

# Efficient inference in a weakly non-linear Bayesian inverse Problem using an affine approximation

Lennart Golks and Colin Fox

November 2025

## 1 Introduction

There are satellites, which orbit around the earth at a height of around 500km above the ground, carry measurement devices to determine trace gas concentrations in the stratosphere. More specifically, atmospheric limb-sounders are pointing through the atmosphere and detect thermal radiation of trace gases, e.g., ozone. Examples of such limb-sounders are the Microwave Limb Sounder (MLS) on NASA’s Aura mission [29] and the Michelson Interferometer for Passive Atmospheric Sounding (MIPAS) on ESA’s Envisat [19].

Measurements of such devices can be described by the radiative transfer equation (RTE), which in this case is a path integral along the satellite’s pointing direction. This path integral includes an absorption term accounting for the re-attenuation of the thermal radiation along the satellite’s line of sight. That makes inferring the trace gas concentration from a set of measurements a non-linear inverse problem.

Existing methods use regularisation methods to retrieve trace gas concentration from measurements [21, 14, 16]. These methods do not include hyper-parameters, such as noise covariance, and produce biased results [9].

Firstly, given some simulated data, we treat this non-linear problem as a linear problem by neglecting the absorption term in the RTE. A linear-Gaussian hierarchical Bayesian framework is employed to infer the ozone concentration in the stratosphere. This includes establishing a hierarchical structure and classifying hyper-parameters and parameters. For efficient inference we employ the marginal-and-then-conditional (MTC) scheme as in [8]. This gives a low-dimensional marginal posterior probability distribution over the hyper-parameters and a high-dimensional conditional posterior probability for the ozone parameter. We show that evaluating the marginal posterior on a grid gives the conditional posterior mean at no additional cost. To quantify the conditional posterior variance, samples from the conditional posterior are drawn via the randomise-then-optimize (RTO) scheme [1].

Using the results of the linearised problem, an affine map approximating the non-linear forward model is obtained. We employ the same hierarchical Bayesian

framework as previously used, but with the approximated forward model, to quantify the unbiased mean and variance of the posterior ozone profile.

All programming and analysis in this paper are done in Python, and the reported computation times are taken on a MacBook Pro from 2019 with a 2.4 GHz quad-core Intel i5 processor.

## 2 Hierarchical Bayesian Modelling

First, the concept of hierarchical Bayesian modelling is introduced. Assume we observe some data

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \boldsymbol{\eta}, \quad (1)$$

based on the forward model  $\mathbf{A}(\mathbf{x})$  with an unknown parameter vector  $\mathbf{x}$  and some additive random noise  $\boldsymbol{\eta}$ . Naturally, due to the noise, the observation process in Eq. 1 is a random process. Hence, in Bayesian modelling, the aim is to determine a probability distribution over the parameter  $\mathbf{x}$  given some data  $\mathbf{y}$ . Further, a hierarchical Bayesian model incorporates (auxiliary) hyper-parameters  $\boldsymbol{\theta}$  (see Fig. 1 for a schematic representation). Within a Bayesian approach all unknown hyper-parameters and parameters are treated as random variables [12, Chapter 3].

According to Bayes' theorem, the joint posterior distribution over the parameters  $\mathbf{x}$  and the hyper-parameter  $\boldsymbol{\theta}$  is given as

$$\pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) = \frac{\pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}, \boldsymbol{\theta})}{\pi(\mathbf{y})} \propto \pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}, \boldsymbol{\theta}), \quad (2)$$

with finite and non-zero  $\pi(\mathbf{y})$ . The likelihood function  $\pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$  is defined by the nature of the noise and the noise-free data  $\mathbf{A}(\mathbf{x})$ , which we read as the distribution over  $\mathbf{y}$  conditioned on  $\mathbf{x}$  and  $\boldsymbol{\theta}$ . Here  $\boldsymbol{\theta}$  describe multiple hyper-parameters, e.g. the noise precision so that  $\boldsymbol{\eta} \sim \pi_{\boldsymbol{\eta}}(\cdot | \boldsymbol{\theta})$ , where  $\sim$  reads as "is distributed as". Further,  $\boldsymbol{\theta}$  may account for some physical properties of  $\mathbf{x}$  such as the smoothness (see Sec. 4). Because all unknown parameter are treated as random variables the joint prior distribution is introduced as  $\pi(\mathbf{x}, \boldsymbol{\theta}) = \pi(\mathbf{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})$  with the parameter prior distribution  $\pi(\mathbf{x} | \boldsymbol{\theta})$  and the hyper-prior distribution  $\pi(\boldsymbol{\theta})$ . Choosing these prior distributions is ultimately a modeller's choice and is crucial, as those shall be as uninformative as possible for regions in hyper-parameter and parameter space where the data is informative. If the data is uninformative, the prior distributions can be informative and may represent a rather restrictive range of (physically) feasible hyper-parameters and parameters.

We can write the hierarchical model as:

$$\mathbf{y} | \mathbf{x}, \boldsymbol{\theta} \sim \pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \quad (3a)$$

$$\mathbf{x} | \boldsymbol{\theta} \sim \pi(\mathbf{x} | \boldsymbol{\theta}) \quad (3b)$$

$$\boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}). \quad (3c)$$

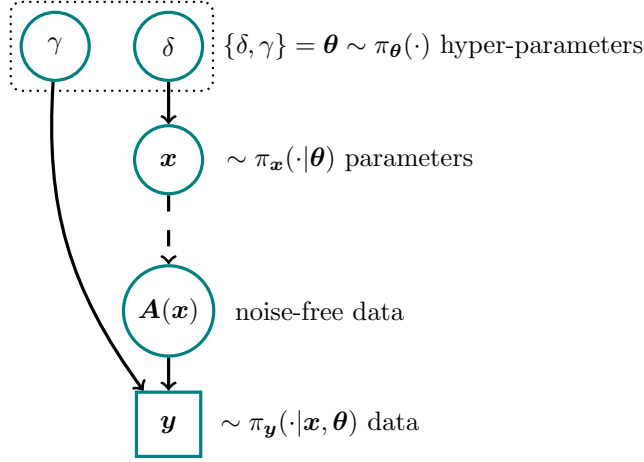


Figure 1: A directed acyclic graph (DAG) for an inverse problem visualises statistical dependencies as solid line arrows and deterministic dependencies as dotted arrows. The hyper-parameters  $\theta$  are distributed as  $(\sim)$  the hyper-prior distribution  $\pi(\theta)$ . The prior distribution  $\pi_x(\cdot|\theta)$  for the parameter  $x$  and the noise  $\eta \sim \pi_\eta(\cdot|\theta)$  are statistically dependent on some of those hyper-parameters. Then a parameter  $x \sim \pi_x(\cdot|\theta)$  is deterministically mapped through the forward model  $A(x)$ . Based on the noise-free data we observe (square box) a data set  $y = A(x) + \eta$  with some additive random noise, which determines the likelihood function  $\pi(y|x, \theta)$ .

Usually, the objective is to calculate the expectation of a function  $h(x)$ , which is defined as

$$\mathbb{E}_{x, \theta|y}[h(x)] = \int \int h(x) \pi(x, \theta|y) dx d\theta. \quad (4)$$

## 2.1 Marginal-then-Conditional Method

Characterising the posterior distribution or quickly generating a representative sample set from the posterior distribution often presents a significant challenge. This is mainly due to the strong correlations that usually exist between the parameters  $x$  and hyper-parameters  $\theta$ , as discussed by Rue and Held in [22].

For models where the full conditional over the latent field has a known form, it is beneficial to factorise the joint posterior distribution

$$\pi(x, \theta|y) = \pi(x|\theta, y) \pi(\theta|y) \quad (5)$$

into the full conditional posterior  $\pi(x|\theta, y)$  over the latent field  $x$  and the marginal posterior  $\pi(\theta|y)$  over hyper-parameter  $\theta$ . This approach, known as the marginal-and-then-conditional (MTC) method [8], is particularly advantageous when  $x \in \mathbb{R}^n$  is high-dimensional, while  $\theta$  is low-dimensional and the evaluation

of the marginal posterior

$$\pi(\boldsymbol{\theta}|\mathbf{y}) = \frac{\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})\pi(\mathbf{y})} \propto \frac{\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \quad (6)$$

as in [8, Lemma 2] is relatively cheap. Applying the law of total expectation [5], Eq. (4) becomes

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}}[h(\mathbf{x})] = \int \int h(\mathbf{x})\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \, d\mathbf{x} \, \pi(\boldsymbol{\theta}|\mathbf{y}) \, d\boldsymbol{\theta} \quad (7)$$

$$= \int \mathbb{E}_{\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}}[h(\mathbf{x})] \, \pi(\boldsymbol{\theta}|\mathbf{y}) \, d\boldsymbol{\theta} \quad (8)$$

$$= \mathbb{E}_{\boldsymbol{\theta}|\mathbf{y}}[\mathbb{E}_{\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}}[h(\mathbf{x})]] . \quad (9)$$

In the case of a linear-Gaussian hierarchical Bayesian model, as in [8], both the marginal distribution  $\pi(\boldsymbol{\theta}|\mathbf{y})$  and the inner expectation  $\mathbb{E}_{\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}}[h(\mathbf{x})]$  are well defined (see Sec. 4 and [8]). If the integral in Eq. 8 is expensive to calculate sample-based methods may be used to calculate the expectations in Eq. (8). To produce samples  $\{(\mathbf{x}, \boldsymbol{\theta})^{(1)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(k)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(N)}\} \sim \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$  one needs an independent sample from  $\boldsymbol{\theta}^{(k)} \sim \pi(\boldsymbol{\theta}|\mathbf{y})$  first and then draws a sample from the full conditional posterior  $\mathbf{x}^{(k)} \sim \pi(\mathbf{x}|\boldsymbol{\theta}^{(k)}, \mathbf{y})$ .

Note that for the affine case, where e.g., the forward model is given by  $\mathbf{Ax} + \mathbf{b}$ , the MTC method as in [8] is still applicable. For Gaussian noise and a Gaussian prior, the form of the posterior of the affine case does not change compared to the linear-Gaussian case, where the forward model may be given by  $\mathbf{Ax}$ .

