# Inferring Galactic Parameters from Chemical Abundances with Simulation-Based Inference

Tobias Buck[1,2], Berkay Günes[1,2], Giuseppe Viterbo[1,2], William H. Oliver[1,2], and Sven Buder[3,4]

[1] Interdisciplinary Center for Scientific Computing (IWR), University of Heidelberg, Im Neuenheimer Feld 205, D-69120 Heidelberg, Germany
[2] Universität Heidelberg, Zentrum für Astronomie, Institut für Theoretische Astrophysik, Albert-Ueberle-Straße 2, D-69120 Heidelberg, Germany
[3] Research School of Astronomy and Astrophysics, Australian National University, Canberra, ACT 2611, Australia
[4] ARC Centre of Excellence for All Sky Astrophysics in 3 Dimensions (ASTRO 3D), Australia
   e-mail: tobias.buck@iwr.uni-heidelberg.de

**ABSTRACT**

*Context.* Galactic chemical abundances provide crucial insights into fundamental galactic parameters, such as the high-mass slope of the initial mass function (IMF) and the normalization of Type Ia supernova (SN Ia) rates. Constraining these parameters is essential for advancing our understanding of stellar feedback, metal enrichment, and galaxy formation processes. However, traditional Bayesian inference techniques, such as Hamiltonian Monte Carlo (HMC), are computationally prohibitive when applied to large datasets of modern stellar surveys.
*Aims.* We leverage simulation-based-inference (SBI) as a scalable, robust, and efficient method for constraining galactic parameters from stellar chemical abundances and demonstrate its the advantages over HMC in terms of speed, scalability, and robustness against model misspecifications.
*Methods.* We combine a Galactic Chemical Evolution (GCE) model, `CHEMPY`, with a neural network emulator and a Neural Posterior Estimator (NPE) to train our SBI pipeline. Mock datasets are generated using `CHEMPY`, including scenarios with mismatched nucleosynthetic yields, with additional tests conducted on data from a simulated Milky Way-like galaxy. SBI results are benchmarked against HMC-based inference, focusing on computational performance, accuracy, and resilience to systematic discrepancies.
*Results.* SBI achieves a $\sim 75,600\times$ speed-up compared to HMC, reducing inference runtime from $\gtrsim 42$ hours to mere seconds for thousands of stars. Inference on 1,000 stars yields precise estimates for the IMF slope ($\alpha_{\mathrm{IMF}} = -2.298 \pm 0.002$) and SN Ia normalization ($\log_{10}(N_{\mathrm{Ia}}) = -2.885 \pm 0.003$), deviating less than 0.05% from the ground truth. SBI also demonstrates similar robustness to model misspecification than HMC, recovering accurate parameters even with alternate yield tables or data from a cosmological simulation.
*Conclusions.* SBI represents a paradigm shift in GCE studies, enabling efficient and precise analysis of massive stellar datasets. By outperforming HMC in speed, scalability, and robustness, SBI is poised to become a cornerstone methodology for future spectroscopic surveys facilitating deeper insights into the chemical and dynamical evolution of galaxies.

**Key words.** Galaxies: fundamental parameters – Galaxies: stellar content – Methods: data analysis – Methods: statistical –

## 1. Introduction

Understanding the chemical enrichment of galaxies is fundamental to deciphering their formation and evolution. Chemical abundances of stars offer a wealth of information about galactic parameters, such as the high-mass slope of the initial mass function (IMF) and the normalization of Type Ia supernova (SN Ia) rates. These parameters critically influence the production of heavy elements (e.g. Romano et al. 2005; Vincenzo et al. 2015; Mollá et al. 2015), stellar feedback, and star formation histories, making their accurate determination essential for realistic hydrodynamical simulations of galaxy formation (e.g. Sawala et al. 2016; Hopkins et al. 2018; Pillepich et al. 2018b; Buck 2020; Buck et al. 2020, 2021; Font et al. 2020; Agertz et al. 2021). Despite their importance, constraining these parameters has proven challenging due to limited observational data and the computational demands of traditional inference techniques.

For example, a range of high-mass IMF slopes have been suggested (Côté et al. 2016, Tab. 7), with a steeper-than-canonical slope being suggested by a range of studies (e.g. Weisz et al.

2015; Rybizki & Just 2015; Chabrier et al. 2014). In addition, the IMF slope may itself be not a constant but rather a function of metallicity, introducing further complexity (e.g. Gutcke & Springel 2019; Martín-Navarro et al. 2019). Similarly, the choice of SN Ia delay-time-distribution and normalization plays a crucial role in the enrichment of the interstellar medium (ISM; e.g. Buck et al. 2021) and is heavily debated (Maoz et al. 2010, 2012; Jiménez et al. 2015).

Recent advances in stellar spectroscopic surveys, such as APOGEE (Abdurro'uf et al. 2022) and GALAH (Buder et al. 2021, 2024), have produced unprecedented datasets of stellar chemical abundances across a third of the period table. These datasets hold the potential to unlock detailed constraints on galactic parameters across diverse environments. However, traditional Bayesian inference methods, such as Markov Chain Monte Carlo (MCMC) and Hamiltonian Monte Carlo (HMC), struggle to scale to these large datasets. Such methods are computationally expensive, requiring hours of runtime for even modest sample sizes, and are susceptible to biases when confronted with high-dimensional posterior distributions.
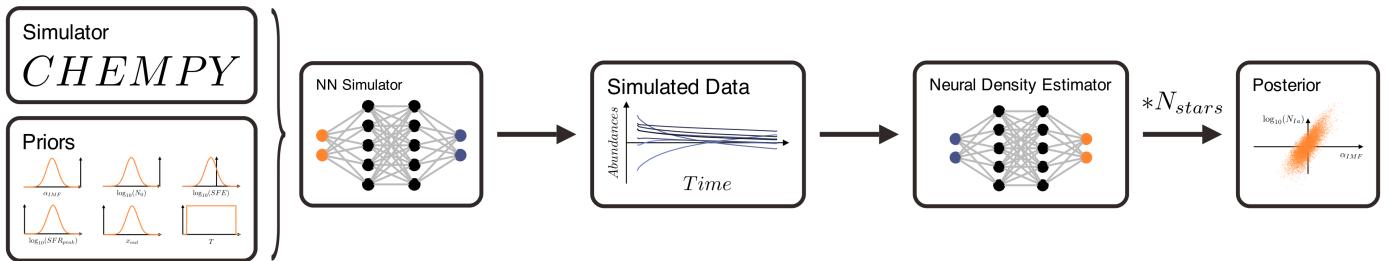
**Fig. 1.** SBI flow chart. From a set of priors we simulate a sample of stellar abundances using CHEMPY (Rybizki et al. 2017; Philcox & Rybizki 2019) which we use to train a *neural network* emulator to speed up the data generation process. Using the *neural network* emulator we produce training data to train the Neural Density Estimator. With this we infer the posterior distribution of the model parameters from a single star. Repeating that for $N_{stars}$ from the same galaxy gives an accurate fit of the IMF slope and Type Ia supernovae normalization.

In this work, we present a novel approach leveraging simulation-based inference (SBI, e.g. Cranmer et al. 2020) to address these limitations. SBI bypasses the need for explicit likelihood functions, enabling efficient and scalable inference of galactic parameters directly from simulated stellar abundances. By combining a neural network emulator for the CHEMPY Galactic Chemical Evolution (GCE) model with a Neural Posterior Estimator (NPE), we achieve rapid and robust inference. Unlike HMC, which requires extensive sampling for each dataset, our method amortizes the computational cost during training, allowing subsequent inference to scale seamlessly to larger datasets.

This study focuses on two critical global galactic parameters: the high-mass slope of the Chabrier (2003, Tab. 1) IMF ($\alpha_{\mathrm{IMF}}$) and the SN Ia normalization, $\log_{10}(N_{\mathrm{Ia}})$, the rate of SN Ia explosions per unit mass. We demonstrate the accuracy, scalability, and robustness of SBI through tests on mock datasets generated by CHEMPY (Rybizki et al. 2017), as well as on data from hydrodynamical simulations. Additionally, we compare our results to those obtained using HMC-based inference on the same datasets (see Philcox & Rybizki 2019), highlighting SBI's superior performance in terms of speed, precision, and resilience to model misspecification.

The structure of this paper is as follows: In Section 2, we outline the methods used, including the GCE model and SBI framework. Section 3 presents our results on both CHEMPY and IllustrisTNG (Pillepich et al. 2018b) mock data, emphasizing SBI's advantages over traditional approaches. Finally, in Section 4, we discuss the broader implications of our findings and outline potential future applications, before concluding in Section 5.

Finally, we publicly release all of our code to reproduce the results of this manuscript via GitHub[1] and refer to Appendix A for a more extended overview of our code availability. All our datasets and network weights are publicly available on Zenodo.[2]

## 2. Methods

In order to establish our new method based on SBI we need two ingredients: A simulator (in our case a GCE model) that simulates observational data from a set of model parameters (in our case the IMF slope and the Type Ia supernovae normalization) and a flexible way of parametrizing the posterior density conditioned on the observation in order to perform our inference (see Fig. 1 for a schematic visual representation of our method). In the next subsections we describe both ingredients in detail.

### 2.1. Galactic chemical evolution models

Our simulator is based on the CHEMPY model (Rybizki et al. 2017). CHEMPY is a simple GCE model that is able to predict stellar chemical abundances throughout cosmic time by using published nucleosynthetic yield tables for three key processes (SN Ia and SN II explosions and AGB stellar feedback) and a small number of parameters controlling simple stellar populations (SSPs) and ISM physics. We refer to the initial CHEMPY paper (Rybizki et al. 2017) for the details of the model.

In particular, we are using the CHEMPYScoring module (Philcox et al. 2018) publicly available as the CHEMPYMulti (Philcox & Rybizki 2019)[3] package a further development of the original CHEMPY model.

