# Bayesian inference of 3D densities of galactic HI and H2

**Laurin Söding,**[a,*] **Philipp Mertsch**[a] **and Vo Hong Minh Phan**[a]

[a]*Institute for Theoretical Particle Physics and Cosmology (TTK), RWTH Aachen University, 52056 Aachen, Germany*

*E-mail:* soeding@physik.rwth-aachen.de, pmertsch@physik.rwth-aachen.de, vhmphan@physik.rwth-aachen.de

Due to our vantage point in the disk of the Galaxy, its 3D structure is not directly accessible. However, knowing the spatial distribution, e.g. of atomic and molecular hydrogen gas is of great importance for interpreting and modelling cosmic ray data and diffuse emission. Using novel Bayesian inference techniques, we reconstruct the 3D densities of atomic and molecular hydrogen in the Galaxy together with (part of) the galactic velocity field. In order to regularise the infinite number of degrees of freedom and obtain information in regions with missing or insufficient data, we incorporate the correlation structure of the gas fields into our prior. Basis for these reconstructions are the data-sets from the HI4PI-survey on the 21-cm emission line and the CO-survey compilation by Dame et al. (2001) on the $(1 \rightarrow 0)$ rotational transition together with a variable gas flow model. We present the preliminary estimated mean surface mass densities and corrections to the prior assumption of the galactic velocity field. In the future, we plan to relax assumptions on the optical thickness and include additional data to further constrain either the galactic velocity field or the gas densities.

*Speaker

PoS(ICRC2023)658

## 1. Introduction

In order to properly interpret measurements of cosmic rays and gamma-ray diffuse emission, it is necessary to understand the emission, propagation and absorption of radiation in the interstellar medium of the Milky Way. This medium is a complex system consisting mainly of gas, magnetic fields, interstellar radiation fields and cosmic rays which permeate the entire Galaxy. While it makes up for only a few percent of the total mass of the Galaxy (the majority is in the form of stars or dark matter), it fills out most of the available volume and thereby defines the dynamics of radiation and particles within.

Due to our vantage point, the 3D-distribution of the constituents of the Galaxy is not easily determined by observations of the sky. No matter in which direction we point our telescopes, we always observe an integrated signal of radiation that has travelled an a priori unknown distance through the Galaxy. However, due to galactic rotation and peculiar motion, light that reaches us will be Doppler-shifted by a certain amount, depending on the relative velocity of its source with respect to us. This is particularly useful when looking at emission lines that have a narrow width as it enables us to determine the relative velocity of its emitter and observer very precisely. This idea is unfortunately somewhat tainted by the fact that we do not know the precise structure of the galactic velocity field - and even circular rotation features a velocity ambiguity for positions within the solar circle. Any attempt to produce 3D maps of some quantity from such data will thus have to specify some rule according to which said quantity is placed when there is an ambiguity. Multiple approaches have been tried, most of them treating every line-of-sight (direction) independently, thereby missing out on a lot of information. This work will attempt to produce 3D maps of the distribution of HI (atomic hydrogen) and $H_2$ (molecular hydrogen) in the Milky Way using novel Bayesian inference techniques, exploiting spatial correlations of the gas structure to regularise ambiguities. Our approach will not only yield maps of the estimated gas densities, but also uncertainty information.

The two observational datasets used are that of the HI4PI-survey [1] mapping the 21-cm emission of atomic hydrogen in the galaxy (see figure 1) and the CO-survey compilation by [2] observing the $1 \rightarrow 0$ rotational transition of CO as a tracer for molecular clouds and thereby $H_2$ (see figure 2).

This work builds on precursory reconstructions (see [3, 4]) with some key differences:

1. A different numerical grid is used trading resolution far away from the observer for a much more refined resolution nearby.

2. The inference of galactic HI and $H_2$ is unified into a common inference process coupled by a common galactic velocity field.

3. The galactic velocity field is partly inferred, modifying our prior assumption by adding a curl-free field.

In the following section, we will formulate this problem in a Bayesian manner and shortly describe the used approach to this very-high-dimensional problem. Thereafter, we will show our preliminary results, i.e. 3D-maps of the distribution of HI and $H_2$ in the galaxy.