### 3 The Forward Model

Here we present the forward model to which we apply the methodology. The forward model describes a Limb-sounder measuring thermal radiation of ozone to determine the atmospheric ozone concentration. We follow the MIPAS handbook [18] and simulate data according to an idealised cloud-free atmosphere in local thermodynamic equilibrium, assuming a measurement instrument with infinite spectral resolution and no pointing errors. This is a simplified forward model. No other instrument-specific details such as sensor area or antenna response are included because they are not available to us.

As displayed in Fig. 2, a satellite at a constant height  $h_{\text{sat}}$  is pointing through the atmosphere (limb-sounding) to measure thermal radiation of ozone. For each measurement  $j = 1, 2, \dots, m$ , the tangent height  $h_{\ell_j}$  and the corresponding line-of-sight  $\Gamma_j$  are defined. Additionally, we introduce the pointing angle  $0 \leq \phi_j < \phi_{\text{max}}$ , so that if  $\phi = 0$  arc sec the satellite points at  $h_{L,0}$  and for a pointing angle  $\phi_{\text{max}}$  at  $h_{L,n}$ . Further, the atmosphere is discretised into  $n$  layers defined by height values  $h_{L,i-1} < h_{L,i}$  with respect to the surface of the Earth, for  $i = 1, \dots, n$ . More specifically, the  $i$ -th layer is defined by two spheres around the centre of the Earth with radii  $r_0 + h_{L,i-1}$  and  $r_0 + h_{L,i}$ , where  $r_0$  is the

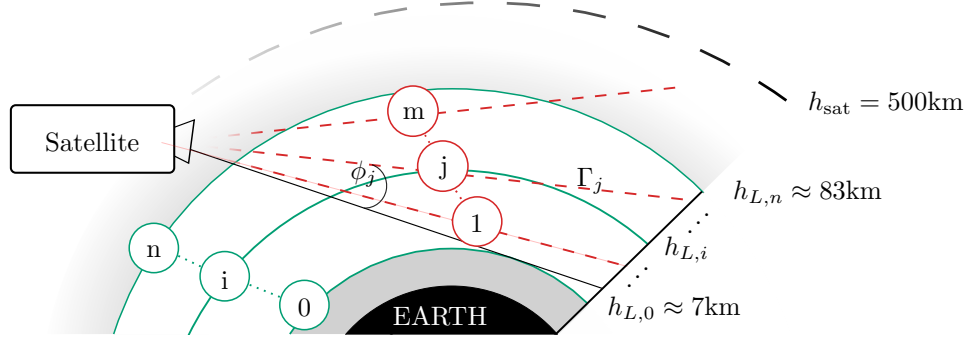


Figure 2: Schematic of measurement and analysis geometry, not to scale. The stationary satellite, at a constant height  $h_{\text{sat}}$  above Earth, takes  $m$  measurements along its line-of-sight defined by the line  $\Gamma_j$ . Each measurement has a pointing angle  $\phi_j$  and a tangent height  $h_{L,j}$ ,  $j = 1, 2, \dots, m$  defined as the closest distance of  $\Gamma_j$  to the Earth's surface. Between  $h_{L,0} \approx 7\text{km}$  and  $h_{L,n} \approx 83\text{km}$ , the atmosphere is discretised into  $n$  layers as illustrated by the solid green lines.

Earth's radius. Within a layer the signal is constant, whereas above  $h_{L,n}$  and below  $h_{L,0}$  no signal can be obtained.

### 3.1 Radiative Transfer Equation

One noise-free measurement of thermal radiation emitted by gas molecules within the atmosphere is described by the radiative transfer equation (RTE) [18]

$$\int_{\Gamma_j} B(\nu, T) k(\nu, T) \frac{p(r)}{k_B T(r)} x(r) \tau(r) dr \quad (10)$$

$$\text{with } \tau(r) = \exp \left\{ - \int_{r_{\text{obs}}}^r k(\nu, T) \frac{p(r')}{k_B T(r')} x(r') dr' \right\}. \quad (11)$$

This is a path integral along the satellite's straight line of sight  $\Gamma_j$  with the ozone volume mixing ratio (VMR)  $x(r)$  at distance  $r$  from the satellite, at the wave number  $\nu$ . Within the atmosphere, the number density  $p(r)/(k_B T(r))$  of molecules is dependent on the pressure  $p(r)$ , the temperature  $T(r)$ , and the Boltzmann constant  $k_B$ . The factor  $\tau(r) \leq 1$  accounts for re-absorption of the radiation along the line-of-sight, which makes the RTE non-linear. The absorption constant is given as

$$k(\nu, T) = L(\nu, T_{\text{ref}}) \frac{Q(T_{\text{ref}})}{Q(T)} \frac{\exp \{-c_2 E''/T\}}{\exp \{-c_2 E''/T_{\text{ref}}\}} \frac{1 - \exp \{-c_2 \nu/T\}}{1 - \exp \{-c_2 \nu/T_{\text{ref}}\}} \quad (12)$$

with Planck's constant  $h$  and speed of light  $c$ . The line intensity  $L(\nu, T_{\text{ref}})$  at reference temperature  $T_{\text{ref}} = 296\text{K}$ , the lower-state energy  $E''$  in  $\text{cm}^{-1}$  of the targeted transition and the second radiation constant  $c_2 := hc/k_B \approx 1.44\text{cmK}$

are provided by the HITRAN database [10]. The total internal partition function is given as

$$Q(T) = g' \exp \left\{ -\frac{c_2 E'}{T} \right\} + g'' \exp \left\{ -\frac{c_2 E''}{T} \right\}, \quad (13)$$

with the statistical weight  $g''$  for the lower and  $g'$  for the upper energy state (also called the degeneracy factors) accounting for the molecule's non-rotational and rotational energy states (see also [25]), and the upper state energy  $E' = E'' + \nu$ . Under the assumption of local thermodynamic equilibrium (LTE), the black body radiation acts as a source function

$$B(\nu, T) = \frac{2hc^2\nu^3}{\exp \left\{ \frac{c_2\nu}{T} \right\} - 1}. \quad (14)$$

For fundamentals on the RTE, we recommend [23, Chapter 1], and for a more comprehensive model, we refer to [17].

When simulating data, we assume an idealised limb-sounder. Since the measurement device has a negligible frequency window, the line broadening around  $\nu$  for the calculations of  $L(\nu, T_{\text{ref}})$  is neglected. Normally, this is modelled as the convolution of the normalised Lorentz profile (collisional/pressure broadening) and the normalised Doppler profile (thermal broadening) [18]. Additionally, we target one specific molecule and calculate  $k(\nu, T)$  accordingly. Usually, this would involve a summation over the individual absorption constants for multiple radiating molecules weighted by their respective VMR [18].

### 3.2 Simulated Data and Ground Truth

As the ground truth for our methodology, we consider an ozone profile at distinct pressure values generated from some data [24] of the MLS on the Aura satellite within the Antarctic region. This ozone profile has a peak in the middle atmosphere and a second peak at higher altitudes, see Fig. 8, which seems to be a typical nighttime profile [13]. For more information on the processes within the atmosphere for ozone, we refer to [13].

We can relate the height  $h$  and the pressure values  $p$  via the hydrostatic equilibrium equation

$$d(\log p) = \frac{dp}{p} = \frac{-gM}{R^*T} dh. \quad (15)$$

Here the acceleration due to gravity is  $g$ , the universal gas constant is  $R^* = 8.31432 \times 10^{-3} \text{Nm/kmol/K}$  and the mean molecular weight of the air is  $M = 28.97 \text{kg/kmol}$ , as in [27]. To enable efficient calculation of the RTE we discretise the atmosphere as in Fig. 2. Then the ozone VMR  $\mathbf{x} = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^n$ , pressure  $\mathbf{p} = \{p_1, p_2, \dots, p_n\} \in \mathbb{R}^n$  and temperature  $\mathbf{T} = \{T_1, T_2, \dots, T_n\} \in \mathbb{R}^n$ , as well as all other height dependent parameters, are discretised profiles with constant values between the heights  $h_{L,i-1} \leq h < h_{L,i}$ , for each layer  $i =$

$1, \dots, n$ . The hydrostatic equilibrium equation for the discretised atmosphere is

$$h_{L,i} = h_{L,i-1} - \frac{\Delta p R^* T_{i-1}}{p_{i-1} g_{i-1} M} \quad (16)$$

with  $\Delta p = p_i - p_{i-1}$  and  $T_{i-1} = T(h_{i-1})$  as in Eq. 18 (see also [4, 20]), for  $i = 1, \dots, n$ . At sea level  $h = 0$ km the mean pressure is  $p_0 = 1013.25$ hPa and the mean temperature is  $T_0 = 288.15$ K [27]. The acceleration due to gravity is

$$g_i = g_0 \left( \frac{r_0}{r_0 + h_{L,i}} \right), \quad (17)$$

where the polar radius of the Earth is  $r_0 \approx 6356$  km, the gravitation at sea level is  $g_0 \approx 9.81$ m/s<sup>2</sup>. For a ground truth temperature profile we follow [27] and form the temperature function

$$T(h) = \begin{cases} T_0 & , h = 0 \\ T_0 + a_0 h & , 0 \leq h < h_{T,1} \\ T_0 + a_0 h_{T,1} & , h_{T,1} \leq h < h_{T,2} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) + a_2(h - h_{T,2}) & , h_{T,2} \leq h < h_{T,3} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h - h_{T,3}) & , h_{T,3} \leq h < h_{T,4} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h_{T,4} - h_{T,3}) + a_4(h - h_{T,4}) & , h_{T,4} \leq h < h_{T,5} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h_{T,4} - h_{T,3}) + a_4(h_{T,5} - h_{T,4}) \\ \quad + a_5(h - h_{T,5}) & , h_{T,5} \leq h < h_{T,6} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h_{T,4} - h_{T,3}) + a_4(h_{T,5} - h_{T,4}) \\ \quad + a_5(h_{T,6} - h_{T,5}) + a_6(h - h_{T,6}) & , h_{T,6} \leq h \lesssim 86 \end{cases} \quad (18)$$

with gradient and height values in Tab. 1 provided by [27]. This function describes the mean temperature in the atmosphere with various height-depending gradients according to the different atmospheric layers, as displayed in Fig. 3. This holds up to a geometric height of 86km, where we ignore a 0.04% non-linear change in  $M$  from 80km to 86km.