CHEMPY parameters: In this work, we allow six CHEMPY parameters to vary freely (see also Tab. 1). These can be categorized into three groups:

1. $\mathbf{\Lambda}$: **Global Galactic Parameters** describe SSP physics and comprise the high-mass Chabrier (2003) IMF slope, $\alpha_{\mathrm{IMF}}$, which effectively sets the number of CC-SNe a SSP generates and (logarithmic) Type Ia SN normalization, $\log_{10}(N_{\mathrm{Ia}})$, which controls the total number of SN Is per SSP. We treat these as star-independent and assume them to be constant across galactic environments and cosmic time[4]. We adopt the same broad priors as (Philcox & Rybizki 2019) for these variables (see also Tab.1).

2. $\{\mathbf{\Theta}_i\}$: **Local Galactic Parameters** describe the local physics of the ISM and are hence specific to each stellar environment, indexed by $i$. As defined in (Rybizki et al. 2017), these include the star-formation efficiency (SFE), $\log_{10}(\mathrm{SFE})$, which qunatifies the star formation rate per unit gas, the peak of the star formation rate (SFR), $\log_{10}(\mathrm{SFR_{peak}})$, and the fraction of stellar outflow that is fed to the gas reservoir, $\mathrm{x_{out}}$. We adopt broad priors for all parameters and, as in (Philcox & Rybizki 2019), fix the SN Ia delay-time distribution, $\log_{10}(\tau_{\mathrm{Ia}})$, to $\log_{10}(\tau_{\mathrm{Ia}}) = -0.80$ (see also Philcox et al. 2018).

3. $\{T_i\}$: **Stellar Birth-Times**. Time in Gyr at which a given star is formed from the ISM. We assume that its proto-stellar abundances match the local ISM abundances at $T_i$.

The separability of local (ISM) parameters and global (SSP) parameters is motivated by recent observational evidence: Ness

---

[3] github.com/oliverphilcox/ChempyMulti
[4] Whilst $\log_{10}(N_{\mathrm{Ia}})$ is constant with respect to time by definition, it being simply a normalization constant, there is some evidence for $\alpha_{\mathrm{IMF}}$ varying as a function of time or metallicity (Chabrier et al. 2014; Clauwens et al. 2016; Gutcke & Springel 2019; Martín-Navarro et al. 2019).

**Table 1.** Free `CHEMPY` parameters for each star, with their prior values and Gaussian widths. Stellar birth-times are set for each star individually from a Uniform prior, based on realistic age estimates.

| Parameter | Description | $\overline{\theta}_{prior} \pm \sigma_{prior}$ | Prior from: |
|---|---|---|---|
| | **$\Lambda$: *Global stellar (SSP) parameters*** | | |
| $\alpha_{IMF}$ | High-mass slope of the (Chabrier 2003) IMF | $-2.3 \pm 0.3$ | (Chabrier 2003, Tab. 1) |
| $\log_{10}(N_{Ia})$ | Number of SN Ia per $M_\odot$ over 15 Gyr | $-2.89 \pm 0.3$ | (Maoz & Mannucci 2012, Tab.1 ) |
| | **$\Theta_i$: *Local ISM parameters*** | | |
| $\log_{10}(SFE)$ | Star formation efficiency governing gas infall | $-0.3 \pm 0.3$ | (Bigiel et al. 2008) |
| $\log_{10}(SFR_{peak})$ | SFR peak in Gyr (scale of $k = 2$ $\Gamma$-distribution) | $0.55 \pm 0.1$ | (van Dokkum et al. 2013, fig. 4b) |
| $x_{out}$ | Stellar feedback fraction | $0.5 \pm 0.1$ | (Rybizki et al. 2017, Tab. 1) |
| | **$T_i$: *Timescale*** | | |
| $T_i$ | Time of stellar birth in Gyr | $[1, 13.8]$ | Observations |

et al. (2019) find that the elemental abundances of red clump stars belonging to the thin disk can be predicted almost perfectly from their age and [Fe/H] abundance. This implies that the key chemical evolution parameters affecting the elemental abundances (SSP parameters and yield tables) are held fixed, whilst ISM parameters vary smoothly over the thin disk (which offsets the metallicity for different galactocentric radii, e.g. Buck 2020; Wang et al. 2024, for a simulated example). Similarly Weinberg et al. (2019) find that ISM parameter variations are deprojected in the [X/Mg] vs [Mg/H] plane (their Fig. 17) and that abundance tracks in that space are independent of the stellar sample's spatial position within the Galaxy (their Fig. 3).

Following Philcox & Rybizki (2019), to avoid unrealistic star formation histories (that are very 'bursty' for early stars), we additionally require that the SFR (parametrized by a $\Gamma$ distribution with shape parameter $a = 2$[5]) at the maximum possible stellar birth-time (13.8 Gyr) should be at least 5% of the mean SFR, ensuring that there is still a reasonable chance of forming a star at this time-step. In our formalism, this corresponds to the constraint $\log_{10}(SFR_{peak}) > 0.294$. For this reason, a truncated Gaussian prior will be used for the SFR parameter. Furthermore, we constrain $T_i$ to the interval $[1, 13.8]$ Gyr (assuming an age of the Universe of 13.8 Gyr), ignoring any stars formed before 1 Gyr, which is justified as these are expected to be rare.

**Nucleosynthetic yield tables:** We adopt the same nucleosynthetic yield tables as in (Philcox & Rybizki 2019), see their Sec. 2.2 for more details. To test our method, we aim further at inferring parameters from a sample of stars taken from a hydrodynamical simulation of a MW type galaxy which we take from the IllustrisTNG project (Pillepich et al. 2018b). To ensure maximal compatibility with TNG, we adopt their nucleosynthetic yield tables in `CHEMPY`, for enrichment by SN Ia, SN II and AGB stars. The utilized yields are summarized in Tab. 2, matching Pillepich et al. (2018c, Tab. 2), and we note that the SN II yields are renormalized such that the IMF-weighted yield ratios at each metallicity are equal to those from the Kobayashi et al. (2006) mass range models alone. `CHEMPY` uses only net yields, such that they provide only newly synthesized material, with the remainder coming from the initial SSP composition. These tables may not well-represent true stellar chemistry, and the effects of this mis-

**Table 2.** Nucleosynthetic yield tables used in this analysis, matching those of the TNG simulation (Pillepich et al. 2018c, Tab. 2).

| Type | Yield Table |
|---|---|
| SN Ia | Nomoto et al. (1997) |
| SN II | Kobayashi et al. (2006); Portinari et al. (1998) |
| AGB | Karakas (2010); Doherty et al. (2014); Fishlock et al. (2014) |

match are examined in Sec. 3.3 by performing inference using an alternative set of yields that does not match the yield set of the training data. For the analysis of observational data, we would want to use the most up-to-date yields, such as Karakas & Lugaro (2016) AGB yields, and carefully choose elements which are known to be well reproduced by our current models (e.g. shown by Weinberg et al. (2019); Griffith et al. (2019)), though this is not appropriate in our context. To facilitate best comparison with Ilustris TNG, we further set the maximum SN II mass as $100\,M_\odot$ (matching the IMF upper mass limit), adopt stellar lifetimes from Portinari et al. (1998) and do not allow for any 'hypernovae' (in contrary to Philcox et al. 2018)).

**Chemical elements:** In our analysis we only track nine elements: C, Fe, H, He, Mg, N, Ne, O and Si since these are the only elements traced by TNG. We principally compare the logarithmic abundances [X/Fe] and [Fe/H] defined by

$$[X/Y] = \log_{10}(N_X/N_Y)_{star} - \log_{10}(N_X/N_Y)_\odot \qquad (1)$$

for number fraction $N_X$ of element X. Here $\odot$ denotes the solar number fractions of Asplund et al. (2009). As is customary, we use H for normalization, thus we are left with $n_{el} = 8$ independent elements which must be tracked by `CHEMPY`[6].

With these modifications, `CHEMPY` allows to predict TNG-like chemical abundances for a given set of galactic parameters. It is important to note that the two GCE models have very different parametrizations of galactic physics, with TNG including vastly more effects, it is thus not certain *a priori* how useful `CHEMPY` will be in emulating the TNG simulation, although its utility was partially demonstrated in Philcox et al. (2018). However, such a test is necessary to prepare for an inference on real data.

---

[5] `CHEMPY` parametrizes the SFR with a $\Gamma$ distribution

$$SFR(t, k, \vartheta) = \frac{1}{\Gamma(k)\vartheta^k} t^{k-1} \exp\left(\frac{-t}{\vartheta}\right), \text{ for } k = 2 \rightarrow \vartheta = SFR_{peak}$$

where the shape parameter is fixed to $k = 2$ such that the scale parameter ($\vartheta$) determines the peak of the SFR

---

[6] In observational contexts, it may be more appropriate to compute abundances relative to Mg rather than Fe, as in (Weinberg et al. 2019), since Mg is only significantly produced by SN II and hence a simpler tracer of chemical enrichment.

## 2.2. Neural network emulator for CHEMPY

Despite the simplifications made by the GCE model CHEMPY, the run-time of CHEMPY and the high-dimensionality of the parameter space incurs some difficulties when sampling the distribution of the global parameters $\Lambda = \{\alpha_{\text{IMF}}, \log_{10}(N_{\text{Ia}})\}$. To alleviate this, we follow Philcox & Rybizki (2019) and implement a *neural network* (NN) emulator of the CHEMPY simulator. We design the NN as a simple feed-forward neural network with 2 hidden layers and 100 neurons in the first and 40 neurons in the second layer. The NN is trained on $\sim 700,000$ data points and validated on $\sim 50,000$ additional data points created with CHEMPY using a uniform prior over the $5\sigma$-range of the original Gaussian prior stated in table 1. The batch size is set to 64 and the learning rate is set to 0.001. We train for 20 epochs using a schedule free optimizer (Defazio et al. 2024). Training this tiny emulator takes about 200s on the CPU.

In essence, instead of computing the full model for each input parameter set, we pass the parameters to the NN which predicts the output abundances to high accuracy. As already argued in Philcox & Rybizki (2019) this has two benefits;

1. **Speed:** The run-time of the CHEMPY function is $\sim 1$ s per input parameter set, which leads to very slow generation of training data for SBI. With the NN emulator, this reduces to $\sim 5 \times 10^{-5}$ s, and is trivially parallelizable, unlike CHEMPY. Having access to a fast simulator opens the possibility for testing Sequential Neural Posterior Estimate (SNPE) as an alternative, but we discuss this possibility in Section 4.