## 2. Method

### 2.1 Bayesian formulation

In the language of probabilities, we want to know the probability of the gas distribution in the galaxy (called signal $s$) given the data $d$ obtained by the sky surveys. Using Bayes' law, this can be written as

$$P(s|d) = \frac{P(s,d)}{P(d)} = \frac{P(d|s)P(s)}{P(d)} \,. \tag{1}$$

Since the datasets we are using measured the brightness temperature $T$ (a measure of the intensity of the observed radiation) as a function of relative line-of-sight velocity and position on the sky, we will attempt to infer the CO and HI volume emissivity $s = (\varepsilon_{\text{HI}}(\vec{x}), \varepsilon_{\text{CO}}(\vec{x}))$ simultaneously and later convert to gas densities.
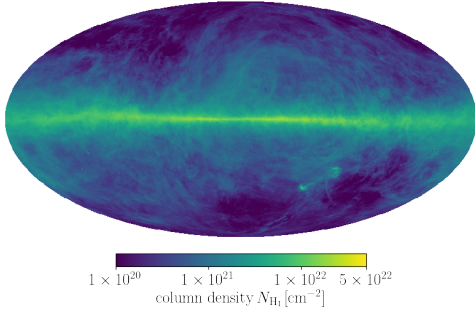
Equation 1 is often solved by creating a model that allows to sample from the prior distribution and compute the likelihood of said sample. Then, algorithms like MCMC sampling can be used to probe the shape of the posterior probability. For problems with many free parameters (usually more than $\gtrsim 10^2$), this becomes computationally unfeasible. A commonly used approach for high-dimensional problems are so-called *Variational Inference*(VI)-methods [5]. The idea of these methods is to approximate the posterior by a family of parametric distributions, for example a multi-variate Gaussian distribution.

The parameters of this approximation can be determined by minimising the "distance" between the approximated posterior and the true posterior, for example via the Kullback-Leibler-divergence [6]. Computing this in theory involves the inversion of the full covariance matrix of all the correlated latent variables which - with millions or more of free parameters - is impossible even to store in common memory modules. As an approximation, it has been suggested to replace the inverse of the full covariance matrix by the inverse Fisher information metric, an approach known as Metric Gaussian Variational Inference (MGVI, [7]). This method can be applied to problems with more than $10^6$ parameters while still being computationally efficient on regularly available hardware.
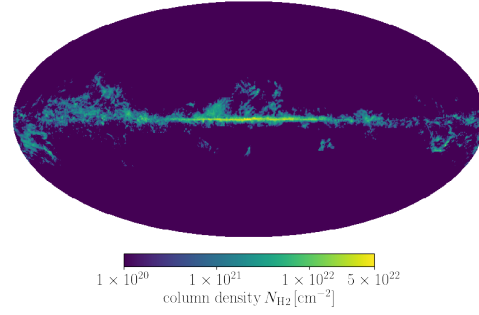
This algorithm is implemented in an iterative scheme in the publicly available code-package NIFTY8[1]. It alternates between estimating the covariance of the probability distribution with the inverse Fisher information metric at the current mean and optimising the mean of the distribution by minimising the Kullback-Leibler divergence to the true posterior with respect to the mean. This does not require explicitly computing the covariance matrix at any point (which would require the inversion of the Fisher information metric): instead the Kullback-Leibler divergence is estimated stochastically by drawing samples from a Gaussian with the appropriate covariance, leading to linear scaling in the model parameters. This can be implemented in terms of implicit operators which apply the Fisher information metric to some vector and then solving a linear system via conjugate gradient methods to obtain the application of the inverse Fisher information metric to some vector (which then features the desired correlation structure).

This way, a set of samples of the posterior distribution is obtained from which the Kullback-Leibler divergence can be calculated and, in turn, minimised. The algorithm has converged once the estimate for the mean and the estimate for the covariance are self-consistent. The result of this

---

[1]Available at https://gitlab.mpcdf.mpg.de/ift/NIFTy

**Figure 1:** HI column density map from the 21cm-data by the HI4PI collaboration [1]



**Figure 2:** H$_2$ column density map from the CO-emission-data compiled by Dame et al. [2]

algorithm is a set of samples of the approximated posterior distribution, implicitly containing the correlation structure between all model parameters. To apply this algorithm, we thus need

1. A model that allows drawing prior samples from a set of latent variables, taking into account the spatial correlation structure of the 3D gas distribution (the *signal*)

2. A connection between the drawn gas realisation and the expected measurement data (the *response*)

## 2.2 Gas model and Signal

Since the data from sky surveys has some fixed angular resolution, it is wise to represent the gas density on a grid that shares this property. If we chose to represent the gas density on a regular x-y-z grid (as in [3, 4]), voxels nearby would occupy almost half of the sky while voxels far away occupy an area on the sky much smaller than the available data resolution. In order to be consistent with the data resolution we thus choose to represent our signal data on a HEALPix-grid on angular direction and a logarithmic grid in radial direction. This ensures a high resolution nearby where we expect to be the most sensitive to the actual gas distribution. This choice will also make the response-function trivial as the otherwise costly line-of-sight integration reduces to a simple sum along the radial direction of the grid.

