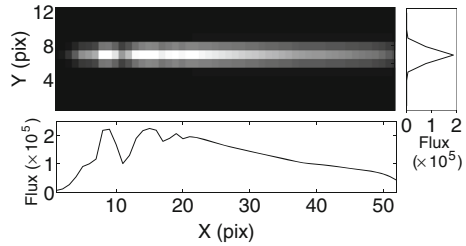
---

<sup>2</sup>[http://www.stsci.edu/institute/software\\_hardware/pyfits](http://www.stsci.edu/institute/software_hardware/pyfits)



**Fig. 1** Flowchart of the EChOSim-DP design. The pipeline is subdivided into five main modules (contained as individual python classes): 1) Object initialisation and data read, collating all input data and parameter files and performing format conversions where necessary, 2) Data reduction, reducing the two dimensional focal plane images to 1D wavelength dependent time series, 3) de-trending all or individual time-series using non-parametric machine learning techniques, 4) model fitting the final lightcurve, 5) collecting all data and calculating the final spectrum

**Fig. 2** *Top left* Focal plane of the mid-IR2 detector as read in by EChOSim-DP. *Bottom and right* cross cuts through the focal plane along the spectral and spatial directions respectively. Most flux is contained within three pixels of the spatial direction



by the user, EChOSim-DP can use these auxiliary data to calculate exact time series normalisation constants and eclipse/transit models to estimate best-case scenarios.

EChOSim-DP specific parameters are specified in a separate ascii file and parsed using the python specific ConfigParser.<sup>3</sup>

### 2.2 Flat-fielding and bad-pixel rejection

Before spectra are extracted, the focal plane data is flat-field subtracted and scanned for bad-pixels. The flat field is provided by EChOSim and constitutes a inter-pixel sensitivity variation map of the detector. In the current implementation no other flat fielding is provided by EChOSim. After flat-fielding, each frame is scanned for  $3\sigma$  flux variant pixels (either from cosmic ray hits or otherwise) which are masked and subsequently excluded from further analysis.

### 2.3 Focal plane binning

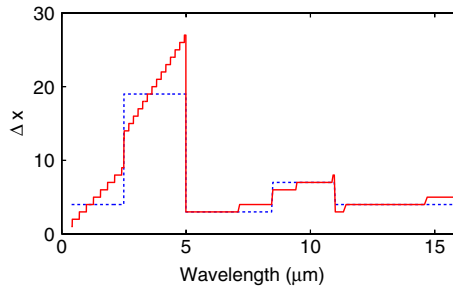
Given current detector design specifications, the native spectral resolution ( $R = \lambda/\Delta\lambda$ ) of EChO can exceed that required by the science case. EChOSim-DP provides two available spectral binning formats: 1) constant R, (1); 2) constant  $\Delta\lambda$ , (2):

$$\Delta x = \frac{\lambda R(\lambda_{mid})}{2\Delta_{pix}} \tag{1}$$

$$\Delta x = \frac{\lambda_{mid} R(\lambda_{mid})}{2\Delta_{pix}} \tag{2}$$

where  $\Delta x$  is the binning interval along the spectral axis in pixels,  $\lambda$  and  $\lambda_{mid}$  the wavelength and central wavelength in  $\mu\text{m}$  and  $\Delta_{pix}$  the pixel size in  $\mu\text{m}$ . Note that EChO spectrometers sample each spectral resolving element with two detector pixels. Figure 3 shows  $\Delta x$  for both binning methods as function of  $\lambda$ . Binning is performed directly on the focal plane before spectral extraction. This increases S/N and avoids potential biasing of the data.

<sup>3</sup><http://docs.python.org/2/library/configparser.html>

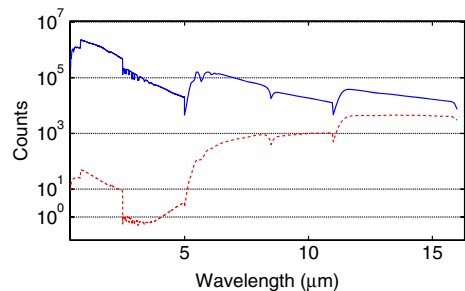


**Fig. 3** Showing binning steps  $\Delta x$  in pixels as function of wavelength for the two spectral binning modes available in EChOSim-DP. The native resolution for all detectors (VIS, NIR, MIR1, MIR2, FIR) are  $R = 330, 530, 52, 103, 62$ . Red-solid line shows the constant in  $R$  binning; blue-discontinuous line shows constant in  $\Delta\lambda$  binning to resolutions of  $R = 50, 50, 30, 30, 30$  respectively

## 2.4 Optimal extraction

After the data has been binned, we extract the raw spectrum along the spatial axis for each individual time stamp. At each integration time, the raw spectrum is extracted from the data by fitting a model of the PSF to the point-like dispersed signal of the star + planet flux. Two standard extraction techniques are available in EChOSim-DP: Photometric window extraction and optimal extraction. The photometric window extraction is the simplest spectral extraction technique which consists of summing detector counts contained in a box of fixed spatial axis width. This method is very robust in low background observations and when the instrument PSF is not known with adequate precision. Optimal extraction weighs individual pixel columns with the optimal PSF of the detector and creates a very tightly fit ‘extraction window’. This method is preferable in high background observations when the instrument PSF is well determined. EChO will have a well characterised PSF across individual detectors. Here optimal extraction techniques are preferable since observations in the mid to far-IR channels can feature significant zodiacal and thermal backgrounds as well as increased dark current rates (Fig. 4). For the remainder of this paper we will only consider optimal extraction techniques.