2. **Differentiability:** The NN is written in pytorch and has additionally a simple closed-form analytic structure (described in the appendix of Philcox & Rybizki 2019), unlike the complex CHEMPY model. Both aspects allow it to be differentiated either via auto-diff or analytically, so one can use it to sample via advanced methods such HMC as done in (Philcox & Rybizki 2019).

Despite the additional complexity introduced by using multiple stellar data-points, our NN simply needs to predict the birth-time abundances for a single star (with index $i$) from a given set of six parameters; $\{\Lambda, \Theta_i, T_i\}$. The same NN can be used for all $n_{\text{stars}}$ stars (and run in parallel), reducing a set of $n_{\text{stars}}$ runs of CHEMPY to a single matrix computation. With the above network parameter choices, the NN predicts abundances with an absolute percentage error of $1.9^{+3.2}_{-1.2}$ % (which translates into a logarithmic error of 0.008 dex) which is far below typical observational errors and even smaller away from the extremes of parameter space (see Fig. 2)[7]. In fact, we will add additional observational uncertainty to our mock observational data later during training of the neural posterior estimator network.

## 2.3. Bayesian model

CHEMPY effectively is a Bayesian model for stellar abundances given a set of parameters $\{\Lambda, \Theta, T\}$ and Philcox & Rybizki (2019) extended the CHEMPY framework to be able to model multiple stellar data-points. Consider a given star with index $i$ that is born in some region of the ISM. This star will carry its own set of parameters $\{\Lambda, \Theta_i, T_i\}$, where $\Lambda$ are star independent and hence taken to be global parameters while the ISM parameters $\Theta_i$ and the birth-time $T_i$ are star specific. Using CHEMPY (or the trained
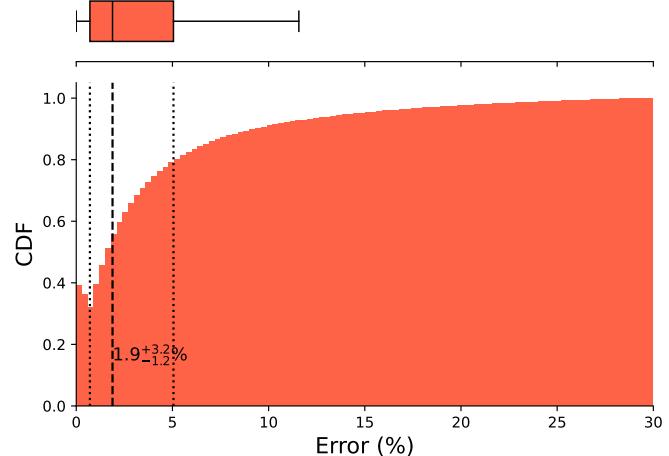
**Fig. 2.** Cumulative absolute percentage error of the NN emulator for the CHEMPY simulator. The orange histogram shows the cumulative distribution of percentage errors with the vertical dashed line indicating the median and the vertical dotted lines indicating the first and third quartile. The box plot on the top of the plot extends from the first quartile to the third quartile of the data, with a line at the median. The whiskers extend from the box to the farthest data point lying within 1.5× the inter-quartile range from the box. The NN predicts abundances with an absolute percentage error far below typical observational errors.

neural network emulator) we can easily model the set of $n_{\text{el}}$ chemical abundances $\{X_i^j\}$ for the $i$-th star as:

$$\{X_i^j\} = \text{CHEMPY}(\Lambda, \Theta_i, T_i), \tag{2}$$

where $j$ indexes the chemical element. These model abundances can then be compared against observations, with measured abundances $d_i^j$ and corresponding Gaussian errors $\sigma_{i,\text{obs}}^j$ jointly denoted as $D_i = \{d_i^j, \sigma_{i,\text{obs}}^j\}$.

**Posteriors for Galactic parameters** As stated in Philcox & Rybizki (2019) the full posterior for this case is given by

$$\mathbb{P}(\Lambda, \{\Theta_i\}, \{T_i\}|\{D_i\}) \propto \left[\prod_{i=1}^{n_{\text{star}}} p_\Theta(\Theta_i) p_{T_i}(T_i)\right] \times p_\Lambda(\Lambda) \tag{3}$$
$$\times \mathcal{L}(\{D_i\}|\Lambda, \{\Theta_i\}, \{T_i\})$$

where $p_\Theta(\Theta_i) p_{T_i}(T_i)$ are the priors on the variables $\Theta_i$ or $T_i$ belonging to a given set of stars.

In order to determine the optimal values of the global galactic parameters ($\Lambda$) one has to sample the posterior of Eq. 5. In practice this is a costly computation, since even with advanced techniques such as HMC sampling the posterior can only be evaluated for a small set of stars ($\lesssim 200$) and requires long compute times ($\sim 42$ hours Philcox & Rybizki 2019). However, recent advances in implicit-likelihood inference or SBI (Cranmer et al. 2020) offer another very efficient approach to approximate the posterior (see next paragraph for more details). These methods train a neural conditional density estimator to represent the conditional posterior, $\mathbb{P}(\Lambda|\{D_i\})$, which can be very efficiently evaluated given observational data $\{D_i\}$.

In particular, if we marginalize over the star specific parameters and solely focus on the global parameters $\Lambda$ we can make the assumption that individual observations of stars are identically and independently distributed (i.i.d.) and factorize the joint

posterior from above to simply express it as:

$$\mathbb{P}(\boldsymbol{\Lambda}|\{D_i\}) \propto \mathbb{P}(\boldsymbol{\Lambda})\mathbb{P}(\{D_i\}|\boldsymbol{\Lambda}) \quad \text{(Bayes rule)} \quad (4)$$

$$= \mathbb{P}(\boldsymbol{\Lambda})\mathbb{P}(D_1, ..., D_{n_{\text{star}}}|\boldsymbol{\Lambda})$$

$$= \mathbb{P}(\boldsymbol{\Lambda})\prod_{j=1}^{n_{\text{star}}} \mathbb{P}(D_j|\boldsymbol{\Lambda}) \quad \text{(i.i.d.)}$$

$$\propto \mathbb{P}(\boldsymbol{\Lambda})\prod_{j=1}^{n_{\text{star}}} \frac{\mathbb{P}(\boldsymbol{\Lambda}|D_j)}{\mathbb{P}(\boldsymbol{\Lambda})} \quad \text{(Bayes rule)}$$

$$= \mathbb{P}(\boldsymbol{\Lambda})^{1-n_{\text{star}}} \prod_{j=1}^{n_{\text{star}}} \mathbb{P}(\boldsymbol{\Lambda}|D_j) \quad (5)$$

This factorization of the joint posterior into single star posteriors $\mathbb{P}(\boldsymbol{\Lambda}|D_j)$ has the advantage that we do not need to specify the exact number of observational data points beforehand. We can simply train our neural density estimator on single star observations and then at evaluation time, we simply combine as many posterior estimates as we have observational data.

**Simulation-based inference:** In a nutshell, SBI (e.g. Cranmer et al. 2020; Papamakarios et al. 2021; Gloeckler et al. 2024) – also called likelihood free inference within a Bayesian inference framework – works as follows: given an assumed generative model $\mathcal{M}$ of parameters $\boldsymbol{\theta}$ (in our case a GCE model) and a set of simulated observations of individual stellar abundances $\boldsymbol{X}$ from that model, we train a mapping between the two to estimate the posterior distribution $p(\boldsymbol{\theta}|\boldsymbol{X}, \mathcal{M})$ of the generative model parameters $\boldsymbol{\theta}$ that reproduce the simulated observations $\boldsymbol{X}$. Once this mapping is trained, we can apply it to observations of stellar abundances $\boldsymbol{X_R}$ to infer $p(\boldsymbol{\theta}|\boldsymbol{X_R}, \mathcal{M})$. We show a schematic visual representation of our method in Fig. 1. Note, we do not need to know anything about $p(\boldsymbol{\theta}|\boldsymbol{X_R}, \mathcal{M})$, we solely need to be able to sample from it.

We use a NPE (Zeghal et al. 2022) which utilises the gradients of the generative model $\mathcal{M}$ with a Masked Autoregressive Flow (MAF; Papamakarios et al. 2018) containing 8 hidden features and 4 transformation layers for the normalizing flow. The expressivity of the MAF allows the NDE to be capture complex distributions, while also maintaining computational tractability. The final model was selected after extensive hyperparameter tuning, varying: the architecture between Neural Spline Flow (NSF; Durkan et al. 2019), MAF, and MAF with rational-quadratic spline (MAF-RQS), the number of neurons between 10, 20, 50, 100, and the number of transformations between 1, 5, 10, 30, optimizing over the test set for both the highest mean log posterior probability and the best calibration as measured by the TARP value (Lemos et al. 2023). For more details on the TARP value see section B in the appendix.

We train our NPE with $10^5$ data points consisting of $n_{\text{elements}} = 8$ simulated with the NN emulator described in the previous section. Inputs are sampled from the Gaussian priors shown in Tab. 1. Training takes $\sim 10$ minutes on an Apple M1 Max. We evaluate the accuracy of the NPE using $5 \times 10^3$ validation data points from the original CHEMPY simulator. In order to mimic observational noise, we add 5% observational uncertainties to the abundances simulated with CHEMPY before feeding them to our NPE. We have made sure that our NPE is well calibrated and posterior distributions are trustable. For more details on simulation-based calibration see appendix B.

Note, the methods presented here can naturally be extended to any fast and flexible GCE model (e.g. Talbot & Arnett 1971;

Johnson & Weinberg 2020; Côté & Ritter 2018, and others), not just CHEMPY and can even be used to infer vastly different galactic parameters such as accretion events from abundance distributions of stars (e.g. Viterbo & Buck 2024).

**Multi-star inference:** Following eq. 5 we can calculate the joint posterior for a combined inference using abundance observations of multiple stars. For this, we will condition our NPE individually on single star observations and sample the posterior. Then, according to eq. 5 we simply have to multiply individual posteriors and account for the inverse weighting by some power of the prior which we take from Table 1.