We target ozone at a frequency of 235.71GHz, which lies within the region where the MLS observes ozone [15, 29]. The corresponding wave number is  $\nu = 7.86$ cm<sup>-1</sup>. The absorption constant  $k(\nu, T)$  is calculated as in Eq. 11, following the high-resolution transmission (HITRAN) database [10]. The HITRAN database provides the line intensity  $L(\nu, T_{\text{ref}})$  for the isotopologue <sup>16</sup>O<sub>3</sub> with the AFGL Code 666.

To compute a data vector, we define an atmosphere between  $h_{L,0} = 6.9$ km and  $h_{L,n} = 83.3$ km with  $n = 45$  equidistant layers and a satellite fixed at a height of  $h_{\text{sat}} = 500$ km (see Fig. 2). We measure  $m = 30$  times between heights of  $\approx 7$ km and  $\approx 68$ km with pointing accuracy 175arc sec and equidistant spaced

subscript $i$	geometric height $h_{T,i}$ in km	gradient $a_i$
0	0	-6.5
1	11	0
2	20.1	1
3	32.2	2.8
4	47.4	0
5	51.4	-2.8
6	71.8	-2

Table 1: Definition of height depending temperature gradients.

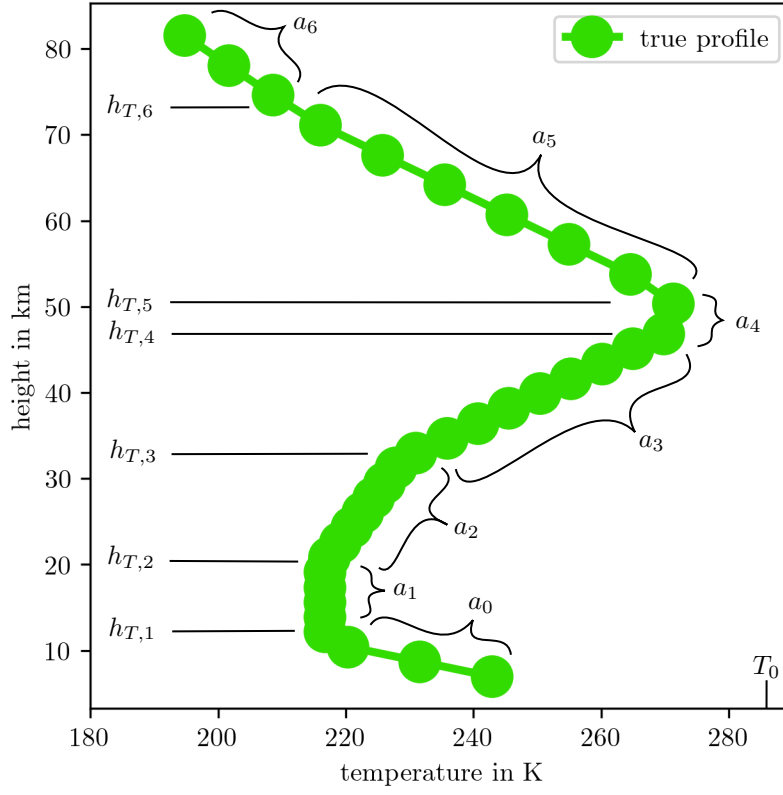


Figure 3: True temperature profile including hyper-parameters as in Eq. 18.

pointing angles

$$\phi_j = (j - 1)175 \text{ arc sec}, \quad \text{for } j = 1, \dots, 30.$$

Above  $\approx 68 \text{ km}$  the data is noise dominated (see Fig. 5), hence no measurement are taken in larger altitudes. Each pointing angle  $\phi_j$  defines a path



$\Gamma_j$  (see Fig. 2). The corresponding path integrals in Eq. 10 and Eq. 11 are evaluated using the trapezoidal rule and define the non-linear forward model  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^m$  for the set of  $m$  noise-free measurements. Here, each entry  $A_j$  of  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^m$  includes multiple evaluations of the integral in Eq. 11 to calculate the absorption  $\tau(r)$ . For brevity we denote the non-linear forward model as  $\mathbf{A}(\mathbf{x}) := \mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T})$ . The simulated data vector

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \boldsymbol{\eta} \quad (19)$$

includes an additive identically-distributed Gaussian noise vector  $\boldsymbol{\eta} \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$ . The noise precision is chosen so that the signal-to-noise ratio (SNR) is approximately 150. The SNR is defined as

$$\text{SNR} := \frac{\max(y)}{\text{STD noise}} = \frac{\text{peak signal}}{\text{RMS noise}}, \quad (20)$$

where STD noise is the standard deviation of the noise. An SNR of 150 is similar to [9], where a signal with a maximal spectral intensity of around 100K and a noise range of 0.4 to 1.6K is reported.

By neglecting the absorption (e.g., set  $\tau = 1$  in Eq. (11)) the RTE is linearised. This denotes the linear forward model matrix  $\mathbf{A}_L \in \mathbb{R}^{m \times n}$ . The integral in Eq. (10) is evaluated using the trapezoidal rule and enables matrix-vector multiplication  $\mathbf{A}_L \mathbf{x}$  to compute noise-free linear data. Since neglecting the absorption changes the measurements only slightly (about 1%, see Sec. 5.3.1), we classify the inverse problem as a weakly non-linear inverse problem. Note that the methods used here will work with different SNRs or other frequencies.

This weakly non-linear forward model may be approximated with an affine map  $\mathbf{M}$ , so that  $\mathbf{A}(\mathbf{x}) \approx \mathbf{M} \mathbf{A}_L \mathbf{x}$ . The affine map  $\mathbf{M} = \mathbf{I} + \epsilon$  captures the linear forward model plus a small change  $\epsilon$ . Note, because of this, we do not necessarily have to find a fixed point to approximate around. In the case of a Taylor series (of first order), one may approximate around the posterior ozone mean, see Eq. 46. The affine map (or the small change  $\epsilon$ ) may be calculated with a linear solver or a least squares fit.

## 4 Linear-Gaussian Hierarchical Model

First this inverse problem is treated as a linear inverse problem by neglecting the absorption term in the RTE (see Eq. 10), so that the linear forward model matrix is defined as  $\mathbf{A} := \mathbf{M} \mathbf{A}_L$  with  $\mathbf{M} = \mathbf{I}$ . The reader is guided through the process of setting up a hierarchical Bayesian framework and establishing a choice of prior distributions. Applying the MTC scheme, we explicitly formulate the respective posterior distributions. A directed acyclic graph (DAG) is used to visualise conditional dependencies between hyper-parameters and the parameter (see Fig. 4), as well as how distributions progress through to an observation (square box). We plot statistical dependencies as solid arrows and deterministic dependencies as dotted arrows.

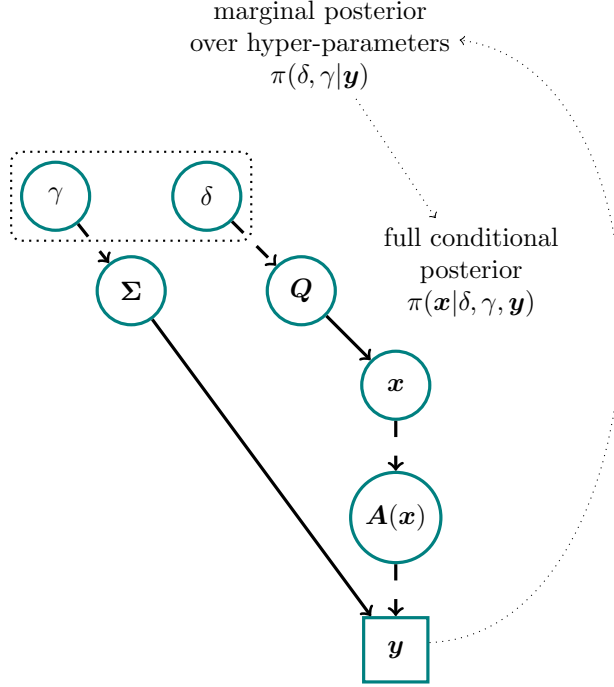


Figure 4: DAG for visualisation of the hierarchical modelling process and the conditional dependency between the parameter and the hyper-parameters. The hyper-parameter  $\gamma$  deterministically (dotted line) sets the noise covariance  $\Sigma = \gamma^{-1}\mathbf{I}$ , which then describes the random (solid line) noise vector  $\boldsymbol{\eta} \sim \mathcal{N}(0, \Sigma)$ . The hyper-parameter  $\delta$  determines (dotted line) the prior precision matrix  $\mathbf{Q} = \delta\mathbf{L}$  for the normally distributed (solid line) prior  $\mathbf{x}|\delta \sim \mathcal{N}(0, \mathbf{Q}^{-1})$ , where  $\mathbf{L}$  is a graph Laplacian, see Eq. 22. The hyper-prior distributions (solid line)  $\pi(\delta, \gamma)$  are defined by  $\boldsymbol{\theta}_\gamma$  and  $\boldsymbol{\theta}_\delta$ . Through the forward model we generate noise-free data  $\mathbf{A}(\mathbf{x})$  and the observed data set  $\mathbf{y}$  includes some added noise  $\boldsymbol{\eta}$ . Within the MTC scheme, we evaluate the marginal posterior over the hyper-parameters  $\pi(\gamma, \delta|\mathbf{y})$  first and then the full conditional posterior  $\pi(\mathbf{x}|\delta, \gamma, \mathbf{y})$ . This breaks the correlation structure of  $\delta$  and  $\gamma$ , and  $\mathbf{x}$ , and allows us to evaluate the marginal posterior independently of  $\mathbf{x}$ .

The distributions that define the hierarchical Bayesian model are:

$$\mathbf{y}|\mathbf{x}, \delta, \gamma \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \gamma^{-1}\mathbf{I}) \quad (21a)$$

$$\mathbf{x}|\delta \sim \mathcal{N}(\mathbf{0}, (\delta\mathbf{L})^{-1}) \quad (21b)$$

$$\delta \sim \mathcal{T}(\alpha_\delta, \beta_\delta) \quad (21c)$$

$$\gamma \sim \mathcal{T}(\alpha_\gamma, \beta_\gamma). \quad (21d)$$

Assuming Gaussian noise  $\boldsymbol{\eta} \sim \mathcal{N}(0, \gamma^{-1}\mathbf{I})$ , the likelihood function is a normal

distribution with mean  $\mathbf{Ax}$  and covariance matrix  $\gamma^{-1}\mathbf{I}$ . We define a normal prior-distribution  $\pi(\mathbf{x}|\delta)$  with zero mean and precision matrix  $\delta\mathbf{L}$ , where  $\delta$  is a smoothness hyper-parameter and  $\mathbf{L}$  is a discrete approximation to the second derivate operator (see Eq. 22). Here the hyper-prior distributions  $\pi(\delta)$  and  $\pi(\gamma)$  are Gamma distributions with shape  $\alpha$  and rate  $\beta$ .