**Fig. 4** Showing the extracted flux for a single frame as function of wavelength. *Blue-continuous line* Optimally extracted flux before background subtraction; *red-discontinuous line* estimated background counts measured on the off-axis spatial direction



Two extraction options are available: 1) Unconstraint PSF, 2) EChOSim PSF with Fine Guidance Sensor (FGS) offset data.

*Option 1:* is the least constraint extraction. Depending on user input, EChOSim-DP fits a Gaussian or Generalised Gaussian Distribution (GGD) PSF along the spatial axis. The GGD is given by

$$PSF_{ggd} = \frac{\beta}{2\alpha\Gamma(1/\beta)} \exp - [ |(\mu_y + \Delta y(t)) - y|/\alpha ]^\beta \quad (3)$$

where  $\mu_y$  is the mean position of the spectrum along the spatial axis  $y$  for all frames,  $\Delta y(t)$  is a time dependent offset from the mean,  $\alpha$  is a scale parameter and in this case equivalent to  $\alpha = 2\sigma_y$  and  $\sigma_y$  signifies the width of the PSF. The shape parameter  $\beta$  introduces a kurtosis argument in the Gaussian distribution. We retrieve the Normal PSF by setting  $\beta = 2$  and obtain leptokurtic and platykurtic distributions for  $\beta < 2$  and  $\beta > 2$  respectively. We do not assume skew of the PSF in the spatial direction. The PSF shape can either be left as free parameter (to be fitted from the data) or specified as user input. Equation (3) is convolved with the detector response function assumed by EChOSim to obtain the extraction profile.

$$\mathcal{P}(y, t) = PSF(y, t) \otimes \mathcal{R}(y) \quad (4)$$

where  $\otimes$  is the convolution operator and the detector response [26] is given by

$$\begin{aligned} \mathcal{R}(y; \Delta_{pix}, l_y) &= \\ &= \frac{\tan^{-1} \left( \tanh \left( \frac{\Delta_{pix} - y}{4l_y} \right) \right) - \tan^{-1} \left( \tanh \left( -\frac{\Delta_{pix} - y}{4l_y} \right) \right)}{\tan^{-1} \left( \tanh \left( \frac{\Delta_{pix}}{4l_y} \right) \right) - \tan^{-1} \left( \tanh \left( -\frac{\Delta_{pix}}{4l_y} \right) \right)} \end{aligned} \quad (5)$$

where  $\Delta_{pix}$  is the pixel size in  $\mu\text{m}$  and  $l_y$  the diffusion length in  $\mu\text{m}$ .

*Option 2:* Here we assume a Gaussian PSF (by setting  $\beta = 2$ ) with a fixed width given by  $\sigma_y = F\#K_y\lambda$  where  $F\#$  is the effective focal length of the telescope in  $\mu\text{m}$ ,  $K_y$  is the PSF aberration parameter and  $\lambda$  the wavelength in  $\mu\text{m}$ . We hold  $\mu_y$  fixed at an EChOSim specified value and obtain the time dependent offset  $\Delta y(t)$  from the EChOSim provided fine guidance sensor (FGS) centroiding.

We note that for current simulations we use a Gaussian PSF. This is through lack of calibration data of the instrument in the current study phase. EChOSim-DP natively supports the inclusion of more realistic PSF functions available in future simulations.

The centroiding is provided as part of the auxiliary information BinaryHDUs and consists of a time series of  $y$ -positional offsets sampled at 1Hz frequency. EChOSim-DP downsamples the positional offsets to the integration times specified in the FITS headers. The downsampling operation correctly reflects the error in the positional offset  $\Delta y(t)$  and the associated flux error.

## 2.5 Background subtraction

EChOSim-DP calculates the background by computing the median (or mean given user input) focal plane illumination  $4\sigma_y$  away from  $\mu_y$ . The background flux is integrated over the area (in pixels) of the extraction profile and subtracted from the extracted flux.

## 2.6 PSF instabilities

Simulations of PSF variability due to pointing jitter have shown to result in an overall flux error of  $\sim 10^{-5} - 1 \times 10^{-4}$  but significantly higher for the spectral ranges of the NIR instrument (2-5  $\mu\text{m}$ ) where uncorrected flux errors can reach  $5 \times 10^{-3}$  levels. This is to be expected as the SWIR instrument features a smaller pixel size. Effects due to telescope and optical bench thermal drifts are found to be negligible in the wavelength ranges below 14  $\mu\text{m}$  and temperatures below 50K. We refer the reader to [2, 25] for further information.