This does however present one issue, that since we only have access to samples from the posterior and not the posterior itself it is difficult to evaluate eq. 5. We circumvent this by approximating each single star posterior by a Gaussian and fit for the parameters of mean and covariance. With this it is straight forward to evaluate eq. 5 analytically. In fact, under our assumption the combined posterior is a product between the prior and the product of multiple Gaussians for the individual star posteriors. The latter product is also a Gaussian with mean $\mu'$ and variance $\sigma'^2$:

$$\mu' = \frac{\sum_{i=1}^{N_{stars}} \frac{\mu_i}{\sigma_i^2}}{\sum_{i=1}^{N_{stars}} \frac{1}{\sigma_i^2}} \quad (6)$$

$$\sigma'^2 = \frac{1}{\sum_{i=1}^{N_{stars}} \frac{1}{\sigma_i^2}} \quad (7)$$

Further, in our case the prior for the galactic parameters $\Lambda$ is Gaussian as well. Therefore the resulting factorized posterior from eq. 5 is again a Gaussian and can be expressed with mean $\mu$ and variance $\sigma$ as:

$$\mu = \frac{\frac{\mu'}{\sigma'^2} - \frac{(1-N)\mu_{prior}}{\sigma_{prior}^2}}{\frac{1}{\sigma'^2} - \frac{(1-N)}{\sigma_{prior}^2}} \quad (8)$$

$$\sigma^2 = \frac{1}{\frac{1}{\sigma'^2} - \frac{(1-N)}{\sigma_{prior}^2}} \quad (9)$$

Given the tiny and simple neural network that represents our NPE, we note that the above assumption of Gaussianity in each of the single star posteriors not expected to notably increase our pipeline's error. We have further empirically verified that single star posteriors are indeed close to Gaussian. However, in future work we plan to alleviate this simplification and directly approximate the joint posterior of a multi-star inference.

## 3. Results

We use our SBI method described in the previous section to infer the global galactic parameters $\boldsymbol{\Lambda} = \{\alpha_{\text{IMF}}, \log_{10}(N_{\text{Ia}})\}$. In order to demonstrate the performance and robustness of our methods we use three mock data-sets:

1. Mock observations drawn from CHEMPY from the same yield set as the training data for the neural network emulator. With this we ensure to test our training strategy and the performance without any systematic distribution shifts.
2. CHEMPY mock data using a different yield set to test for potential biases through model misspecification in our SBI setup. (see 3)
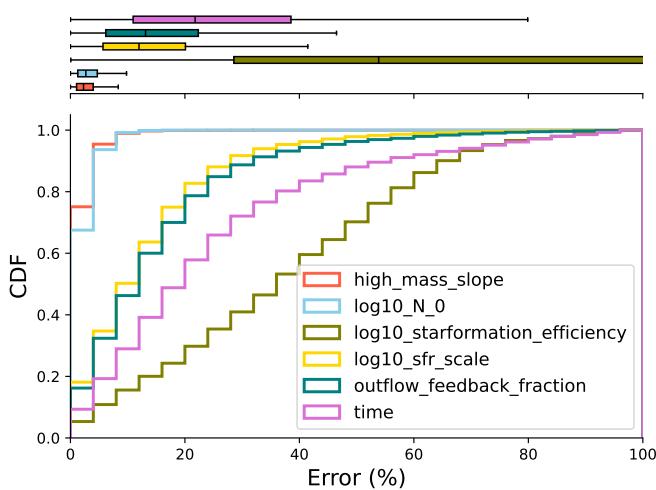
**Fig. 3.** Absolute percentage error of the neural posterior density estimate for a single star. Different colored histograms show the full error distribution for all 6 parameters of interest with the median values highlighted by the vertical dashed lines. The box plots show again the first and third quantiles with the median represented by a vertical line and the whiskers extending from the box to the farthest data point lying within 1.5× the inter-quartile range from the box. The global parameters of main interest for this work are shown by the light blue and red histogram.

3. Simulated data from a large-scale hydrodynamical simulation taken from the IllustrisTNG suite (Pillepich et al. 2018b) but with the same yield set as our `CHEMPY` training data to ensure that we recover the correct parameters even for models with a completely different treatment of ISM physics.

### 3.1. CHEMPY Mock Observational Data

Our analysis uses mock observations drawn from the neural network emulator with fixed values of the global galactic parameters $\alpha_{\mathrm{IMF}} = -2.3$ and $\log_{10}(N_{\mathrm{Ia}}) = -2.89$ and varying local parameters $\Theta_i$ using priors from Tab. 1. Additionally we draw $T_i$ uniformly in the range [2, 12.8] Gyr to minimize overlap with the neural network training birth-time limits when observational uncertainties are included. Each set of parameters is passed to our `CHEMPY` emulator, producing eight true chemical abundances. In order to fully represent observational data, we augment this data with observational uncertainties by adding a Gaussian error of 0.05 dex for the abundances representative for typical APOGEE data (Majewski et al. 2016).

For each individual observation consisting of the abundance measurements of a single star we sample the posterior estimate for all six parameters $\{\Lambda, \Theta_i, T_i\}$ with 1,000 points. This takes around $0.3s$ for each star, making it extremely fast to infer the parameters of a large amount of stars. Our method takes around 10 seconds to build a posterior for all six parameters for a dataset size of 1,000 stars (each time sampling the single star posterior with 1,000 points and fitting for the Gaussian parameters). Our combined runtime for the NPE training plus sampling is then 610 seconds on an Apple M1 Max for 1,000 stars hence our methods making it more than 1240 times faster than current HMC methods which take around 42h for only 200 stars (Philcox & Rybizki 2019). Note though, that since our approach is amortized and we do not need to retrain our NPE model each time we want to make an inference, we in fact achieve an inference speed-up of a factor of $\sim 75,600$. Hence, shorter computing times make it feasible to use orders of magnitudes more observations.

**Validation through absolute percentage error:** We start evaluating our method by using `CHEMPY` to produce a mock observational dataset to ensure no systematical shift between training and testing data in terms of physics parameters.

To evaluate our method quantitatively, we compare the posterior mean to the ground truth value for each observation and calculate the absolute percentage error (APE, see Fig. 3). For a single star observation, our NPE has an APE of $9.3^{+17.1}_{-6.3}$ % when looking at all 6 parameters of interest. If we restrict ourselves to only the two global parameters $\Lambda$, our NPE achieves an APE of $2.3^{+1.7}_{-1.2}$ % as shown in Fig. 3. There we can also see, that the accuracy for an individual star of the NPE is not particularly high.

However, currently our NPE network is not particularly good at estimating ages from abundances alone. On average our inference for ages results in an APE of $\sim 21\%$ which is slightly larger than the observational noise of 20% that we add during the mock up of our data but well in agreement with current uncertainties of stellar age inference that range from 15% to 30% for turn-off stars and giants respectively.

We accompany the APE analysis by showing the full posterior inference results for a single star in Fig. 4. This figure shows that the SBI approach is well able to infer correct parameters and their cross-correlations. In partciular, there is a strong correlation between the two global parameters as seen in the top left corner as well as for time (or stellar age) and all other five parameters as visible from the bottom row.

Finally, as already mentioned in the method section for the global parameters $\Lambda$ we can boost the accuracy by combining the inference for many stars of the same galaxy.

### 3.2. Inference on mock data from CHEMPY with TNG yield set

Combining the inference on multiple observed stars gives us higher accuracy and precision of the global galactic parameters $\Lambda$. We therefore perform inference using a range of stars $n \in [1, 1000]$ to show how inference accuracy increases with number of observations (see Fig. 5). We find that in the limit of less than $\sim 100$ stars SBI shows a larger uncertainties than the HMC results of Philcox & Rybizki (2019). But when using more than a few hundred stars the accuracy and precision is superior compared to HMC. In particular, after using a few tens of stars, the NPE estimate is already less biased than the HMC results. Finally, given our computational advantage we will be able to use orders of magnitude more stars for our inference. This is particularly important since sample variances play a large role when using small samples of stars as also noted in Philcox & Rybizki (2019).

In Fig. 6 we show the joint posterior for $\alpha_{IMF}$ and $\log_{10}(N_{Ia})$ for our inference. The red star indicates the ground truth value, the black dot shows our posterior mean and white contours indicate $1 - 3\sigma$ levels. Using a sample of 1,000 stars we infer $\alpha_{\mathrm{IMF}} = -2.299 \pm 0.002$ and $\log_{10}(N_{\mathrm{Ia}}) = -2.890 \pm 0.003$ which deviates less than $\sim 0.04\%$ from the ground truth value. We have further checked the accuracy of our inference for a vastly different mock observational dataset created with shifted parameters of $\alpha_{IMF} = -2.1$ and $\log_{10}(N_{Ia}) = -3$ and found that also in this case our model is well able to lead to correct inferences (see sec. D in the appendix).

Note that our analysis is in principle also able to infer the local parameters $\Theta_i$ and $T_i$. This would allow us to estimate/infer stellar ages as well. But note, as discussed above at the end of sub-section 3.1 our NPE currently is not well calibrated to estimate stellar ages accurately enough.