To model our prior, we generate samples of correlated lognormal random fields. These are obtained by drawing an – initially white-noise – sample $\xi(\vec{x})$ of latent variables and correlating it using a method called Iterative Charted Refinement [8] according to a Matérn-covariance function. We infer the parameters of this correlation structure at the same time as the gas density. The result is a correlated Gaussian random field $g(\vec{x})$. Upon exponentiation, we obtain a lognormal correlation structure. This ensures positive gas densities while also allowing for large density differences as are expected to be present in the interstellar medium. This is not yet a very good prior assumption for gas in the galaxy as most of the gas is tightly constrained to the galactic disk which has a small scale height ($\approx 150$ pc) compared to its diameter ($\approx 15$ kpc). This can be immediately seen in the data-sets (figures 1 and 2): most of the gaseous emission is concentrated around latitude zero. The fidelity of the reconstruction can be increased by explicitly modelling the inhomogeneous large-scale variations. We therefore multiply the correlated field with a profile in z-direction and in

radial direction:

$$P_z(\vec{x}) = P_z(z) = \exp\left(\frac{-|z|}{z_h}\right), \tag{2}$$

$$P_{\text{rad}}(\vec{x}) = P_{\text{rad}}(r_{\text{gal}}) = \exp\left(\frac{R_{\text{cutoff}} - r_{\text{gal}}}{R_{\text{scale}}}\right), \text{ for } r_{\text{gal}} > R_{\text{cutoff}}, \text{ else } 1. \tag{3}$$

Using this, we obtain

$$\epsilon(\vec{x}) = A \cdot P_z(\vec{x}) \cdot P_{\text{rad}}(\vec{x}) \cdot \exp\left(g(\vec{x})\right) \tag{4}$$

for HI and H$_2$ respectively. For the HI-profile, we choose $z_h = z_h(r_{\text{gal}}) = 150\,\text{pc} \cdot \exp\left(\frac{r_{\text{gal}} - R_\odot}{9.8\,\text{kpc}}\right)$ for $r_{\text{gal}} > 5\,\text{kpc}$, $R_{\text{scale}} = 3.15\,\text{kpc}$ and $R_{\text{cutoff}} = 7.0\,\text{kpc}$ as suggested by [9]. For the H$_2$-profile, we use $z_h = 50\,\text{pc}$, $R_{\text{scale}} = 1.0\,\text{kpc}$ and $R_{\text{cutoff}} = 8.0\,\text{kpc}$. This does not prevent the inference from reconstructing fields that differ from this profile, but ensures that the drawn prior samples feature a galactic-disk-like gas distribution.

## 2.3 Data and Response

The second ingredient is the response function that connects the signal to the observation by modelling the generation of synthetic data from a signal sample. For simplicity, we work in the optically thin limit and ignore any absorption effects. In this case, the measured brightness temperature $T(\hat{n}, v)$ in some direction $\hat{n}$ Doppler-shifted by a velocity difference $v$ is related to the volume emissivity $\epsilon(\vec{x})$ by a linear response map $R$ via
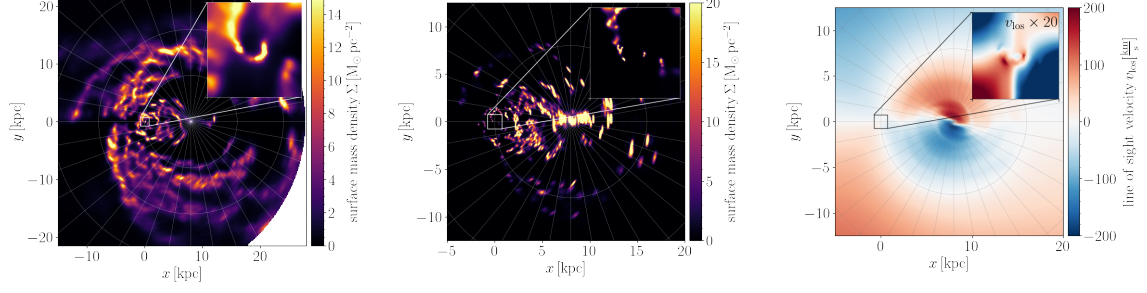
$$T(\hat{n}, v) = R(\epsilon(\vec{x})) = \int_0^\infty \mathrm{d}r\, \epsilon(\vec{x}) \delta(v - v_{\text{LSR}}(\vec{x})), \tag{5}$$