Intra-visit (i.e. within the observation of an eclipse/transit event) thermal-mechanical distortions and/or other external forcing functions can introduce additional noise on the FGS centroiding information. This has been accounted for by adding a Gaussian centroiding error with a 10 milli-arcsecond rms amplitude, following the outcome of the industrial studies (priv. com.). Inter-visit (i.e. from observation of one eclipse/transit to the next) variations in the FGS PSF are not considered as drifts can be calibrated upon acquisition of the target.

## 2.7 Normalisation

The final step is the normalisation of the data to the out of transit (OOT) baseline. Similarly to Section 2.7 the normalisation can either be estimated from the data itself by calculating the OOT mean or normalised using noiseless stellar fluxes provided by EChOSim

$$F_{norm}(\lambda, t) = \frac{F_{total}(\lambda, t)}{F_{star}(\lambda, t)}. \quad (6)$$

The noiseless flux measure provided by EChOSim allows the idealised case to be calculated where a perfect knowledge of the stellar spectrum (and activity) is assumed. We discuss the more complex case in Section 4.2.

## 3 Data de-trending

After the data has been reduced to 1D time series, EChOSim-DP can attempt a de-correlation of wavelength correlated non-Gaussian systematics. These systematics tend to be due to array wide fluctuations of quantum efficiencies, insufficient flat-fielding, slit-loss effects and pointing jitter. These complex non-Gaussian signals have shown to be important effects in real instruments [14, 27–29]. EChOSim implements inter and intra-pixel variations and non-Gaussian pointing jitter noise. Other



non-linear noise sources such as correlated astrophysical noise (e.g. such as stellar pulsation, stellar spots and faculae noise) will be included in future releases of EChOSim.

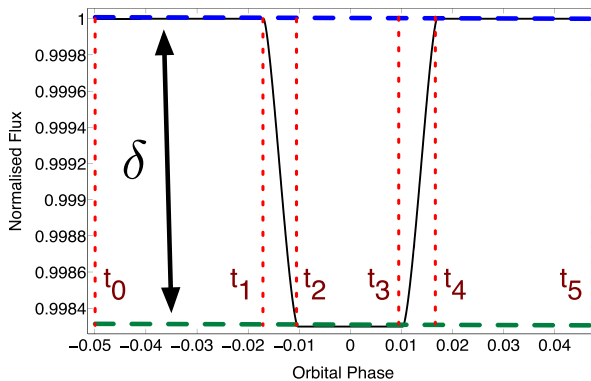
Here we implement the ACICA de-trending algorithm [29]. Based on blind-deconvolution using Independent Component Analysis [27, 30, 31], we estimate the common non-Gaussian time and wavelength correlated signals and construct a systematic noise model which is then used to correct each individual time series. The advantage of these types of de-trending algorithms over others such as Gaussian Processes [32] are their non-parametric nature. This guarantees a high degree of objectivity in the de-trending as well as a simple implementation into existing code (due to the lack of parameterisation required).

## 4 Lightcurve modelling

Once the data is reduced and de-correlated, the pipeline provides several means of model fitting the resulting lightcurves. The modelling is divided into two main modes: Radiometric and Dynamic. In the simplest model assumption, the radiometric case, we simply calculate the error bar from the out-of-transit (OOT) scatter of the time series and estimate the transit depth by taking the ratio of in-transit (IT) and OOT data. For the Dynamic case we use a full transiting planet model [33] and iteratively fit for the transit depth parameter using a simplex-downhill algorithm as well as a Markov Chain Monte Carlo (MCMC) routine.

### 4.1 Radiometric data analysis

For most cases, and for the sake of computational efficiency, the simplistic radiometric model results are desired for EChOSim observations. Let us assume a secondary



**Fig. 5** Single Mandel & Agol (2002) eclipse model. The discontinuous blue line marks the out of transit baseline. The discontinuous green line marks the in-transit flux and  $\delta$  defines the transit depth. Discontinuous red lines note the contact points  $t_{1-4}$

eclipse measurement of an exoplanet. In the simplest radiometric case, we calculate the transit depth via the simple relation

$$\delta = F_{out} - F_{in} \quad (7)$$

where  $\delta$  is the transit depth,  $F_{out}$  is the baseline flux (blue line in Fig. 5) and is defined as

$$F_{out} = \frac{1}{N(t_{0-1}, t_{4-5})} \left( \sum_{t=t_0}^{t_1} F_t + \sum_{t=t_4}^{t_5} F_t \right) \quad (8)$$

where  $t$  is the time index,  $N$  the number of observations in time,  $t_{0-1}$  defines pre-ingress baseline time and  $t_{4-5}$  post-egress timeline (see Fig. 5). Similarly we define the in-transit flux as