In summary, our SBI pipeline is quite capable of correctly and precisely inferring global parameters of chemical enrichment
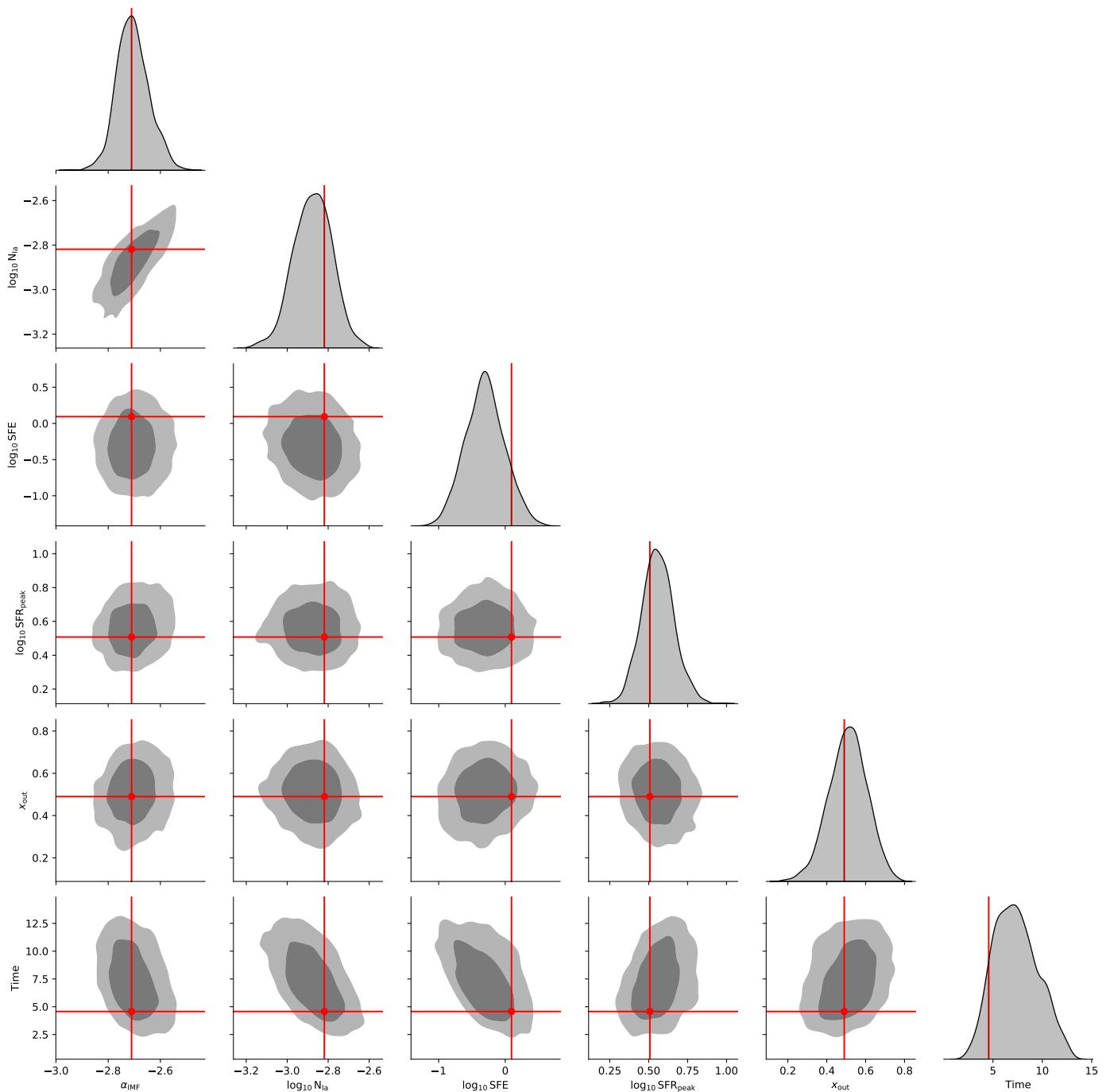
**Fig. 4.** Corner plot of the posteriors for all six parameters for a single star from the validation set. The gray contours show a kde-estimate of the posterior from our SBI inference and the red dot and lines show the ground truth parameter values. Gray histograms on the diagonal show a kde estimate of the marginals.

models from stellar abundance alone when using the same physical model and yield tables as during training. Next, we will check what happens if the training data is generated with a different yield set to that use at inference time.

### 3.3. Inference on mock CHEMPY data with incorrect yield set

There is an extensive discussion in the literature about stellar nucleosynthesis with various different yield sets proposed (see e.g. the discussion in Rybizki et al. 2017). In fact, all tabulated yield sets currently differ from reality and during an application of our inference pipeline it will not be clear which tabulated yield set most closely matches reality and hence which should be used. In order to investigate how sensitive our method is to model misspecification by using an incorrect yield set, we create another set of mock data using CHEMPY with a different yield set (Tab. 3) than during training of our NPE. For better cross-comparison we decided to use the same alternative yield sets as presented in Tab. 5 of (Philcox & Rybizki 2019).

By choosing this set of yields we have made sure that contributions to all three processes differ by $O(10\%)$. For more details see also Sec. 6.2 of Philcox & Rybizki (2019).

Our mock data generation and inference then follows the one of Sec. 3.2. This means we apply our NPE that we trained
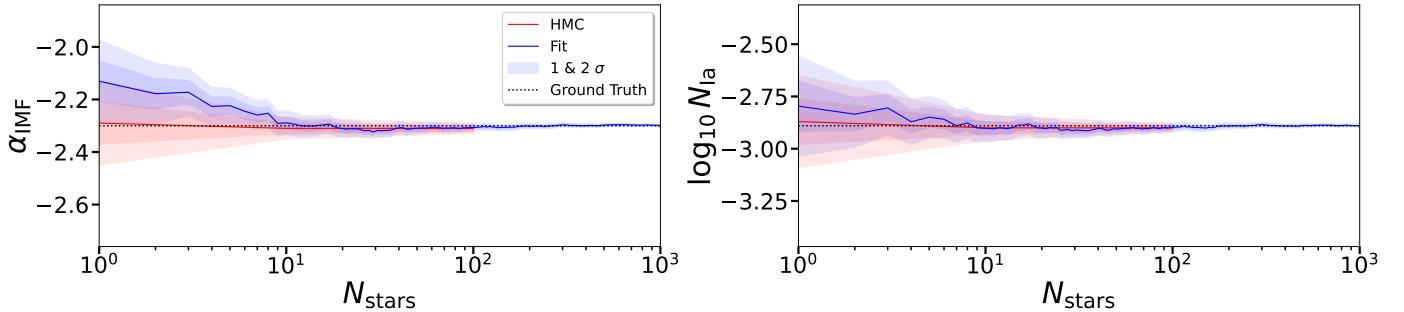
**Fig. 5.** Accuracy of inferred global galactic parameters $\alpha_{IMF}$ and $\log_{10}(N_{Ia})$ as a function of number of observed stars. We compare our SBI results (blue line) to the inferred values using HMC (red line) as done by Philcox & Rybizki (2019) and the ground truth values (black dashed line) for various test cases as described in Sec. For the SBI analysis we show $1\sigma$ and $2\sigma$ contours while HMC results only show $1\sigma$ statistical uncertainties as reported in Tab. 3 of Philcox & Rybizki (2019) (blue/red shaded regions). See Sec. 3.2 for a full description.



**Fig. 6.** Joint posterior for global galactic parameters $\alpha_{IMF}$ and $\log_{10}(N_{Ia})$ of 100 stars. The ground truth value is shown by the red star, the posterior mean of the SBI inference is shown with the black dot and white ellipses show the $1-3\sigma$ contours of our inference. The gray ellipses show the results for HMC inference performed by Philcox & Rybizki (2019). See Sec. 3.2 for a full description.

**Table 3.** Alternative nucleosynthetic yield tables used for model mis-specification tests.

| Type | Yield Table |
|------|-------------|
| SN Ia | Thielemann et al. (2003) |
| SN II | Nomoto et al. (2013) |
| AGB | Karakas & Lugaro (2016) |

on `CHEMPY` stellar abundances simulated with the Ilustris TNG yield set to a sets of stellar abundances simulated with `CHEMPY` but using the alternative yield sets mentioned in Tab. 3. This effectively probes the effect of model misspecification on the inference results.

In Fig. 7 we show our inference results (blue) for this setup for varying number of stars and compare them against HMC results (red) from Philcox & Rybizki (2019). We see that the SBI

results are similar biased as the HMC results. In fact with 100 stars HMC inferences are about 4 and 3 $\sigma$ away from the ground truth value for $\alpha_{IMF}$ and $\log_{10}(N_{Ia})$, respectively. Again, the joint posterior is shown in Fig. 9 and shows that while individually parameter inferences look good, jointly taken the SBI results are on the edge of being $3\sigma$ away from the ground truth.

Nevertheless, the performance of our inference is still very good. We see that an increasing number of observations helps to decrease the models uncertainty just as before. However, we also note that our inference is now slightly biased as the observational data does not match the training data. Looking at Fig. 9 we see that our inference is inconsistent with the ground truth within several $\sigma$ levels. Comparing this to the HMC results (red band in Fig. 7), we see that SBI is performing similarly bad as HMC when the model is misspecified. Part of the problem here is that through our assumption of a factorized posterior, we decrease the inference uncertainty as we increase the number of observations. In future we will improve upon this through different model architectures. In summary, the drastically reduced compute times offer a key advantage of our SBI method compared to more standard approaches such as HMC. Nevertheless, in the limit of large star counts our current approach becomes over-confident and model uncertainities are underestimated.

In an accompanying paper we more closely investigate measures of model misspecification and inference of best fitting models next to just parameter inference.

### 3.4. Inference on mock data from a IllustrisTNG simulated galaxy

As a GCE code, `CHEMPY` is a one-zone model with simplified ISM physics that only approximately describes star formation and feedback as well as metal mixing in the ISM. In the parametrization of `CHEMPY` as used here, we can assign each star to its own ISM environment, but we cannot exchange gas between environments and do not model sudden star formation or infall events. Hence, this section is dedicated to investigating whether this significantly biases our inference of the SSP parameters (noting that results from Weinberg et al. (2019) justify the treatment of ISM parameters as latent variables).