where $v_{\text{LSR}}(\vec{x})$ is the relative line-of-sight velocity at the position $(\vec{x})$ in the local standard of rest as dictated by our velocity model and $r = |\vec{x}|$ is the distance from Earth. We approximate the Dirac-delta by a Gaussian with a width of $\sigma_{\text{HI}} = 10\frac{\text{km}}{s}$ [10] and $\sigma_{\text{H}_2} = 5\frac{\text{km}}{s}$ [11] in order to account for velocity dispersion inside gas clouds.

For the velocity model, we use a fixed component based on a smoothed particle hydrodynamics simulation by [12], extended beyond 8 kpc using a flat rotation curve. On top of this velocity field, we add another component computed as the gradient of a scalar velocity potential:

$$v_{\text{LSR}}(\vec{x}) = \vec{v}_0(\vec{x}) + \nabla S(\vec{x}). \tag{6}$$

This scalar field will be modelled as a correlated Gaussian random field and reconstructed at the same time as the gas emissivities . This opens the possibility for the model to adjust the velocity-field during the reconstruction. Additionally, we can hope to learn something about the true velocity field in areas, where the data is very constraining. However, there is no direct velocity-information in the data if one does not demand that the resulting gas densities should follow a certain correlation structure. Even then, this introduces many ambiguities as it greatly expands the possibilities, where gas clouds can be mapped. In the future, we will have to add additional data to either constrain the gas densities tighter (e.g. using correlations with dust [13]), thereby learning about the velocities or constrain the velocities tighter (e.g. using parallax information of masers or young stars [14]), thereby learning about the gas densities.

**Figure 3:** Results of the inference with zoom-in on the local neighbourhood. Left panel: HI surface mass density. Middle panel: H$_2$ surface mass density. Right panel: Line-of-sight velocity at zero latitude

## 2.4 Noise and Likelihood

Taking into account additive noise in the observations, we alter our model for the relation between brightness temperature and volume emissivity to

$$T(\hat{n}, v) = R(\epsilon(\vec{x})) + n \,, \tag{7}$$

where we assume the noise $n$ to be normal distributed and uncorrelated (white) with some diagonal covariance $N$. The likelihood can then be written as

$$p(T|\epsilon) = \int \mathrm{d}n\, p(T|\epsilon, n)p(n) = \int \mathrm{d}n\, \delta(T - R(\epsilon) - n)\mathcal{G}(n, N) = \mathcal{G}(T - R(\epsilon), N) \,. \tag{8}$$

## 3. Results

We run our reconstruction with a resolution of NSIDE = 32 in angular direction and 500 radial pixels between $r_{\min} = 50\,\mathrm{pc}$ and $r_{\max} = 28\,\mathrm{kpc}$. The sample-average surface mass densities resulting from the 3D-maps can be seen in figure 3 for HI and H$_2$. This figure also shows the (partly reconstructed) sample-average line-of-sight velocity in a zero-latitude slice. The reconstructed gas densities show disk-like structures of gas clusters with imprints of galactic arms that are particularly visible in the HI-gas reconstruction. Both gas reconstructions suffer from a set of problems that we will discuss in the following.

## 3.1 Inferred HI gas density

The outside of the solar circle is well populated with HI-gas, whereas in the inside of the solar circle, the distance ambiguity seems to be resolved very one-sidedly towards the nearer solution. This could be due to the logradial grid giving the algorithm the opportunity to place gas nearby with a much higher fidelity than far away. This can then reproduce the data with a much higher fidelity as well leading to a much higher likelihood. This could be tested and perhaps solved by e.g. modifying the grid to have a uniform resolution inside the solar circle or by increasing the total resolution until saturation. The nearby gas shows a circular structure at the same radius as the reconstructed velocity as well as a tilted line-like structure in negative $x$-direction.

### 3.2  Inferred H$_2$ gas density

The quality of reconstruction of H$_2$ appears to be worse than that of HI, mainly due to the fact that the gas is very concentrated at the plane $z = 0$ and the amount of grid-points that are far away and very close to the galactic plane become very few. One can clearly see circular structures in the gas projection stemming from the (too) low angular resolution. Clearly visible is a bar-like structure in the galactic centre as well as two "wall"-like structures in-between us and the galactic centre also seen in precious reconstructions on a regular grid [3]. The excellent local resolution lets us see fine structures in the nearby gas showing a similar structure as the HI-gas in negative $x$-direction and small nearby clouds of gas towards the galactic centre.