$$F_{in} = \frac{1}{N(t_{2-3})} \sum_{t=t_2}^{t_3} F_t \quad (9)$$

Equation (9) is valid for the secondary eclipse case and mid-IR transit cases where limb-darkening is negligible. To avoid the effect of limb-darkening in the case of primary eclipses in the near-IR, we borrow the ‘correct’ transit depths from EChOSim’s auxiliary output files. Note that this is a valid procedure since we are dealing with an over simplistic model here. The dynamic model fitting does not assume auxiliary data. Given (7), we calculate the error on  $\delta$  as the sum of squares of the time series error

$$\sigma_{total}/\sqrt{N} = \sqrt{\sigma_{out}^2/N_{out} + \sigma_{in}^2/N_{in}} = \frac{\sqrt{2}\sigma}{\sqrt{N}} \quad (10)$$

where  $N$  is the number of observations and we assume that  $N_{out} = N_{in} = 2N$  as well as  $\sigma_{out} = \sigma_{in}$ .

#### 4.1.1 Interpretation of radiometric model

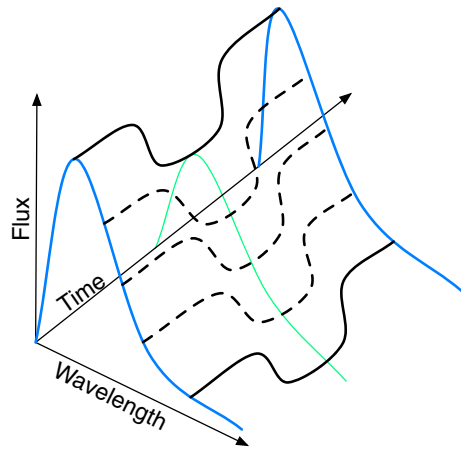
The assumption  $\sigma_{out} = \sigma_{in}$  seems straight forward as one expects the photometric stability not to vary significantly between out-of-eclipse and in-eclipse times. The radiometric error as in (10) is the correct error treatment for the observation of a single lightcurve at a single wavelength with equal lengths of out-of-transit and in-transit measurements. It assumes that no additional knowledge of the baseline (out-of-transit) flux is available and describes the state of largest ignorance, i.e.  $\sigma \rightarrow \sqrt{2}\sigma$ . Should additional knowledge of the baseline flux be available (via the calibration of the wavelength dependent stellar spectrum), we can reduce the normalisation error on the baseline. Hence for a perfect knowledge of the baseline flux level  $\sigma_{total} \rightarrow \sigma$ .

#### 4.2 Dynamic data analysis

Going beyond the radiometric model assumptions, EChOSim-DP has two additional time-resolved lightcurve model modes: 1) Simplex and 2) MCMC (Fig. 6).

In the simplex case, we fit an analytical lightcurve model [33] to each individual lightcurve in wavelength space,  $\lambda$ . It fully supports eccentric orbit calculations

**Fig. 6** Schematic outline of EChO observations illustrating the changing baseline flux levels. Here *blue curves* illustrate the stellar out of transit spectra and the *green curve* the in-transit spectrum of the star. In the case of a secondary eclipse, the green curve represents the stellar spectrum only whilst the blue curve is star+planetary flux



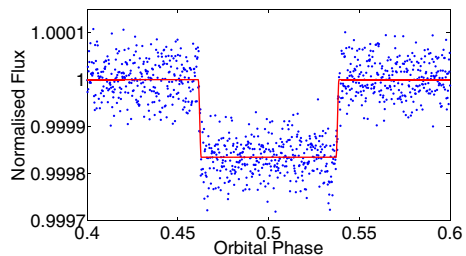
following [34] and allows all model parameters to be fitted. For lightcurves in wavelength ranges below  $5\mu\text{m}$  we assume stellar limb-darkening for primary eclipses. Here we linearly interpolate the quadratic limb-darkening coefficients of [35] or read the limb-darkening coefficient grid provided by EChOSim to provide an exact match. For the model minimisation we use a simplex-downhill algorithm [36, 37]. In this simple minimisation scheme, we obtain the error bar on the model fit using (10). Each modelling run creates a new model-fitting object in the data pipeline which allows multiple model runs (radiometric as well as dynamic) to be executed in the same instance of the EChOSim-DP.

We furthermore include a more computationally intensive Markov Chain Monte Carlo routine in EChOSim-DP. This routine allows us to investigate more complex scenarios and potential prior dependence (should prior knowledge on the exoplanetary or stellar spectrum be known). The posterior on the model parameter  $\theta$  can be written as