In order to explore what effect this simplified treatment of ISM physics has on the inference we now turn to a more complex model of the formation and chemical enrichment of a Milky Way-type galaxy taken from the IllustrisTNG simulations. Note, that by now also the NIHAO simulations (Wang et al. 2015; Buck 2020; Buck et al. 2020) have implemented `CHEMPY` supported yield tables including the TNG yield set (Buck et al. 2021) and
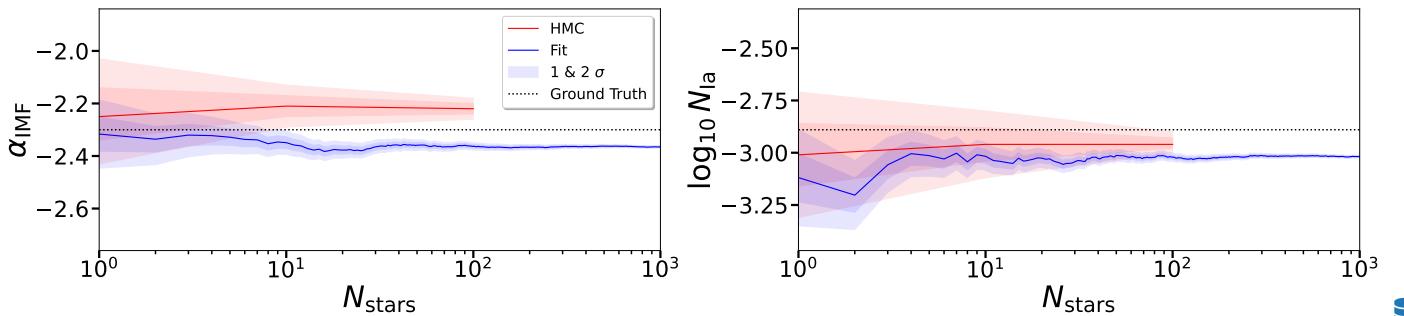
**Fig. 7.** Same as Fig. 5 but for the mock data created with a different yield set than the training data. See Sec. 3.3 for a full description.
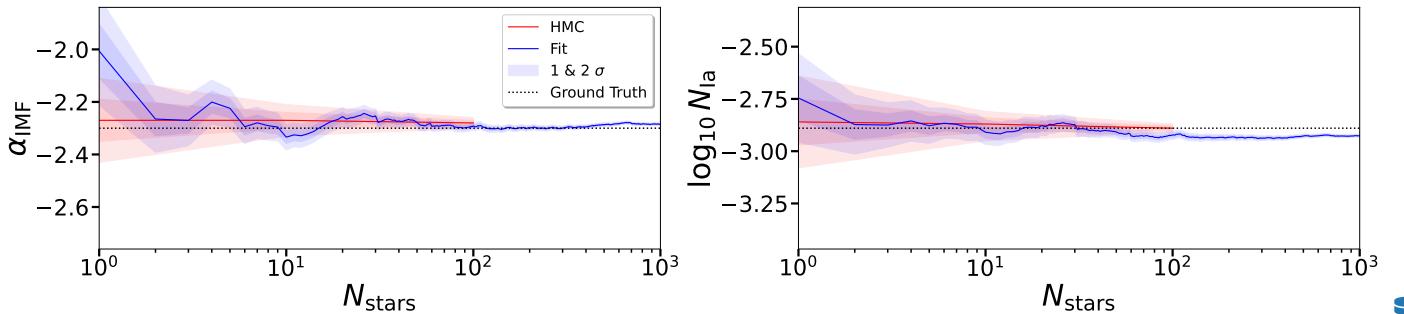


**Fig. 8.** Same as Fig. 5 but for the mock data taken from an IllustrisTNG Milky Way-like galaxy. See Sec. 3.4 for a full description.
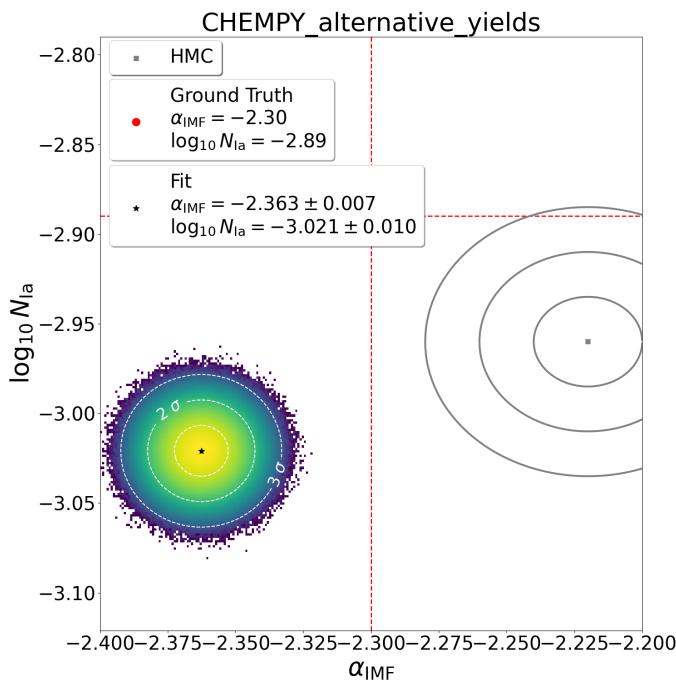


**Fig. 9.** Same as Fig. 6 but for the mock data created with a different yield set than the training data. See Sec. 3.3 for a full description.
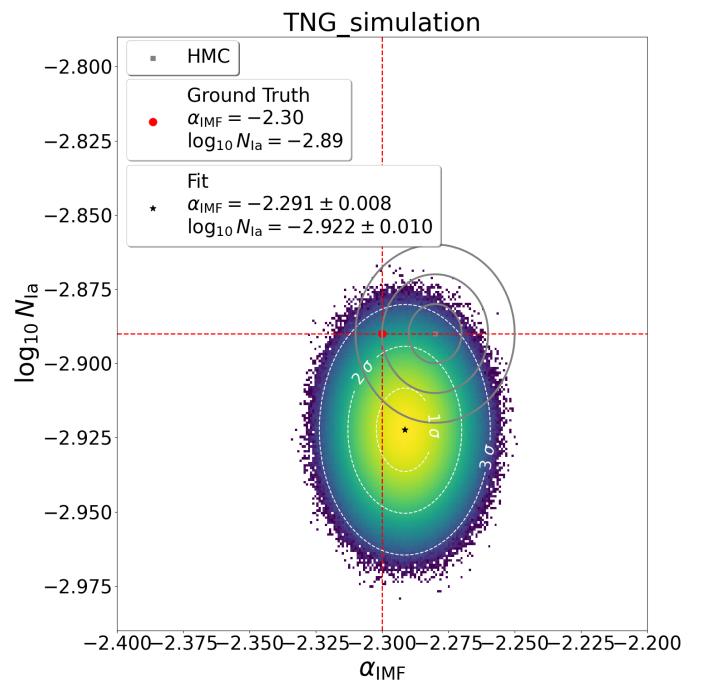
**Fig. 10.** Same as Fig. 10 but for the mock data taken from an IllustrisTNG Milky Way-like galaxy. See Sec. 3.4 for a full description.

hence would make for a nice dataset for our analysis. However, we have decided to use the exact same galaxy as in Philcox & Rybizki (2019) for better comparison of our results.

In detail, we use a single galaxy from the $z = 0$ snapshot of the highest-resolution TNG100-1 simulation. We choose a subhalo (index 523071) with mass close to $10^{12}\,\mathrm{M}_\odot$ to select a Milky Way-like galaxy. From this, we extract a set of 1,000 random 'stellar particles' from a total of $\sim 40,000$. Each star particle has a mass of $\sim 1.4 \times 10^6\,\mathrm{M}_\odot$ (Nelson et al. 2019). These act as proxies

for stellar environments, giving the elemental mass fractions, $\{d_i^j\}$, and cosmological scale factor, $a_i$, at the time of stellar birth. Mass fractions are converted to [X/Fe] abundance ratios using Asplund et al. (2009) solar abundances as in `CHEMPY`, with the scale-factor ($a_i$) to birth-time ($T_i$) conversion performed using `astropy` (Astropy Collaboration et al. 2013; Price-Whelan et al. 2018),[8] assuming a ΛCDM cosmology with Planck Collaboration

---

[8] http://www.astropy.org

et al. (2016) parameters, as in TNG (Pillepich et al. 2018a).[9] Observational errors are incorporated as above, giving a full dataset that is identical in structure to the CHEMPY mock data. For more details and a plot of the [Mg/Fe] vs. [Fe/H] plane for this galaxy see Sec. 6.3 and Fig. 5 of Philcox & Rybizki (2019).

We note that this TNG galaxy was deliberately chosen by Philcox & Rybizki (2019) to have both a high-$\alpha$ and low-$\alpha$ chemical evolution sequence to test their inference on a mock galaxy with Milky Way-like properties. There is still some debate on the exact formation of this bimodality but it is generally attributed to gas-rich mergers and different modes of star formation (e.g. Grand et al. 2018; Mackereth et al. 2018; Clarke et al. 2019; Buck 2020; Buck et al. 2023). Similarly, in chemo-dynamical models, Milky Way-like bimodalities can also be achieved by a combination of radial migration and selection effects without the need for mergers or starbursts (e.g. Schönrich & Binney 2009; Minchev et al. 2013; Andrews et al. 2017).

We show our inference results for the TNG data set in Fig. 8 and Fig. 10, respectively. Again, we find that inference becomes better with increasing number of stars. Despite the drastic difference in chemical enrichment model between training and testing data, our SBI pipeline is impressively capable of inferring the correct posterior values. SBI results for $\log_{10}(N_{Ia})$ are almost perfect while $\alpha_{IMF}$ is slightly biased high. This is very similar to the results of HMC where inference for $\alpha_{IMF}$ is biased and results for $\log_{10}(N_{Ia})$ are more in agreement. Hence, our SBI inference is on par with the HMC results and given their extreme computational advantage they actually supersede HMC.

Looking at the joint posteriors in Figs. 6, 9 and 10 we see that also in this case SBI recovers the true parameters only with a few $\sigma$ bias similar to HMC. Again, part of this problem is our assumption of a factorized posterior that underestimated model uncertainties as for the alternative yield case.

## 4. Discussion

Our study demonstrates that simulation-based inference (SBI) provides a powerful and efficient alternative to conventional methods such as Hamiltonian Monte Carlo (HMC) for inferring global galactic parameters. By leveraging neural density estimators and neural network emulators, we achieved remarkable computational efficiency without compromising precision or accuracy.

The key assumption underpinning our methodology – that individual stars are independently sampled from their respective stellar environments – is foundational for tractability. However, this assumption warrants closer examination. In reality, stellar abundances are correlated due to shared star formation histories, cumulative enrichment, and dynamical interactions. Future efforts should explore methodologies capable of incorporating such correlations to further refine the accuracy of SBI, potentially through hierarchical modeling or graph-based methods. In particular this assumption leads to over-confident models in the limit of large observational samples that drastically underestimate model uncertainties. One very promising avenue to circumvent this problem will be compositional score models (Geffner et al. 2022) which we will explore in future work.