We can visualise this hierarchical structure and the conditional dependencies between hyper-parameters and parameters through a DAG, as in Fig. 4. The hyper-parameter  $\gamma$  sets the noise covariance deterministically (dotted line), but is itself statistically (solid line) defined by the hyper-prior distribution  $\pi(\gamma)$ . This is a Gamma distribution, where  $\theta_\gamma$  determines the shape and rate of  $\pi(\gamma)$ . Similarly  $\theta_\delta$  defines the Gamma distribution  $\pi(\delta)$ . Through the linear forward model the space  $\Omega$  is determined by all measurable noise-free data sets  $\mathbf{Ax}$ . From that space we observe (square box) a data set  $\mathbf{y}$  including some additive noise  $\boldsymbol{\eta}$ .

Given that data, we “reverse the arrows” to determine the posterior distribution  $\pi(\mathbf{x}, \delta, \gamma|\mathbf{y})$  over the parameter  $\mathbf{x}$  and the hyper-parameters  $\delta$  and  $\gamma$ . Usually, due to underlying correlation structures between the parameter and the hyper-parameters, evaluating this posterior poses a significant challenge. The MTC scheme breaks this correlation and provides the marginal posterior  $\pi(\delta, \gamma|\mathbf{y})$  first and then the full conditional posterior  $\pi(\mathbf{x}|\delta, \gamma, \mathbf{y})$ .

#### 4.0.1 Prior Modelling

Completing this Bayesian framework one has to define prior distributions over the hyper-parameters and parameters. Ideally, we define the prior distributions as uninformative as possible, and include functional dependencies and physical properties.

By choosing a normally distributed prior  $\pi(\mathbf{x}|\delta)$  with zero mean and no other restrictions, it is clear that our model does not take into account that ozone values cannot be negative. The precision matrix of that prior distribution is

$$\delta\mathbf{L} = \delta \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \quad (22)$$

which is a discrete approximation to the second derivative operator with Dirichlet boundary condition and defines a 1-dimensional Graph Laplacian as in [28, 8]. We reduce the dimension of  $\mathbf{x}$  from 45 to  $n = 34$  by discarding every second ozone VMR over a height of  $\approx 47\text{km}$ . Doing that, while not changing  $\mathbf{L}$  effectively induces a larger correlation between points at higher altitude.

For  $\delta$  and  $\gamma$  we pick relatively uninformative Gamma distributions so that  $\gamma \sim \mathcal{T}(\theta_\gamma) \propto \gamma^{\alpha_\gamma-1} \exp(-\beta_\gamma\gamma)$  and  $\delta \sim \mathcal{T}(\theta_\delta)$  with  $\theta_\gamma = \{\alpha_\gamma, \beta_\gamma\} = \{\alpha_\delta, \beta_\delta\} = \theta_\delta = (1, 10^{-35})$  similar to [8]. Because of those Gamma distributions,  $\pi(\gamma|\lambda, \mathbf{y}) \sim \mathcal{T}(\cdot)$  is a Gamma distribution with  $\lambda = \delta/\gamma$  and easy to sample from.

## 4.1 Posterior Distribution

As explained in Sec. 2.1, we factorise the posterior

$$\pi(\mathbf{x}, \delta, \gamma | \mathbf{y}) \propto \pi(\mathbf{y} | \mathbf{x}, \delta, \gamma) \pi(\mathbf{x}, \delta, \gamma) \quad (23)$$

into

$$\pi(\mathbf{x}, \delta, \gamma | \mathbf{y}) = \pi(\mathbf{x} | \delta, \gamma, \mathbf{y}) \pi(\delta, \gamma | \mathbf{y}) \quad (24)$$

the marginal posterior  $\pi(\delta, \gamma | \mathbf{y})$  and full conditional posterior  $\pi(\mathbf{x} | \delta, \gamma, \mathbf{y})$  (see Eq. 5). As discussed in [8], for the linear-Gaussian case,  $\mathbf{x}$  cancels in the marginal posterior over the hyper-parameters. Following the MTC scheme, we characterise the marginal posterior first and then the full conditional posterior.

### 4.1.1 Marginal Posterior

For the hierarchical model specified in Eq. 21a to Eq. 21d, the marginal posterior distribution over the hyper-parameters is given by

$$\pi(\lambda, \gamma | \mathbf{y}) \propto \lambda^{n/2 + \alpha_\delta - 1} \gamma^{m/2 + \alpha_\delta + \alpha_\gamma - 1} \exp \left\{ -\frac{1}{2} g(\lambda) - \frac{\gamma}{2} f(\lambda) - \beta_\delta \lambda \gamma - \beta_\gamma \gamma \right\}, \quad (25)$$

with the regularisation parameter  $\lambda = \delta/\gamma$ , and

$$f(\lambda) = \mathbf{y}^T \mathbf{y} - (\mathbf{A}^T \mathbf{y})^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L})^{-1} \mathbf{A}^T \mathbf{y}, \quad (26a)$$

$$\text{and } g(\lambda) = \log \det(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L}). \quad (26b)$$

Note that when changing variables from  $\delta = \lambda \gamma$  to  $\lambda$  the hyper-prior distribution changes to  $\pi(\lambda | \gamma) \propto \lambda^{\alpha_\delta - 1} \gamma^{\alpha_\delta} \exp(-\beta_\delta \lambda \gamma)$ , due to  $d\delta/d\lambda = \gamma$ . Here,  $\lambda$  introduces as the regularisation parameter [8]. Because of the chosen Gamma priors the conditional marginal posterior

$$\gamma | \lambda, \mathbf{y} \sim \mathcal{T} \left( \frac{m}{2} + \alpha_\delta + \alpha_\gamma, \frac{1}{2} f(\lambda) + \beta_\gamma + \beta_\delta \lambda \right) \quad (27)$$

is a Gamma distribution.

### 4.1.2 Full Conditional Posterior

As in [26], consider the joint Gaussian distribution

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} \end{pmatrix}, \begin{pmatrix} \mathbf{Q} + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} & -\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \\ \boldsymbol{\Sigma}^{-1} \mathbf{A} & \boldsymbol{\Sigma}^{-1} \end{pmatrix}^{-1} \right] \quad (28)$$

with  $\boldsymbol{\Sigma}^{-1} = \gamma \mathbf{I}$  and  $\mathbf{Q} = \delta \mathbf{L}$  and  $\boldsymbol{\mu} = \mathbf{0}$ . Then the full conditional posterior distribution of ozone

$$\mathbf{x} | \delta, \gamma, \mathbf{y} \sim \mathcal{N} \left( \underbrace{(\mathbf{A}^T \mathbf{A} + \delta/\gamma \mathbf{L})^{-1} \mathbf{A}^T \mathbf{y}}_{\mathbf{x}_\lambda}, \underbrace{(\gamma \mathbf{A}^T \mathbf{A} + \delta \mathbf{L})^{-1}}_{\gamma \mathbf{B}_\lambda} \right), \quad (29)$$

is a normal distribution with  $\lambda = \delta/\gamma$  and samples can be drawn via the randomise then optimise (RTO) method (see Sec. 5.2).

Alternatively, the posterior mean

$$\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}} = \int \mathbf{x}_\lambda \pi(\lambda|\mathbf{y}) \, \mathrm{d}\lambda \approx \sum \mathbf{x}_{\lambda_i} \pi(\lambda_i|\mathbf{y}), \quad (30)$$

and posterior covariance

$$\boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{y}} = \int \gamma^{-1} \pi(\gamma|\mathbf{y}) \, \mathrm{d}\gamma \int \mathbf{B}_\lambda^{-1} \pi(\lambda|\mathbf{y}) \, \mathrm{d}\lambda \approx \sum \gamma_i^{-1} \pi(\gamma_i|\mathbf{y}) \sum \mathbf{B}_{\lambda_i}^{-1} \pi(\lambda_i|\mathbf{y}) \quad (31)$$

of  $\pi(\mathbf{x}|\mathbf{y})$  can be computed as weighted expectations over the marginal posterior  $\pi(\lambda, \gamma|\mathbf{y})$  by quadrature [6, Sec. 2.1] with  $\sum \pi(\lambda_i|\mathbf{y}) = \sum \pi(\gamma_i|\mathbf{y}) = 1$ .

## 5 Results

Given the simulated data the problem is treated as a linear inverse problem, neglecting the absorption in the RTE. The marginal posterior is evaluated on a predefined grid, where the posterior mean is obtained at no additional computational cost. Samples from the conditional posterior are obtained via the RTO method. An affine map approximates the non-linear forward model from the linear forward model. Again, with the approximated forward model, the marginal posterior is calculated on the same predefined grid. Lastly, the mean and the sample-based STD of the posterior ozone profile are provided. Recall that  $\mathbf{x} \in \mathbb{R}^n$  with  $n = 34$  and  $\mathbf{y} \in \mathbb{R}^m$  with  $m = 30$ .

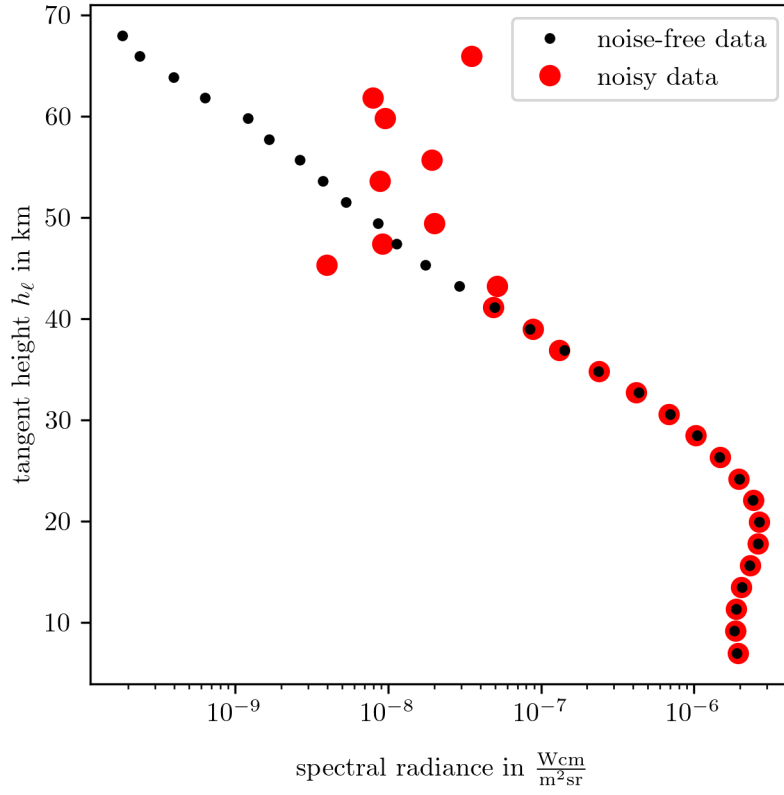


Figure 5: Non-linear noise-free data and non-linear noisy data at different tangent heights. Negative noisy data values are not shown due to logarithmic scaling.