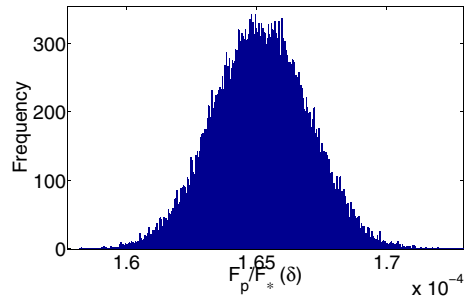
$$p(\hat{\theta}|F) \propto \mathcal{L}(\theta)\pi(\theta). \tag{11}$$

where  $\mathcal{L}(\theta)$  is the model likelihood and  $\pi(\theta)$  the prior distribution on the parameter  $\theta$  (Figs. 7 and 8). Whilst we here only consider  $\theta$  to be the transit depth parameter, EChOSim-DP natively supports the inclusion of other free parameters, such as orbital (e.g. ephemeris, eccentricity, orbital inclination) as well as free-floating limb-darkening parameters. In a typical EChO observation, these additional parameters are

**Fig. 7** Normalised *lightcurve* of secondary eclipse of 55 Cnc e (5 eclipses co-added). *Red line* analytic lightcurve model [33] with the eclipse depth  $\delta$  as only free parameter. Note the lack of stellar limb-darkening in secondary eclipses and hence a very discrete ingress and egress



**Fig. 8** Histogram of MCMC chain run for 50,000 iterations. The histogram approximates the posterior distribution of the transit depth parameter  $\delta$  for the model fit shown in Fig. 7



thought to be well determined by previous studies and are assumed to be fixed. The likelihood is here assumed to be Gaussian and is given by

$$\mathcal{L}(\theta, \mathbf{d}) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \sum_{t=0}^{t=N} \left( \frac{d_t - \Phi_t(\theta)}{\sigma} \right)^2 \right] \quad (12)$$

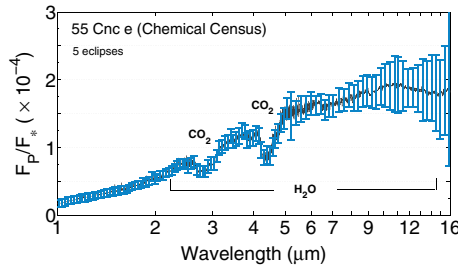
where  $\mathbf{d}$  is the data column vector, and  $d_t$  and  $\Phi_t(\theta)$  are the datum and lightcurve model at given time-stamp  $t$ .

We use the PyMC<sup>4</sup> package implementing the adaptive Metropolis Hastings algorithm of [38]. The MCMC chains are typically run with 20,000 iterations taking the minimised result of the simplex-downhill algorithm as starting value to minimise burn-in time [39] which we restrict to 1000 iterations. We here present the univariate version of the likelihood as in most cases all transit parameters but the depth,  $\delta$ , are fixed. To minimise parameter covariances for multiple free parameters one can follow parameterisation by [40] or [41]. Using a Bayesian approach, we can investigate more complex model solutions such as the impact of the stellar variability on the normalisation of individual lightcurves. Figure 6 illustrates a time series observation of a transiting exoplanet over a wide range of wavelengths. Here the blue curves represent the stellar spectrum, the black curves the time dependent flux variation due to the transiting extrasolar planet with the green line marking the minimum flux. As discussed in Section 4.1.1, if all time series measurements are assumed to be independent of each other (i.e. not correlated in wavelength), we must assume an error of  $\sqrt{2}\sigma$  on the measurement, given the uncertainty of the OOT normalisation. However, it is clear from Fig. 6 that OOT flux of individual time series is correlated in  $\lambda$  through the stellar spectrum. For a perfect correlation (i.e. absolute knowledge on the correct normalisation of the individual time series) the measurement error hence reduces to  $\sigma$ . Hence the normalisation error,  $\sigma_{norm}$ , is bound by  $0 \leq \sigma_{norm} \leq \sqrt{2}$ .

We can now express the likelihood of our observation,  $\mathcal{L}$ , as product of the likelihood of the lightcurve model,  $\mathcal{L}(\theta)$  and the stellar spectrum model  $\mathcal{L}(\varphi)$ . Note that by taking the product we implicitly assume statistical independence between lightcurve and stellar spectra models and below we explicitly assume a Gaussian noise model

$$\mathcal{L} = \mathcal{L}(\theta)\mathcal{L}(\varphi) = e^{-\frac{1}{2}\chi^2(\theta)} e^{-\frac{1}{2}\chi^2(\varphi)} \quad (13)$$

<sup>4</sup><https://github.com/pymc-devs/pymc>

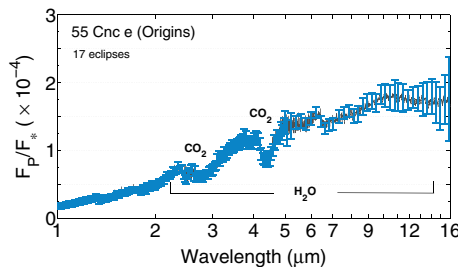


**Fig. 9** Final spectrum generated from EChOSim-DP outputs for 55 Cnc e secondary eclipse run in Chemical census mode (i.e. 5 eclipses stacked,  $R = 50$  for  $\lambda < 5\mu\text{m}$  and  $R = 30$  for  $\lambda > 5\mu\text{m}$ ). Blue error bars derived from EChOSim-DP. Grey: planetary emission spectrum read into EChOSim. We marked prominent emission/absorption features