An exciting avenue for future research involves expanding the SBI method presented here to be able to cope with missing data points by switching to a transformer model architecture.

Also, with slight modifications, SBI could be used to directly infer empirical nucleosynthetic yields, which remain a major uncertainty in GCE models. In the current framework, CHEMPY relies on tabulated yields from theoretical studies, which may not fully represent the complex processes driving stellar enrichment. Adapting SBI to simultaneously infer galactic parameters and refine empirical yield tables would require changes to the simulator. Specifically, CHEMPY would need to incorporate parameterized yield modifications as part of its input space, allowing for flexible adjustments to enrichment rates during inference. This would also necessitate larger training datasets and enhanced validation techniques to ensure convergence. Such an approach could provide a unified framework for calibrating galactic models directly against observational data. Research in this direction is currently performed and will be part of a future paper.

In the meantime, SBI can also be used to perform model comparison (e.g. Spurio Mancini et al. 2023; Zhou et al. 2024) which can be used to determine which of the tabulated yield sets best match observational constraints.

Compared to previous HMC-based studies (e.g. Philcox & Rybizki 2019), our results highlight SBI's resilience to certain types of model misspecification, such as mismatched yield tables. The robustness of SBI in these cases stems from its ability to approximate complex posterior distributions efficiently. Notably, SBI retained high accuracy in scenarios where HMC struggled, particularly for the high-mass slope of the IMF ($\alpha_{ISM}$). This suggests SBI's potential for real-world applications, where the underlying models may deviate from observational data.

The applicability of our method extends to surveys with more observed elements especially neutron capture elements such as GALAH DR4 (Buder et al. 2024) and future spectroscopic surveys, such as those planned by the 4MOST (de Jong et al. 2014) or WEAVE (Dalton et al. 2018) consortia, which will provide orders of magnitude more data than current datasets. Our findings indicate that SBI can seamlessly scale to such large datasets, offering significant advantages in terms of speed and computational cost.

As anticipated in Section 2, the speed of our NN emulator would allow us to perform SNPE, giving up the amortized property in favor of a stronger constraining power, since the Sequential techniques have been shown empirically to outperform the respective amortized version (Ho et al. 2024). This method trains the NDE on a fraction of the initial simulations budget, retrieving a first estimate of the posterior, and at inference time simulate new observations to train on in regions of high posterior density, obtained from the first estimate. In this way we can obtain posteriors that are deliberately optimized for a singular data point. We have decided to leave this approach for future work because the results are quite promising and the amortized property can be crucial for the scalability.

Nevertheless, limitations remain. The simplified physics of the CHEMPY model, while advantageous for computational efficiency, omits the complexities of dynamical gas flows, feedback, and metal mixing present in cosmological simulations. While our tests on IllustrisTNG data affirm the robustness of SBI, integrating more sophisticated models into the inference pipeline represents an exciting avenue for future work. On this line of reasoning we refer also to the discussion in Philcox & Rybizki (2019).

## 5. Summary and conclusions

This study introduces simulation-based inference (SBI) as an innovative framework for constraining galactic parameters us-

---

[9] Note, as for the CHEMPY mock data, we exclude any particles with $T_i \notin [2, 12.8]$ Gyr to ensure that the true times are well separated from our training age limits, avoiding neural network errors. This removes $\sim 5\%$ of the stars.

ing stellar chemical abundances. By training neural posterior estimators on forward simulations from the CHEMPY model, we achieved precise and accurate inferences for two key parameters: the high-mass slope of the initial mass function ($\alpha_{\mathrm{IMF}}$) and the normalization of Type Ia supernova rates ($\log_{10}(N_{\mathrm{Ia}})$).

Our results underscore the transformative advantages of SBI over traditional methods like Hamiltonian Monte Carlo (HMC), marking a paradigm shift in galactic parameter inference:

- Orders-of-magnitude speed-up: SBI dramatically reduces computational requirements, achieving runtime improvements exceeding 75,000-fold compared to HMC. For instance, while HMC requires ∼ 42 hours to infer parameters from just 200 stars, SBI completes inference on thousands of stars in mere minutes. This efficiency makes SBI uniquely suited for analyzing the massive datasets expected from next-generation spectroscopic surveys.
- Scalability: SBI's amortized nature allows it to scale seamlessly with the size of the dataset. By training a neural posterior estimator once, the method can be applied repeatedly at virtually no additional computational cost. This scalability is essential for leveraging the millions of stars that future surveys like 4MOST (de Jong et al. 2014) and WEAVE (Dalton et al. 2018) will provide, enabling precise population-level inferences.
- Robustness to model misspecifications: Unlike HMC, which shows significant biases when faced with discrepancies between model assumptions and data, SBI demonstrates remarkable robustness. Even under conditions of mismatched yield tables or data generated from hydrodynamical simulations, SBI provides accurate and reliable results. This robustness ensures SBI's applicability in real-world scenarios where exact model fidelity cannot be guaranteed.

In addition to its immediate advantages, SBI lays the foundation for future advancements in galactic modeling. Its flexibility can enable the direct inference of empirical nucleosynthetic yields and facilitate integration with more complex galaxy formation models. These enhancements will further solidify SBI as a cornerstone method in galactic archaeology.

We note though, that the SBI pipeline presented here is not perfect and suffers from simplified assumptions that have been made in order to make the problem tractable. In particular, the assumption that the posterior factorizes is a strong assumption that leads to over-confidence of the models at large observational samples. This is certainly something that needs to be improved in future iterations and methodological work in this direction is currently undergoing.

In conclusion, SBI represents a breakthrough in simulation-based analysis, delivering unparalleled speed, precision, and scalability. By overcoming the computational limitations of traditional techniques like HMC, SBI paves the way for extracting deeper insights into the chemical and dynamical evolution of galaxies in the era of massive spectroscopic surveys.

# References

Abdurro'uf, Accetta, K., Aerts, C., et al. 2022, ApJS, 259, 35
Agertz, O., Renaud, F., Feltzing, S., et al. 2021, MNRAS, 503, 5826
Andrews, B. H., Weinberg, D. H., Schönrich, R., & Johnson, J. A. 2017, ApJ, 835, 224
Asplund, M., Grevesse, N., Sauval, A. J., & Scott, P. 2009, ARA&A, 47, 481
Astropy Collaboration, Robitaille, T. P., Tollerud, E. J., et al. 2013, aap, 558, A33
Bigiel, F., Leroy, A., Walter, F., et al. 2008, AJ, 136, 2846
Buck, T. 2020, MNRAS, 491, 5435
Buck, T., Obreja, A., Macciò, A. V., et al. 2020, MNRAS, 491, 3461
Buck, T., Obreja, A., Ratcliffe, B., et al. 2023, MNRAS, 523, 1565
Buck, T., Rybizki, J., Buder, S., et al. 2021, MNRAS, 508, 3365
Buder, S., Kos, J., Wang, E. X., et al. 2024, arXiv e-prints, arXiv:2409.19858
Buder, S., Sharma, S., Kos, J., et al. 2021, MNRAS, 506, 150
Chabrier, G. 2003, PASP, 115, 763
Chabrier, G., Hennebelle, P., & Charlot, S. 2014, ApJ, 796, 75
Clarke, A. J., Debattista, V. P., Nidever, D. L., et al. 2019, MNRAS, 484, 3476
Clauwens, B., Schaye, J., & Franx, M. 2016, MNRAS, 462, 2832
Cook, S. R., Gelman, A., & Rubin, D. B. 2006, Journal of Computational and Graphical Statistics, 15, 675
Côté, B. & Ritter, C. 2018, OMEGA: One-zone Model for the Evolution of GAlaxies, Astrophysics Source Code Library, record ascl:1806.018
Côté, B., Ritter, C., O'Shea, B. W., et al. 2016, ApJ, 824, 82
Cranmer, K., Brehmer, J., & Louppe, G. 2020, Proceedings of the National Academy of Sciences, 117, 30055
Dalton, G., Trager, S., Abrams, D. C., et al. 2018, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 10702, Ground-based and Airborne Instrumentation for Astronomy VII, ed. C. J. Evans, L. Simard, & H. Takami, 107021B
de Jong, R. S., Barden, S., Bellido-Tirado, O., et al. 2014, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 9147, Ground-based and Airborne Instrumentation for Astronomy V, ed. S. K. Ramsay, I. S. McLean, & H. Takami, 91470M
Defazio, A., Yang, X. A., Mehta, H., et al. 2024, arXiv e-prints, arXiv:2405.15682
Doherty, C. L., Gil-Pons, P., Lau, H. H. B., Lattanzio, J. C., & Siess, L. 2014, MNRAS, 437, 195
Durkan, C., Bekasov, A., Murray, I., & Papamakarios, G. 2019, Neural Spline Flows
Fishlock, C. K., Karakas, A. I., Lugaro, M., & Yong, D. 2014, ApJ, 797, 44
Font, A. S., McCarthy, I. G., Poole-Mckenzie, R., et al. 2020, MNRAS, 498, 1765
Geffner, T., Papamakarios, G., & Mnih, A. 2022, arXiv e-prints, arXiv:2209.14249
Gloeckler, M., Deistler, M., Weilbach, C., Wood, F., & Macke, J. H. 2024, ArXiv, abs/2404.09636
Grand, R. J. J., Bustamante, S., Gómez, F. A., et al. 2018, MNRAS, 474, 3629
Griffith, E., Johnson, J. A., & Weinberg, D. H. 2019, arXiv e-prints, arXiv:1908.06113
Gutcke, T. A. & Springel, V. 2019, MNRAS, 482, 118
Ho, M., Bartlett, D. J., Chartier, N., et al. 2024, The Open Journal of Astrophysics, 7, 54
Hopkins, P. F., Wetzel, A., Kereš, D., et al. 2018, MNRAS, 480, 800
Jiménez, N., Tissera, P. B., & Matteucci, F. 2015, ApJ, 810, 137
Johnson, J. W. & Weinberg, D. H. 2020, MNRAS, 498, 1364
Karakas, A. I. 2010, MNRAS, 403, 1413
Karakas, A. I. & Lugaro, M. 2016, ApJ, 825, 26
Kobayashi, C., Umeda, H., Nomoto, K., Tominaga, N., & Ohkubo, T. 2006, ApJ, 653, 1145
Lemos, P., Coogan, A., Hezaveh, Y., & Perreault-Levasseur, L. 2023, 40th International Conference on Machine Learning, 202, 19256
Luger, R., Bedell, M., Foreman-Mackey, D., et al. 2021, arXiv e-prints, arXiv:2110.06271
Mackereth, J. T., Crain, R. A., Schiavon, R. P., et al. 2018, MNRAS, 477, 5072
Majewski, S. R., APOGEE Team, & APOGEE-2 Team. 2016, Astronomische Nachrichten, 337, 863
Maoz, D. & Mannucci, F. 2012, PASA, 29, 447
Maoz, D., Mannucci, F., & Brandt, T. D. 2012, MNRAS, 426, 3282
Maoz, D., Sharon, K., & Gal-Yam, A. 2010, ApJ, 722, 1879
Martín-Navarro, I., Lyubenova, M., van de Ven, G., et al. 2019, A&A, 626, A124
Minchev, I., Chiappini, C., & Martig, M. 2013, A&A, 558, A9
Mollá, M., Cavichia, O., Gavilán, M., & Gibson, B. K. 2015, MNRAS, 451, 3693
Nelson, D., Springel, V., Pillepich, A., et al. 2019, Computational Astrophysics and Cosmology, 6, 2
Ness, M. K., Johnston, K. V., Blancato, K., et al. 2019, arXiv e-prints, arXiv:1907.10606
Nomoto, K., Iwamoto, K., Nakasato, N., et al. 1997, Nucl. Phys. A, 621, 467
Nomoto, K., Kobayashi, C., & Tominaga, N. 2013, ARA&A, 51, 457
Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., & Lakshminarayanan, B. 2021, Normalizing Flows for Probabilistic Modeling and Inference
Papamakarios, G., Pavlakou, T., & Murray, I. 2018, Masked Autoregressive Flow for Density Estimation
Philcox, O., Rybizki, J., & Gutcke, T. A. 2018, The Astrophysical Journal, 861, 40
Philcox, O., Rybizki, J., & Gutcke, T. A. 2018, ApJ, 861, 40
Philcox, O. H. E. & Rybizki, J. 2019, The Astrophysical Journal, 887, 9
Pillepich, A., Nelson, D., Hernquist, L., et al. 2018a, MNRAS, 475, 648
Pillepich, A., Springel, V., Nelson, D., et al. 2018b, MNRAS, 473, 4077