## 5.1 Marginal Posterior

A  $N_{\text{grid}} \times N_{\text{grid}}$  grid for  $\lambda = [10^{-5}, 8 \times 10^{-4}]$  and  $\gamma = [0.8 \times 10^{15}, 1.2 \times 10^{16}]$  with  $N_{\text{grid}} = 25$  grid points is established. The marginal posterior

$$\pi(\lambda, \gamma | \mathbf{y}) \propto \lambda^{34/2} \gamma^{30/2+1} \exp \left\{ -\frac{1}{2} g(\lambda) - \frac{\gamma}{2} f(\lambda) - 10^{-35} \gamma (\lambda - 1) \right\}, \quad (32)$$

is evaluated on that grid. For each evaluation of the marginal posterior, most of the computational effort lies in  $f(\lambda)$  and  $g(\lambda)$ . For each  $\lambda$  on the grid we do a Cholesky decomposition of  $\mathbf{B}_\lambda = \mathbf{A}^T \mathbf{A} + \delta / \gamma \mathbf{L} = \mathbf{C}^T \mathbf{C}$  via the Python function `numpy.linalg.cholesky`. This immediately gives us  $g(\lambda) = 2 \sum \log \text{diag}(\mathbf{C})$ . We calculate  $\mathbf{x}_\lambda = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L})^{-1} \mathbf{A}^T \mathbf{y}$  via `scipy.linalg.cho_solve` which then gives us  $f(\lambda) = \mathbf{y}^T \mathbf{y} - (\mathbf{A}^T \mathbf{y})^T \mathbf{x}_\lambda$ . Then for each  $\lambda$  the cost to compute the marginal posterior  $\pi(\lambda, \cdot | \mathbf{y})$  for varying  $\gamma = [0.8 \times 10^{15}, 1.2 \times 10^{16}]$  is negligible.

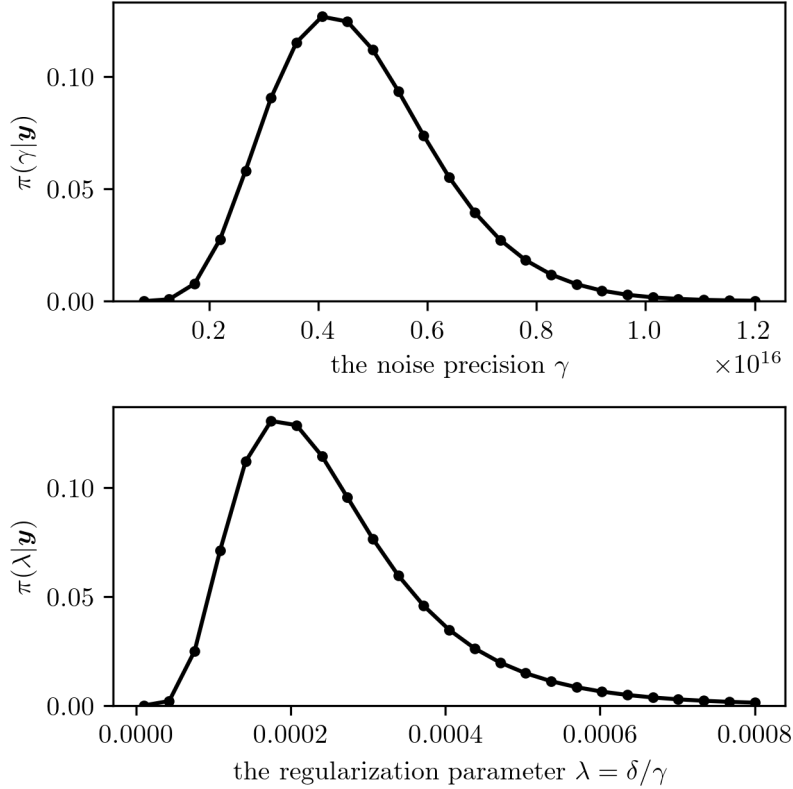


Figure 6: Marginal posterior for linear forward model on a  $20 \times 20$  grid.

Further, the posterior mean in Eq. 30 is obtained by weighted expatiation

over the predefined grid with  $\pi(\lambda|\mathbf{y}) = \int \pi(\lambda, \gamma|\mathbf{y}) d\gamma$  at no additional computational cost. We normalise numerically over the grid so that  $\sum \pi(\lambda_i|\mathbf{y}) = 1$ .

### 5.1.1 Sampling from Marginal posterior

Independent samples  $\lambda^{(k)}, \gamma^{(k)} \sim \pi_{\lambda, \gamma}(\cdot|\mathbf{y})$  from the Marginal posterior are needed if one wants to draw an ozone sample from the full conditional posterior  $\mathbf{x}^{(k)} \sim \pi_{\mathbf{x}}(\cdot|\lambda^{(k)}, \gamma^{(k)}, \mathbf{y})$ .

We sample from the marginal posterior via the inverse Rosenblatt transport (IRT) scheme as in [7]. A seed  $\mathbf{u}^{(k)} \sim \mathcal{U}(0, 1)^d$ , where  $d = 2$  and  $\mathbf{u}^{(k)} = \{u_1^{(k)}, u_2^{(k)}\}$ , from the uniform distribution is projected onto the hyper-parameter space via the cumulative distribution function (CDF)

$$F(\lambda) = \int_{-\infty}^{\lambda} \pi(\lambda'|\mathbf{y}) d\lambda', \quad (33)$$

where  $\pi(\lambda|\mathbf{y}) = \sum_{\gamma} \pi(\lambda, \gamma|\mathbf{y})$  on the predefined grid. The inverse

$$\lambda^{(k)}|\mathbf{y} = F^{-1}(u_1^{(k)}) \quad (34)$$

provides a sample from  $\pi(\lambda|\mathbf{y})$ . In order to draw an samples  $\gamma^{(k)}|\lambda^{(k)}, \mathbf{y}$  we repeat this procedure.  $\pi(\gamma|\lambda^{(k)}, \mathbf{y})$  (see Eq. 32) is calculated for the every  $\gamma$  value on the grid with fixed  $\lambda^{(k)}$ . Here we use piecewise linear interpolation via the Python function `numpy.linalg.interp` to obtain the values  $f(\lambda^{(k)})$  and  $g(\lambda^{(k)})$  (see Fig. 7). Then the CDF

$$F(\gamma) = \int_{-\infty}^{\gamma} \pi(\gamma|\lambda^{(k)}, \mathbf{y}) d\gamma' \quad (35)$$

is calculated and an independent sample is given as:

$$\gamma^{(k)}|\lambda^{(k)}, \mathbf{y} = F^{-1}(u_2^{(k)}). \quad (36)$$

Note that alternatively one could draw an independent sample using Eq. 27, which is a Gamma distribution due to the choice of prior. The presented IRT works for arbitrary priors.

In theory, this scheme provides independent samples, but in practice, every system produces correlated samples. To assess how efficient this scheme is, we define the Integrated Autocorrelation Times (IACT)

$$\tau_{\text{int}} := \left( 1 + 2 \sum_{t=1}^W \frac{\Gamma(t)}{\Gamma(0)} \right) \quad (37)$$

as in [8] with an autocorrelation coefficient  $\Gamma(t)$ . This is twice the value of the IACT in [31, pp. 103-105] and [30, 11], as commonly defined within the physics community. U. Wolff [30] (and the Python implementation by D. Hesse [11]) provide a way to not only calculate the IACT safely but also to quantify the errors of



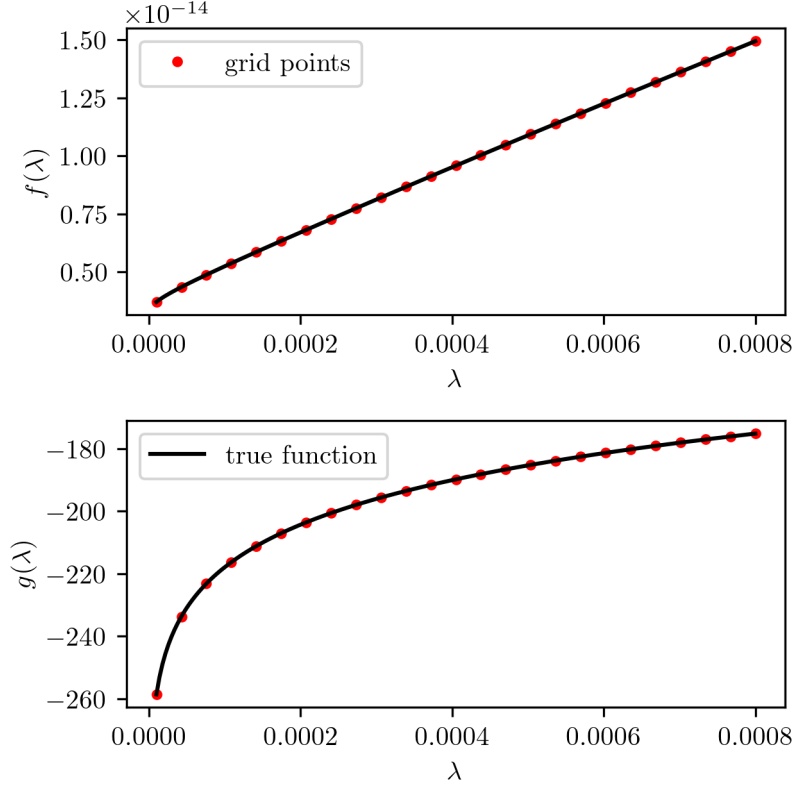


Figure 7: Functions  $f(\lambda)$  and  $g(\lambda)$  for  $\lambda = [10^{-5}, 8 \times 10^{-4}]$ , grid points are marked red. True function values between grid points are plotted in black.

the estimated IACT. The IATCs are  $\tau_{\text{int},\lambda} = 1.05 \pm 0.04$  and  $\tau_{\text{int},\gamma} = 0.99 \pm 0.03$  based on a chain with 10000 samples. So every second sample via the IRT scheme presents an independent marginal posterior sample. Conditioned on an independent sample from the marginal posterior, one can draw an independent sample from the full conditional posterior.