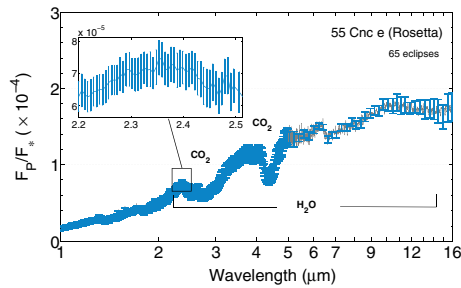
where  $\chi^2$  is the chi-squared distribution. We can now write the log-likelihood as follows

$$\log \mathcal{L} = -\frac{1}{2} \sum_{t=1}^N \left( \frac{F_{t,\lambda} - \Phi(\theta_t)}{\sigma_t} \right)^2 - \frac{1}{2} \sum_{\lambda=1}^M \left( \frac{\bar{F}_{t=t_2-3,\lambda} - \Psi(\theta_\lambda)}{\sigma_\lambda} \right)^2 \tag{14}$$

where  $\Phi(\theta_t)$  is the lightcurve model for given time index  $t$ ,  $\Psi(\theta_\lambda)$  is the stellar model for given wavelength index  $\lambda$ ,  $M$  is the number of resolution elements in the spectrum and  $\sigma_t$  and  $\sigma_\lambda$  are the flux uncertainties on the time series and the stellar spectrum respectively. Note that these error terms are not equivalent and also note that  $\bar{F}_{t=t_2-3,\lambda}$  is the averaged stellar spectrum from time interval  $t_2 - t_3$ .



**Fig. 10** Final spectrum generated from EChOSim-DP outputs for 55 Cnc e secondary eclipse run in Origin mode (i.e. 17 eclipses stacked,  $R = 100$  for  $\lambda < 5\mu\text{m}$  and  $R = 30$  for  $\lambda > 5\mu\text{m}$ ). Blue error bars derived from EChOSim-DP. Grey: planetary emission spectrum read into EChOSim. We marked prominent emission/absorption features



**Fig. 11** Final spectrum generated from EChOSim-DP outputs for 55 Cnc e secondary eclipse run in Rosetta mode (i.e. 65 eclipses stacked,  $R = 300$  for  $\lambda < 5\mu\text{m}$  and  $R = 30$  for  $\lambda > 5\mu\text{m}$ ). Blue error bars derived from EChOSim-DP. Grey: planetary emission spectrum read into EChOSim. Inset is a zoom into the 2.2 – 2.5  $\mu\text{m}$  wavelength region. We marked prominent emission/absorption features

## 5 Outputs

Two types of outputs are provided: spectra in ascii format and python-pickle<sup>5</sup> objects. For each individual lightcurve fitting, EChOSim-DP provides an ascii file containing wavelength, measured flux and error. The pickle file contains all parameters, intermediate and final data products allowing for an exact reproducibility of results. Figure 9 shows the final spectrum for 55 Cnc e in the Chemical Census mode (blue error bars). Figures 10 and 11 show the same simulation for the Origins and Rosetta stone observing modes of EChO.

## 6 Discussion & conclusion

EChOSim-DP is a custom built data reduction and analysis pipeline for the EChOSim end-to-end mission simulator of the EChO mission concept.

Despite its customised nature, we have developed the pipeline with easy adaptability (through its fully object-orientated programming) to other instruments and data-sets in mind. The pipeline features state of the art data de-correlation algorithms as well as a full Bayesian analysis implementation via adaptive MCMC. Both these aspects, the de-trending as well as the exploration of stellar variability are not required for the current version of EChOSim (version 3.x) but included with future releases. These releases will have special emphasis on realistic stellar noise simulations [42] as well as more advanced non-Gaussian instrument systematics.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution License which permits any use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