### 3.3  Inferred velocity field

The reconstructed curl-free modification of the velocity prior is in general very small in amplitude and negligible for distances larger than 1 kpc. For distances smaller than that, there is an almost circular, positive (amplitude up to $10\frac{km}{s}$) velocity correction being reconstructed. The position and amplitude coincide nicely with estimates for the expansion velocity of the local bubble [15]. It is not clear, how much information about the velocity field itself is contained in the data but the combination of two data-sets having to respect the local correlation structure at the same time appears to provide at least some information.

## 4.  Conclusion

We present new preliminary 3D-maps of galactic HI and H$_2$ inferred in conjunction using the same velocity field. We also present a partly reconstructed 3D line-of-sight velocity map featuring a circular structure with outwards-pointing velocities in the local neighbourhood. The ingredients for our inference were the HI4PI-survey from [1] measuring 21cm-emission, the CO-survey compilation by [2] measuring rotational CO-transitions and a VI-algorithm capable of inferring millions of parameters [7]. We have assumed the optically thin limit. In the future we plan to improve upon these shortcoming by including additional data and thereby further constraining either the velocity field or the gas distributions, by lifting our assumption on the optical thinness of the gas; and by improving the angular resolution of our reconstructions.

## References

[1] HI4PI Collaboration, "HI4PI: A full-sky H I survey based on EBHIS and GASS," *Astronomy and Astrophysics*, vol. 594, p. A116, Oct. 2016.

[2] T. M. Dame, D. Hartmann, and P. Thaddeus, "The Milky Way in Molecular Clouds: A New Complete CO Survey," *Astrophysical Journal*, vol. 547, pp. 792–813, Feb. 2001.

[3] P. Mertsch and A. Vittino, "Bayesian inference of three-dimensional gas maps. I. Galactic CO," *Astronomy and Astrophysics*, vol. 655, p. A64, Nov. 2021.

[4] P. Mertsch and V. H. M. Phan, "Bayesian inference of three-dimensional gas maps. II. Galactic HI," *Astronomy and Astrophysics*, vol. 671, p. A54, Mar. 2023.

[5] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, "Variational inference: A review for statisticians," *Journal of the American Statistical Association*, vol. 112, no. 518, pp. 859–877, 2017.

[6] S. Kullback and R. A. Leibler, "On Information and Sufficiency," *The Annals of Mathematical Statistics*, vol. 22, no. 1, pp. 79–86, 1951.

[7] J. Knollmüller and T. A. Enßlin, "Metric Gaussian Variational Inference," *arXiv e-prints*, p. arXiv:1901.11033, Jan. 2019.

[8] G. Edenhofer, R. H. Leike, *et al.*, "Sparse Kernel Gaussian Processes through Iterative Charted Refinement (ICR)," *arXiv e-prints*, p. arXiv:2206.10634, June 2022.

[9] P. M. W. Kalberla and L. Dedes, "Global properties of the H I distribution in the outer Milky Way. Planar and extra-planar gas," *Astronomy and Astrophysics*, vol. 487, pp. 951–963, Sept. 2008.

[10] H. Nakanishi and Y. Sofue, "Three-Dimensional Distribution of the ISM in the Milky Way Galaxy: I. The H I Disk," *Publications of the Astronomical Society of Japan*, vol. 55, pp. 191–202, Feb. 2003.

[11] M. Pohl, P. Englmaier, and N. Bissantz, "Three-dimensional distribution of molecular gas in the barred milky way," *The Astrophysical Journal*, vol. 677, p. 283, apr 2008.

[12] N. Bissantz, P. Englmaier, and O. Gerhard, "Gas dynamics in the Milky Way: second pattern speed and large-scale morphology," *Monthly Notices of the RAS*, vol. 340, pp. 949–968, Apr. 2003.

[13] R. H. Leike, G. Edenhofer, *et al.*, "The Galactic 3D large-scale dust distribution via Gaussian process regression on spherical coordinates," *arXiv e-prints*, p. arXiv:2204.11715, Apr. 2022.

[14] M. J. Reid, K. M. Menten, *et al.*, "Trigonometric parallaxes of high-mass star-forming regions: Our view of the milky way," *The Astrophysical Journal*, vol. 885, p. 131, nov 2019.

[15] C. Zucker, A. A. Goodman, *et al.*, "Star formation near the Sun is driven by expansion of the Local Bubble," *Nature*, vol. 601, pp. 334–337, Jan. 2022.