## 5.2 Full Conditional Posterior

To sample an posterior ozone profile from the full conditional  $\pi(\mathbf{x}|\delta, \gamma, \mathbf{y})$  with  $\delta = \lambda\gamma$  we use the RTO method, as in [1, 2, 3, 8]. Rewrite the full conditional

posterior as

$$\pi(\mathbf{x}|\delta, \gamma, \mathbf{y}) \propto \pi(\mathbf{y}|\mathbf{x}, \gamma)\pi(\mathbf{x}|\delta) \quad (38)$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{A}\mathbf{x} - \mathbf{y})^T \boldsymbol{\Sigma}^{-1}(\mathbf{A}\mathbf{x} - \mathbf{y})\right) \exp\left(-\frac{1}{2}(\boldsymbol{\mu} - \mathbf{x})^T \mathbf{Q}(\boldsymbol{\mu} - \mathbf{x})\right), \quad (39)$$

$$= \exp\left(-\frac{1}{2}\left\|\hat{\mathbf{A}}\mathbf{x} - \hat{\mathbf{y}}\right\|_{L^2}^2\right), \quad (40)$$

where

$$\hat{\mathbf{A}} := \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2} \mathbf{A} \\ \mathbf{Q}^{1/2} \end{bmatrix}, \quad \hat{\mathbf{y}} := \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2} \mathbf{y} \\ \mathbf{Q}^{1/2} \boldsymbol{\mu} \end{bmatrix}, \quad (41)$$

$\mathbf{Q} = \delta \mathbf{L}$  is the prior precision,  $\boldsymbol{\mu} = \mathbf{0}$  the prior mean and  $\boldsymbol{\Sigma} = \gamma^{-1} \mathbf{I}$  the noise covariance. A sample  $\mathbf{x}^{(k)}$  from the full conditional posterior  $\pi(\mathbf{x}|\delta, \gamma, \mathbf{y})$  is obtained by minimising the following equation:

$$\mathbf{x}^{(k)} = \arg \min_{\mathbf{x}} \|\hat{\mathbf{A}}\mathbf{x} - (\hat{\mathbf{y}} + \mathbf{b})\|_{L^2}^2 \quad (42)$$

with a random additive perturbation  $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . This expression becomes

$$(\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} + \mathbf{Q}) \mathbf{x}^{(k)} = \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} + \mathbf{Q} \boldsymbol{\mu} + \mathbf{v}_1 + \mathbf{v}_2, \quad (43)$$

with  $\mathbf{v}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A})$  and  $\mathbf{v}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1})$ , representing independent Gaussian random variables [1, 8].

More explicitly, conditioned on an independent  $\lambda^{(k)}, \gamma^{(k)} \sim \pi_{\lambda, \gamma}(\cdot | \mathbf{y})$ , one independent full conditional posterior sample is given as

$$\mathbf{x}^{(k)} = \left( \underbrace{\gamma^{(k)} \mathbf{A}^T \mathbf{A} + \delta^{(k)} \mathbf{L}}_{\mathbf{B}^{(k)}} \right)^{-1} \left( \gamma^{(k)} \mathbf{A}^T \mathbf{y} + \sqrt{\gamma^{(k)}} \mathbf{A}^T \tilde{\mathbf{v}}_1 + \sqrt{\delta^{(k)}} \mathbf{L}^{1/2} \tilde{\mathbf{v}}_2 \right) \quad (44)$$

with  $\tilde{\mathbf{v}}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\tilde{\mathbf{v}}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\mathbf{L}^{1/2}$  is the Cholesky decomposition of  $\mathbf{L}$  [1]. Note that  $\mathbf{v}_1 \in \mathbb{R}^m$  and  $\mathbf{v}_2 \in \mathbb{R}^n$ . The Cholesky factorisation of  $\mathbf{B}^{(k)}$  and  $\mathbf{L}$  is obtained via the Python function `numpy.linalg.cholesky` and `scipy.linalg.cho_solve` is used to solve for  $\mathbf{x}^{(k)}$ . We draw  $n-1 = 29$  posterior samples in  $\approx 0.003$ . The full conditional posterior samples are plotted in Fig. 8 with negative ozone values set to zero. The fact that we have to deal with negative ozone values is due to the poor prior choice in  $\pi(\mathbf{x}|\delta)$ . Those full conditional ozone samples are used to find an affine map.

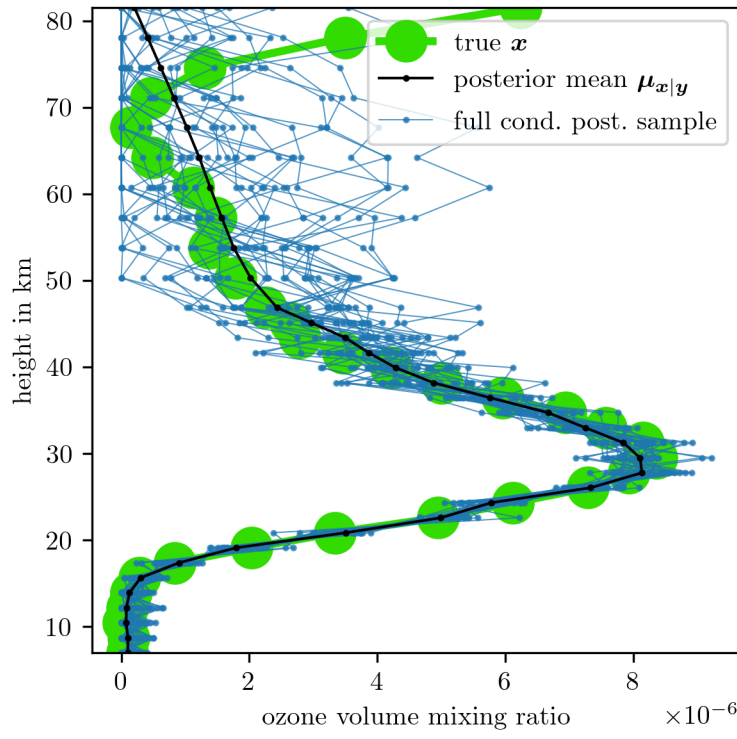


Figure 8: 29 Ozone samples via RTO method and posterior ozone mean via quadrature (see Eq. 30).

### 5.3 Affine Approximation of the Non-Linear Model

The forward map poses a weakly non-linear forward problem. One could tackle this non-linear inverse problem by fixing the absorption at a previously obtained parameter state and treating this as a linear inverse problem. After each parameter sample the absorption is then iteratively updated. Instead, as in Fig. 9 illustrated, we approximate the non-linear model using an affine map, which is a linear map with a translation, e.g.,  $\mathbf{A}\mathbf{x} + \mathbf{b}$ . An affine map  $\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b}$  maps a Gaussian  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  onto a Gaussian  $\mathbf{z} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}^T \boldsymbol{\Sigma} \mathbf{A})$ .

Here we find an affine map  $\mathbf{M}$  that provides a mapping from the linear forward model  $\mathbf{A}_L$  to the non-linear model  $\mathbf{A}(\mathbf{x})$  for parameters  $\mathbf{x}$  near the posterior mean  $\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}$ .

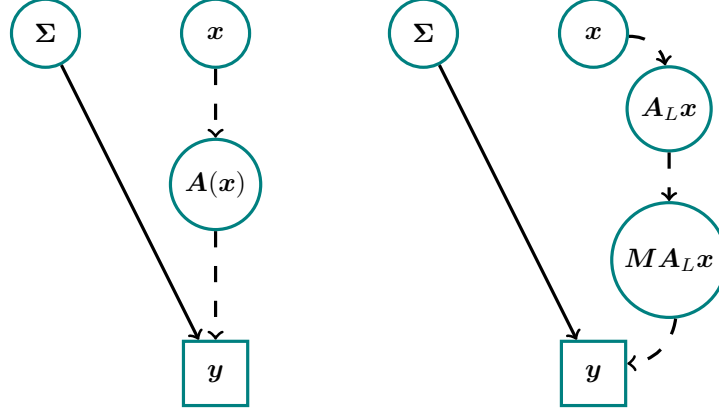


Figure 9: Schematic of how the affine map and the linear forward model approximates the non-linear forward model.