<sup>5</sup><http://docs.python.org/2/library/pickle.html>

## References

1. Tinetti, G., Beaulieu, J.P., Henning, T., Meyer, M., Micela, G., Ribas, I., Stam, D., Swain, M., Krause, O., Ollivier, M., Pace, E., Swinyard, B., Aylward, A., Boekel, R., Coradini, A., Encrenaz, T., Snellen, I., Zapatero-Osorio, M.R., Bouwman, J., Cho, J.Y.K., de Foresto, V.C., Guillot, T., Lopez-Morales, M., Mueller-Wodarg, I., Pallé, E., Selsis, F., Sozzetti, A., Ade, P.A.R., Achilleos, N., Adriani, A., Agnor, C.B., Afonso, C., Prieto, C.A., Bakos, G., Barber, R.J., Barlow, M., Batista, V., Bernath, P., Bézard, B., Bordé, P., Brown, L.R., Cassan, A., Cavarroc, C., Ciaravella, A., Cockell, C., Coustenis, A., Danielski, C., Decin, L., De Kok, R., Demangeon, O., Deroo, P., Doel, P., Drossart, P., Fletcher, L.N., Focardi, M., Forget, F., Fossey, S., Fouqué, P., Frith, J., Galand, M., Gaulme, P., Hernández, J.I.G., Grasset, O., Grassi, D., Grenfell, J.L., Griffin, M.J., Griffith, C.A., Grözing, U., Guedel, M., Guio, P., Hainaut, O., Hargreaves, R., Hauschildt, P.H., Heng, K., Heyrovsky, D., Hueso, R., Irwin, P., Kaltenecker, L., Kervella, P., Kipping, D., Koskinen, T.T., Kovacs, G., Barbera, A., Lammer, H., Lellouch, E., Leto, G., Lopez-Morales, M., Valverde, M.A.L., Lopez-Puertas, M., Lovis, C., Maggio, A., Maillard, J.P., Prado, J.M., Marquette, J.B., Martin-Torres, F.J., Maxted, P., Miller, S., Molinari, S., Montes, D., Moro-Martin, A., Moses, J.I., Mousis, O., Tuong, N.N., Nelson, R., Orton, G.S., Pantin, E., Pascale, E., Pezzuto, S., Pinfield, D., Poretti, E., Prinja, R., Prisinzano, L., Rees, J.M., Reiners, A., Samuel, B., Sanchez-Lavega, A., Forcada, J.S., Sasselov, D., Savini, G., Sicardy, B., Smith, A., Stixrude, L., Strazzulla, G., Tennyson, J., Tessenyi, M., Vasisht, G., Vinatier, S., Viti, S., Waldmann, I., White, G.J., Widemann, T., Wordsworth, R., Yelle, R., Yung, Y., Yurchenko, S.N.: *Exp. Astron.* **34**(2), 311 (2012)
2. Pascale, E., Waldmann, I., et al.: in prep.
3. Beaulieu, J.P., Kipping, D.M., Batista, V., Tinetti, G., Ribas, I., Carey, S., Noriega-Crespo, J.A., Griffith, C.A., Campanella, G., Dong, S., Tennyson, J., Barber, R.J., Deroo, P., Fossey, S.J., Liang, D., Swain, M.R., Yung, Y., Allard, N.: *ApJ* **409**, 963 (2010)
4. Beaulieu, J.P., Tinetti, G., Kipping, D.M., Ribas, I., Barber, R.J., Cho, J.Y.K., Polichtchouk, I., Tennyson, J., Yurchenko, S.N., Griffith, C.A., Batista, V., Waldmann, I., Miller, S., Carey, S., Mousis, O., Fossey, S.J., Aylward, A.: *ApJ* **731**, 16 (2011)
5. Charbonneau, D., Knutson, H.A., Barman, T., Allen, L.E., Mayor, M., Megeath, S.T., Queloz, D., Udry, S.: *ApJ* **686**, 1341 (2008)
6. Brogi, M., Snellen, I.A.G., de Kok, R.J., Albrecht, S., Birkby, J., de Mooij, E.J.W.: *Nature* **486**(7404), 502 (2012)
7. Bean, J.: *Nature* **478**(7367), 41 (2011)
8. Swain, M.R., Vasisht, G., Tinetti, G.: *Nature* **452**, 329 (2008)
9. Swain, M.R., Bouwman, J., Akeson, R.L., Lawler, S., Beichman, C.A.: *ApJ* **674**, 482 (2008)
10. Swain, M.R., Vasisht, G., Tinetti, G., Bouwman, J., Chen, P., Yung, Y., Deming, D., Deroo, P.: *ApJL* **690**, L114 (2009)
11. Crouzet, N., McCullough, P.R., Burke, C., Long, D.: *Astrophys. J.* **761**(1), 7 (2012)
12. Deming, D., Wilkins, A., McCullough, P., Burrows, A., Fortney, J.J., Agol, E., Dobbs-Dixon, I., Madhusudhan, N., Crouzet, N., Desert, J.M., Gilliland, R.L., Haynes, K., Knutson, H.A., Line, M., Magic, Z., Mandell, A.M., Ranjan, S., Charbonneau, D., Clampin, M., Seager, S., Showman, A.P.: *Astrophys. J.* **774**(2), 95 (2013)
13. Grillmair, C.J., Burrows, A., Charbonneau, D., Armus, L., Stauffer, J., Meadows, V., Van Cleve, J., von Braun, K., Levine, D.: *Nature* **456**, 767 (2008)
14. Thatte, A., Deroo, P., Swain, M.R.: *ApJ* **523**, A35 (2010)
15. Tinetti, G., Vidal-Madjar, A., Liang, M.C., Beaulieu, J.P., Yung, Y., Carey, S., Barber, R.J., Tennyson, J., Ribas, I., Allard, N., Ballester, G.E., Sing, D.K., Selsis, F.: *Nature* **448**(7150), 169 (2007)
16. Pont, F., Knutson, H., Gilliland, R.L., Moutou, C., Charbonneau, D.: *MNRAS* **385**(1), 109 (2008)
17. Swain, M., Deroo, P., Tinetti, G., Hollis, M., Tessenyi, M., Line, M., Kawahara, H., Fujii, Y., Showman, A., Yurchenko, S.: *ApJ astro-ph*.EP (2012)
18. Knutson, H.A., Madhusudhan, N., Cowan, N.B., Christiansen, J.L., Agol, E., Deming, D., Desert, J.M., Charbonneau, D., Henry, G.W., Homeier, D., Langton, J., Laughlin, G., Seager, S.: *ApJ* **735**, 27 (2011)
19. Sing, D.K., Pont, F., Aigrain, S., Charbonneau, D., Desert, J.M., Gibson, N., Gilliland, R., Hayek, W., Henry, G., Knutson, H., des Etangs, A.L., Mazeh, T., Shporer, A.: *MNRAS* **416**, 1443 (2011)

20. Tinetti, G., Deroo, P., Swain, M.R., Griffith, C.A., Vasisht, G., Brown, L.R., Burke, C., McCullough, P.: *ApJ* **712**, L139 (2010)
21. de Mooij, E.J.W., Brogi, M., de Kok, R.J., Koppenhoefer, J., Nefs, S.V., Snellen, I.A.G., Greiner, J., Hanse, J., Heinsbroek, R.C., Lee, C.H., van der Werf, P.P.: *Astron. Astrophys.* **538**, 46 (2012)
22. Bean, J.L., Desert, J.M., Kabath, P., Stalder, B., Seager, S., Miller-Ricci Kempton, E., Berta, Z.K., Homeier, D., Walsh, S., Seifahrt, A.: *ApJ* **743**, 92 (2011)
23. Stevenson, K.B., Harrington, J., Nymeyer, S., Madhusudhan, N., Seager, S., Bowman, W.C., Hardy, R.A., Deming, D., Rauscher, E., Lust, N.B.: *ApJ* **464**, 1161 (2010)
24. Tinetti, G., Encrenaz, T., Coustenis, A.: *Astron. Astrophys. Rev.* **21**(1), 63 (2013)
25. Waldmann, I.P., Pascale, E., Swinyard, B., Tinetti, G., Amaral-Rogers, A., Spencer, L., Tessenyi, M., Ollivier, M., Coudé du Foresto, V.: *arXiv* p. 6425 (2013)
26. Barron, N., Borysow, M., Beyerlein, K., Brown, M., Lorenzon, W., Schubnell, M., Tarlé, G., Tomasch, A., Weaverdyck, C.: vol. 119
27. Waldmann, I.P.: *ApJ* **747**(1), 12 (2012)
28. Waldmann, I.P., Tinetti, G., Deroo, P., Hollis, M.D.J., Yurchenko, S.N., Tennyson, J.: *Astrophys J* **766**(1), 7 (2013)
29. Waldmann, I.P.: *ApJ* **780**(1), 23 (2013)
30. Hyvärinen, A., Oja, E.: *Neural Netw.* (2000)
31. Comon, P., Jutten, C.: *Handbook of Blind Source Separation*. Academic Press (2010)
32. Rasmussen, C.E., Williams, C.K.I.: *Gaussian Processes for Machine Learning*. The MIT Press (2006)
33. Mandel, K., Agol, E.: *ApJL* **580**, L171 (2002)
34. Kipping, D.M.: *ApJ* **389**, 1383 (2008)
35. Claret, A.: *A&A* **363**, 1081 (2000)
36. Nelder, J.A., Mead, R.: *Comput. J.* **7**(4), 308 (1965)
37. Press, W.H., Teukolsky, S.A., Vetterling, W.T., Flannery, B.P. *Numerical Recipes*, 3rd. Cambridge University Press (2007)
38. Haario, H., Laine, M., Mira, A., Saksman, E. *Stat. Comput.* **16**(4), 339 (2006)
39. Brooks, S., Gelman, A., Jones, G., Meng, X.L.: *Handbook of Markov Chain Monte Carlo*. CRC Press (2011)
40. Bakos, G.Á., Pál, A., Torres, G., Sipocz, B., Latham, D.W., Noyes, R.W., Kovacs, G., Hartman, J., Esquerdo, G.A., Fischer, D.A., Johnson, J.A., Marcy, G.W., Butler, R.P., Howard, A., Sasselov, D.D., Stefanik, R.P., Lazar, J., Papp, I., Sari, P.: *ApJ* (2008)
41. Burke, C.J., McCullough, P.R., Valenti, J.A., Johns-Krull, C.M., Janes, K.A., Heasley, J.N., Summers, F.J., Stys, J.E., Bissinger, R., Fleenor, M.L., Foote, C.N., Garcia-Melendo, E., Gary, B.L., Howell, P.J., Mallia, F., Masi, G., Taylor, B., Vanmunster, T.: *ApJ* **671**, 2115 (2007)
42. Ballerini, P., Micela, G., Lanza, A.F., Pagano, I.: *Astron. Astrophys.* **539**, A140 (2012)