#### 5.3.1 Finding an Affine Map

We find an affine map by creating the vector spaces  $\mathbf{W}$  based on the linear forward model and  $\mathbf{V}$  based on the non-linear forward model with ground truth pressure and temperature. More specifically  $m-1$  samples  $\mathbf{x}^{(j)} \sim \pi_{\mathbf{x}}(\cdot | \delta^{(j)}, \gamma^{(j)}, \mathbf{y})$ , for  $j = 2, \dots, m$ , from the posterior and the posterior mean  $\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}$  generate,

$$\mathbf{W} = \begin{bmatrix} | & | & & | & & | \\ \mathbf{A}_L \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}} & \mathbf{A}_L \mathbf{x}^{(2)} & \dots & \mathbf{A}_L \mathbf{x}^{(j)} & \dots & \mathbf{A}_L \mathbf{x}^{(m)} \\ | & | & & | & & | \end{bmatrix} \in \mathbb{R}^{m \times m}$$

and

$$\mathbf{V} = \begin{bmatrix} \mathbf{A}(\mu_{\mathbf{x}|\mathbf{y}}) & \mathbf{A}(\mathbf{x}^{(2)}) & \dots & \mathbf{A}(\mathbf{x}^{(j)}) & \dots & \mathbf{A}(\mathbf{x}^{(m)}) \end{bmatrix} = \begin{bmatrix} - & v_1 & - \\ & \vdots & \\ - & v_j & - \\ & \vdots & \\ - & v_m & - \end{bmatrix} \in \mathbb{R}^{m \times m}.$$

Then the non-linear forward model is approximated as

$$\mathbf{A}(\mathbf{x}) \approx \mathbf{M}\mathbf{A}_L\mathbf{x}, \quad (45)$$

where we solve  $v_j = r_j \mathbf{W}$  for each row  $r_j$  in

$$\mathbf{V}\mathbf{W}^{-1} = \mathbf{M} = \begin{bmatrix} - & r_1 & - \\ & \vdots & \\ - & r_j & - \\ & \vdots & \\ - & r_m & - \end{bmatrix} \in \mathbb{R}^{m \times m}.$$

using the Python function `numpy.linalg.solve`. This is feasible since every noise-free measurement is independent of each other, and then every row  $v_j$  of  $\mathbf{V} \in \mathbb{R}^{m \times m}$  is independent of each other as well. For an  $\mathbf{x} = \mu_{\mathbf{x}|\mathbf{y}} + \Delta\mathbf{x}$  we rewrite Eq. 45 to

$$\mathbf{A}(\mathbf{x}) \approx \underbrace{\mathbf{M}\mathbf{A}_L\mu_{\mathbf{x}|\mathbf{y}}}_{=\mathbf{A}(\mu_{\mathbf{x}|\mathbf{y}})} + \underbrace{\mathbf{M}\mathbf{A}_L\Delta\mathbf{x}}_{=\mathbf{A}'(\mu_{\mathbf{x}|\mathbf{y}})\Delta\mathbf{x}} \quad (46)$$

$$= \underbrace{\mathbf{A}'(\mu_{\mathbf{x}|\mathbf{y}})\mathbf{x}}_{\mathbf{A}\mathbf{x}} + \underbrace{\mathbf{A}(\mu_{\mathbf{x}|\mathbf{y}}) - \mathbf{A}'(\mu_{\mathbf{x}|\mathbf{y}})\mu_{\mathbf{x}|\mathbf{y}}}_{\mathbf{b}} \quad (47)$$

to show that  $\mathbf{M} : \mathbf{A}_L\mathbf{x} \rightarrow \mathbf{A}(\mathbf{x})$  is an affine map.

The relative RMS difference  $\|\text{vec}(\mathbf{M}\mathbf{W}) - \text{vec}(\mathbf{V})\|_{L^2} / \|\text{vec}(\mathbf{M}\mathbf{W})\|_{L^2}$  between the mapped linear noise-free data and the non-linear noise-free data is approximately 0.001%. This is much smaller than the relative RMS difference between  $\mathbf{W}$  and  $\mathbf{V}$  of about 1%. Here  $\text{vec}(\mathbf{V})$  vectorises the matrix  $\mathbf{V}$ . Fig. 10 shows the mapping for one posterior ozone sample with a relative RMS error  $\approx 0.01\%$ . This posterior ozone sample has not been used to create this mapping; in other words, this is an unseen event not occurring in the training data. Consequently, from here onwards the approximated forward map is used. This takes  $\approx 0.1\text{s}$ .

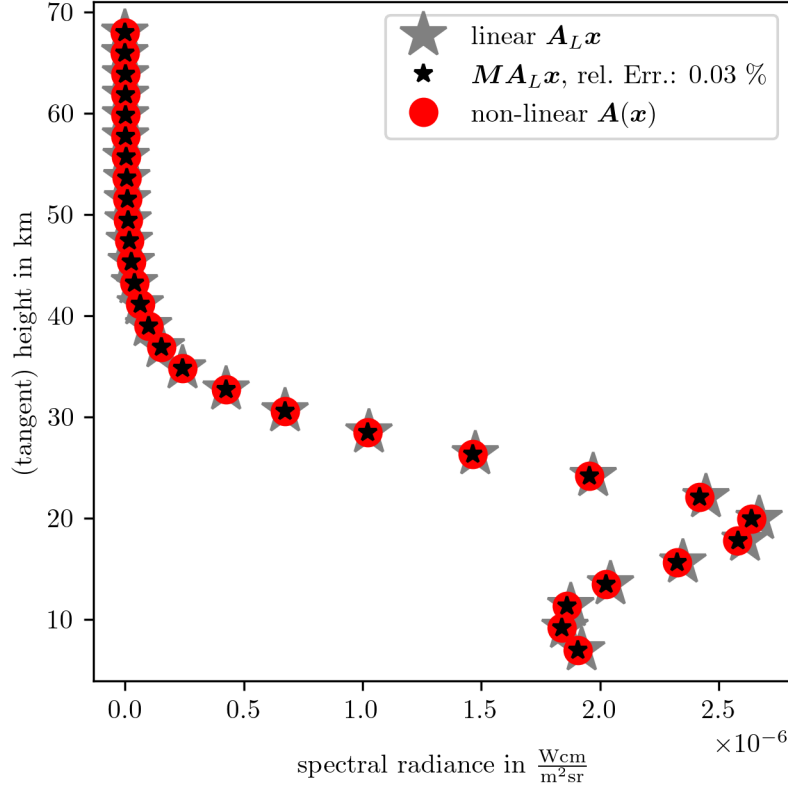


Figure 10: Assessment of how well the affine map  $\mathbf{M}$  approximates noise-free non-linear data  $\mathbf{A}(\mathbf{x})$  (red circles) from noise-free linear data  $\mathbf{A}_L \mathbf{x}$  (grey stars). The approximated noise-free data (black stars) has a relative RMS error of  $\approx 0.01\%$  compared to the true non-linear noise-free data. The ozone profile  $\mathbf{x}$  to generate this noise-free data has not been used to create the affine map.

## 5.4 Posterior ozone

The linear forward model that approximates the non-linear forward model is defined as

$$\mathbf{A} := \mathbf{M}\mathbf{A}_L. \quad (48)$$

We use the exact same grid and setup as in Sec. 5.1 to define the marginal posterior (see Fig. 11). As already mentioned, that gives us the posterior mean

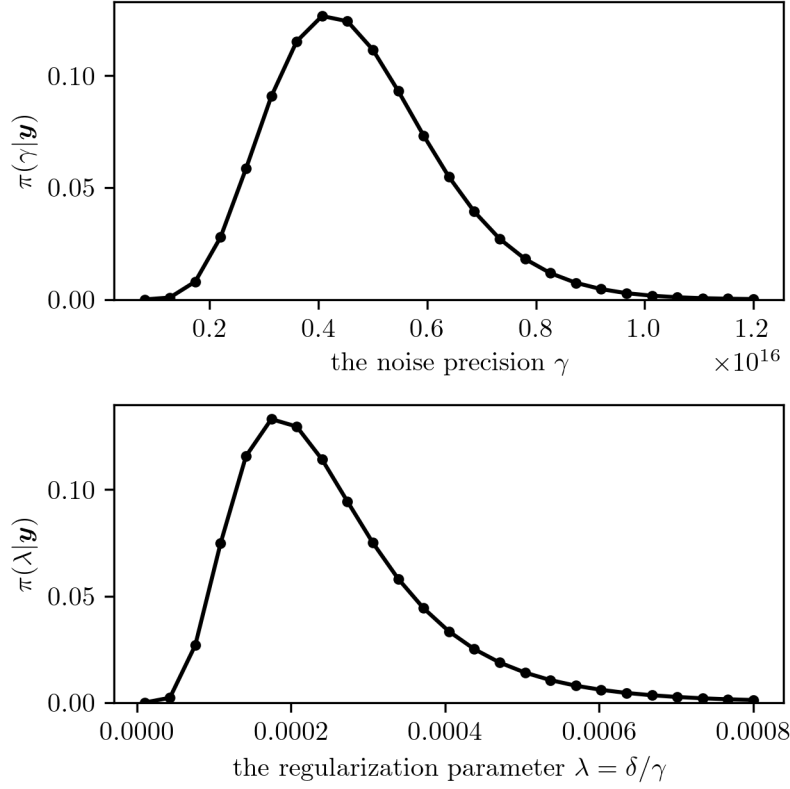


Figure 11: Marginal posterior for approximated forward model on a  $20 \times 20$  grid.

via weighted expectation for free. The IATCs of the marginal posterior samples are  $\tau_{\text{int},\lambda} = 0.97 \pm 0.03$  and  $\tau_{\text{int},\gamma} = 0.95 \pm 0.03$  based on a chain with length of 10000. Hence, every second sample is an independent sample from the marginal posterior.

Samples from the conditional posterior are drawn via the RTO method. To estimate the number of samples needed for a good enough estimate of the STD,

we calculate the average Coefficient of Variation ( $\overline{CV}$ ),

$$\overline{CV} = \frac{1}{n} \left\| \frac{\sqrt{\text{var}_N(\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}})}}{\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}} \right\|_1, \quad (49)$$

recall that  $\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}} \in \mathbb{R}^n$  with  $n = 34$ . The  $\overline{CV}$  gives a measure of how much the STD varies compared to the mean.

If the STD is large compared to the mean the  $\overline{CV}$  is large, if the STD is small compared to the mean the  $\overline{CV}$  is small and a rather accurate estimate of the STD is needed. The sample-based estimate of the variance

$$\text{var}_N(\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}) = \frac{1}{N-1} \sum_{k=1}^N (\mathbf{x}^{(k)} - \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}})^T (\mathbf{x}^{(k)} - \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}) \quad (50)$$

is calculated with the posterior mean  $\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}$  given by weighted expectations as in Eq. 30. Alternatively, one could evaluate the posterior covariance as a weighted expectation over the marginal posterior grid (see Eq. 31).

The  $\overline{CV}$  is calculated for each number of samples ranging from 1 to  $10^3$  and plotted in Fig. 12. Fig. 12 shows that the  $\overline{CV}$  is rather large, indicating that

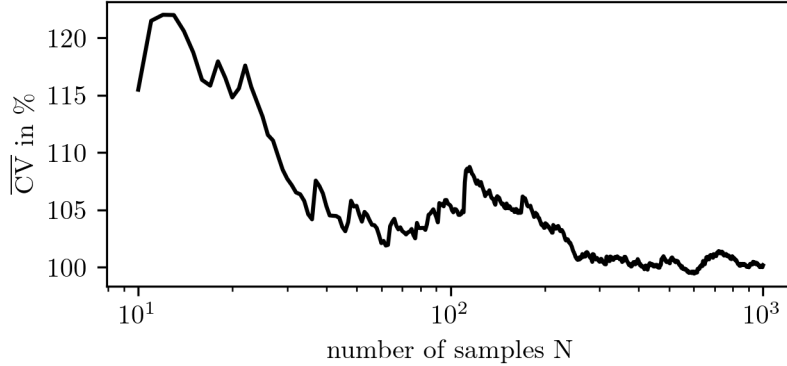


Figure 12: Average Coefficient of Variation in % as calculated in Eq. 49.

the STD is large compared to the posterior mean. Further, we conclude that 50 full conditional posterior ozone samples are enough for a sufficient estimate of the STD. Approximating the marginal posterior on a  $20 \times 20$  grid and taking 50 independent full posterior samples of ozone takes  $\approx 0.008$ s. The sample-based STD with the posterior mean (via weighted expectation) is plotted in Fig. 13. Note that the posterior ozone profile does not capture the second ozone peak at around 80km.



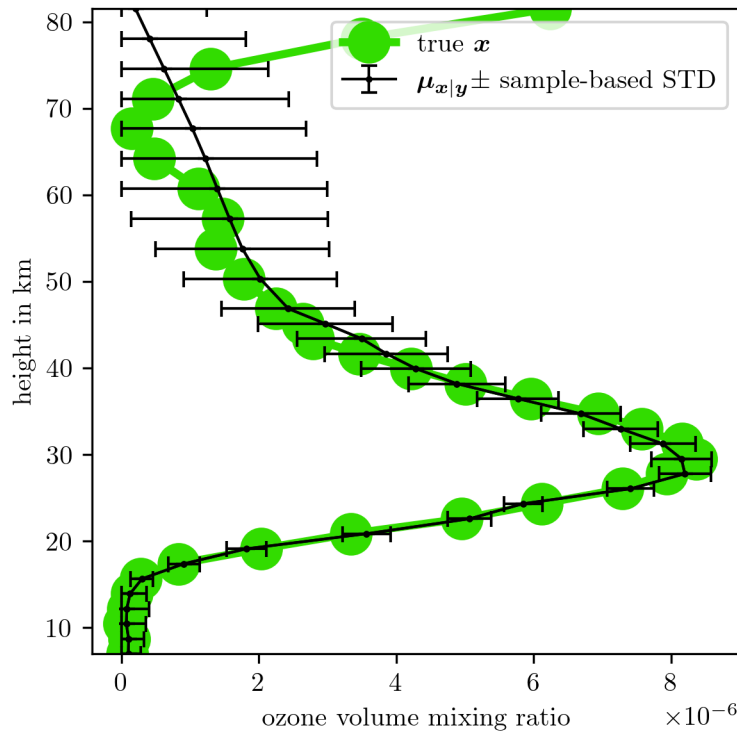


Figure 13: Posterior ozone mean via quadrature (see Eq. 30) and STD based on 20 samples via the RTO method.

## 6 Conclusion

- 20 solves of  $x_\lambda$  for marginal posterior and full posterior mean (see Eq. 29 and Eq. 30)
- 29 samples of full conditional posterior are enough for finding an affine map, this includes 29 solves of “randomised”  $x_\lambda$  via the RTO method (see Eq. 44)
- 50 samples of full conditional posterior are enough for the STD estimate of posterior ozone, this includes 50 solves of “randomised”  $x_\lambda$  via the RTO method (see Eq. 44)

## References

- [1] Bardsley, J. “MCMC-based image reconstruction with uncertainty quantification”. In: *SIAM Journal on Scientific Computing* 34.3 (2012), A1316–A1332.
- [2] Bardsley, J., Solonen, A., Haario, H., and Laine, M. “Randomize-then-optimize: A method for sampling from posterior distributions in nonlinear inverse problems”. In: *SIAM Journal on Scientific Computing* 36.4 (2014), A1895–A1910.
- [3] Bardsley, J. and Cui, T. “A Metropolis-Hastings-Within-Gibbs sampler for nonlinear hierarchical-Bayesian inverse problems”. In: *2017 MATRIX Annals*. Vol. 2. MATRIX Book Series. Switzerland: Springer, 2019, pp. 2–12.
- [4] Carlotti, M. and Ridolfi, M. “Derivation of temperature and pressure from submillimetric limb observations”. In: *Applied Optics* 38.12 (Apr. 1999), pp. 2398–2409.
- [5] Champ, C. W. and Sills, A. V. “The generalized law of total covariance”. In: *preprint* (2022).
- [6] Dick, J., Kuo, F. Y., and Sloan, I. H. “High-dimensional integration: The quasi-Monte Carlo way”. In: *Acta Numerica* 22 (2013), pp. 133–288.
- [7] Dolgov, S., Anaya-Izquierdo, K., Fox, C., and Scheichl, R. “Approximation and sampling of multivariate probability distributions in the tensor train decomposition”. In: *Statistics and Computing* 30 (2020), pp. 603–625.
- [8] Fox, C. and Norton, R. A. “Fast sampling in a linear-Gaussian inverse problem”. In: *SIAM/ASA Journal on Uncertainty Quantification* 4.1 (2016), pp. 1191–1218.
- [9] Froidevaux, L. et al. “Validation of Aura Microwave Limb Sounder stratospheric ozone measurements”. In: *Journal of Geophysical Research: Atmospheres* 113.D15S20 (2008).
- [10] Gordon, I. E et al. “The HITRAN2020 molecular spectroscopic database”. In: *Journal of Quantitative Spectroscopy and Radiative Transfer* 277.107949 (2022).
- [11] Hesse, D. *py-uwerr; Python implementation of Monte Carlo error analysis a la Wolff*. <https://github.com/dhesse/py-uwerr>. [Online; accessed 09/09/25].
- [12] Kaipio, J. P. and Somersalo, E. *Statistical and Computational Inverse Problems*. New York: Springer-Verlag New York, 2005.
- [13] Lee, J. N. and Wu, D. L. “Solar cycle modulation of nighttime ozone near the mesopause as observed by MLS”. In: *Earth and Space Science* 7.4 (2020).

- [14] Livesey, N. J., Van Snyder, W., Read, W. G., and Wagner, P. A. “Retrieval algorithms for the EOS Microwave limb sounder (MLS)”. In: *IEEE Transactions on Geoscience and Remote Sensing* 44.5 (2006), pp. 1144–1155.
- [15] Livesey, N. J. et al. “Validation of Aura Microwave Limb Sounder O<sub>3</sub> and CO observations in the upper troposphere and lower stratosphere”. In: *Journal of Geophysical Research: Atmospheres* 113.D15S02 (2008).
- [16] Raspollini, P. et al. “Level 2 processor and auxiliary data for ESA Version 8 final full mission analysis of MIPAS measurements on ENVISAT”. In: *Atmospheric Measurement Techniques Discussions* 2021 (2021), pp. 1–46.
- [17] Read, W., Shippony, Z., Schwartz, M., Livesey, N. J., and Van Snyder, W. “The clear-sky unpolarized forward model for the EOS aura microwave limb sounder (MLS)”. In: *IEEE Transactions on Geoscience and Remote Sensing* 44.5 (2006), pp. 1367–1379.
- [18] Readings, C. *Envisat MIPAS An Instrument for Atmospheric Chemistry and Climate Research*. Noordwijk: ESA Publications Division, 2000.
- [19] Readings, C. and Harris, R. A. *Envisat MIPAS an Instrument for Atmospheric Chemistry and Climate Research*. <https://earth.esa.int/eogateway/documents/20142/37627/envisat-mipas-instrument-description.pdf>. [Online; accessed 16/07/22]. 2000.
- [20] Ridolfi, M. et al. “Optimized forward model and retrieval scheme for MIPAS near-real-time dataprocessing”. In: *Applied Optics* 39.8 (2000), pp. 1323–1340.
- [21] Rodgers, C. D. “Retrieval of atmospheric temperature and composition from remote measurements of thermal radiation”. In: *Reviews of Geophysics* 14.4 (1976), pp. 609–624.
- [22] Rue, H. and Held, L. *Gaussian Markov Random Fields: Theory and Applications*. London: CRC Press, 2005.
- [23] Rybicki, G. B. and Lightman, A. P. *Radiative Processes in Astrophysics*. Weinheim: Wiley-VCH, 2004.
- [24] Schwartz, M., Froidevaux, L., Livesey, N., and Read, W. *MLS/Aura Level 2 Ozone (O<sub>3</sub>) Mixing Ratio V005*. [https://disc.gsfc.nasa.gov/datasets/ML203\\_005/summary?keywords=mls\%20o3](https://disc.gsfc.nasa.gov/datasets/ML203_005/summary?keywords=mls\%20o3). [Online; accessed 25/04/24]. NASA Goddard Earth Sciences Data and Information Services Center, 2020.
- [25] Šimečková, M., Jacquemart, D., Rothman, L. S., Gamache, R. R., and Goldman, A. “Einstein A-coefficients and statistical weights for molecular absorption transitions in the HITRAN database”. In: *Journal of Quantitative Spectroscopy and Radiative Transfer* 98.1 (2006), pp. 130–155.
- [26] Simpson, D., Lindgren, F., and Rue, H. “Think continuous: Markovian Gaussian models in spatial statistics”. In: *Spatial Statistics* 1 (2012), pp. 16–29.

- [27] *U.S. Standard Atmosphere, 1976*. Washington, D.C.: United States. National Oceanic and Atmospheric Administration, United States Committee on Extension to the Standard Atmosphere, 1976.
- [28] Wang, Y.-X., Sharpnack, J., Smola, A. J., and Tibshirani, R. J. “Trend filtering on graphs”. In: *Journal of Machine Learning Research* 17.105 (2016), pp. 1–41.
- [29] Waters, J. et al. “The earth observing system microwave limb sounder (EOS MLS) on the Aura satellite”. In: *IEEE Transactions on Geoscience and Remote Sensing* 44.5 (2006), pp. 1075–1092.
- [30] Wolff, U. “Monte Carlo errors with less errors”. In: *Computer Physics Communications* 156.2 (2004), pp. 143–153.
- [31] Wolff, U., Bunk, B. hard, Korzec, T., Knechtli, F., and Bär, O. *Lecture Notes on Computational Physics II [in german]*. <https://www.physik.hu-berlin.de/de/com/teachingandseminars/previousCPII>. [Online; accessed 29/08/25]. Humboldt University, Berlin, 2016.