Pillepich, A., Springel, V., Nelson, D., et al. 2018c, MNRAS, 473, 4077

Planck Collaboration, Ade, P. A. R., Aghanim, N., et al. 2016, A&A, 594, A13

Portinari, L., Chiosi, C., & Bressan, A. 1998, A&A, 334, 505

Price-Whelan, A. M., Sip'ocz, B. M., G"unther, H. M., et al. 2018, aj, 156, 123

Romano, D., Chiappini, C., Matteucci, F., & Tosi, M. 2005, A&A, 430, 491

Rybizki, J. & Just, A. 2015, MNRAS, 447, 3880

Rybizki, J., Just, A., & Rix, H.-W. 2017, A&A, 605, A59

Sawala, T., Frenk, C. S., Fattahi, A., et al. 2016, MNRAS, 457, 1931

Schönrich, R. & Binney, J. 2009, MNRAS, 396, 203

Spurio Mancini, A., Docherty, M. M., Price, M. A., & McEwen, J. D. 2023, RAS Techniques and Instruments, 2, 710

Talbot, Jr., R. J. & Arnett, W. D. 1971, ApJ, 170, 409

Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A. 2018, arXiv e-prints, arXiv:1804.06788

Thielemann, F. K., Argast, D., Brachwitz, F., et al. 2003, Nucl. Phys. A, 718, 139

van Dokkum, P. G., Leja, J., Nelson, E. J., et al. 2013, ApJ, 771, L35

Vincenzo, F., Matteucci, F., Recchi, S., et al. 2015, MNRAS, 449, 1327

Viterbo, G. & Buck, T. 2024, arXiv e-prints, arXiv:2411.17269

Wang, K., Carrillo, A., Ness, M. K., & Buck, T. 2024, MNRAS, 527, 321

Wang, L., Dutton, A. A., Stinson, G. S., et al. 2015, MNRAS, 454, 83

Weinberg, D. H., Holtzman, J. A., Hasselquist, S., et al. 2019, ApJ, 874, 102

Weisz, D. R., Johnson, L. C., Foreman-Mackey, D., et al. 2015, ApJ, 806, 198

Zeghal, J., Lanusse, F., Boucaud, A., Remy, B., & Aubourg, E. 2022, Neural Posterior Estimation with Differentiable Simulators

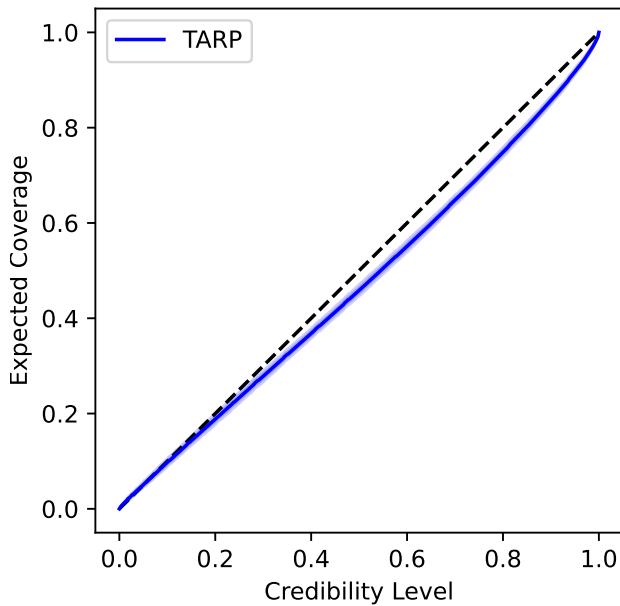Zhou, L., Radev, S. T., Oliver, W. H., et al. 2024, arXiv e-prints, arXiv:2410.10606

**Fig. A.1.** TARP plot showing the expected coverage probability vs. the credibility level $\alpha$. The dashed black 1:1-line shows an ideal calibrated posterior and the blue solid line shows the TARP value for our NPE.

## Appendix A: Code and Data Availability

To facilitate a wider community's usage and contributions, we make use of the reproducibility software *show your work!* (Luger et al. 2021), which leverages continuous integration to programmatically download the data from zenodo.org, create the figures, and compile the manuscript. Each figure caption contains two links: one to the dataset stored on zenodo used in the corresponding figure, and the other to the script used to make the figure (at the commit corresponding to the current build of the manuscript). The git repository associated to this study is publicly available at `https://github.com/TobiBu/sbi-chempy`, and the release v0.4.1 allows anyone to re-build the entire manuscript including rerunning all analysis. The datasets and neural network weights are stored at `https://zenodo.org/records/14925307`. The training and validation data can be found on zenodo as well under this link: `https://zenodo.org/records/14507221`.

## Appendix B: Neural Posterior Calibration

Since SBI relies in neural networks to approximate posterior densities one important point is to check that neural network hyperparameters are well chosen and that posterior estimates are trustable.

After a density estimator has been trained with simulated data to obtain a posterior, the estimator should be made subject to several diagnostic tests. This needs to be performed before being used for inference given the actual observed data. Posterior Predictive Checks provide one way to "critique" a trained estimator based on its predictive performance. Another important approach to such diagnostics is simulation-based calibration as developed by Cook et al. (2006) and Talts et al. (2018).

**Simulation-based calibration** Simulation-based calibration (SBC) provides a (qualitative) view and a quantitive measure to check, whether the variances of the posterior are balanced, i.e. it is neither over-confident nor under-confident. As such, SBC can be viewed as a necessary condition (but not sufficient) for a valid inference algorithm: If SBC checks fail, this tells you that your inference is invalid. If SBC checks pass, this is no guarantee that the posterior estimation is working.

To perform SBC, we sample some $\theta_i^o$ values from the parameter prior of the problem at hand and simulate "observations" from these parameters:

$$x_i = \text{simulator}(\theta_i^o) \tag{B.1}$$

Then we perform inference given each observation $x_i$ which produces a separate posterior $p_i(\theta|x_i)$ for each $x_i$. The key step for SBC is to generate a set of posterior samples $\{\theta\}_i$ from each posterior. We call this $\theta_i^s$, referring to $s$ samples from the posterior $p_i(\theta|x_i)$. Next, we rank the corresponding $\theta_i^o$ under this set of samples. A rank is computed by counting how many samples $\theta_i^s$ fall below their corresponding $\theta_i^o$ value (see section 4.1 in Talts et al. 2018). These ranks are then used to perform the SBC check itself.

The core idea behind SBC is two fold: (i) SBC ranks of ground truth parameters under the inferred posterior samples follow a uniform distribution (If the SBC ranks are not uniformly distributed, the posterior is not well calibrated); and (ii) samples from the data averaged posterior (ensemble of randomly chosen posterior samples given multiple distinct observations $x_o$) are distributed according to the prior.

Hence, SBC can tell us whether the SBI method applied to the problem at hand produces posteriors that have well-calibrated uncertainties, and if the posteriors have uncalibrated uncertainties, SBC surfaces what kind of systematic bias is present: negative or positive bias (shift in the mean of the predictions) or over- or under-dispersion (too large or too small variance).

In the top panel of Fig. A.2 we show the distribution of ranks (depicted in red) in each dimension. Highlighted with black lines, you see the 99% confidence interval of a uniform distribution given the number of samples provided. In plain english: for a uniform distribution, we would expect 1 out of 100 (blue) bars to lie outside the grey area. This figure shows that overall our posteriors are decently calibrated. Only for the parameter $\log_{10}(\text{SFE})$ and Time we see a slight bmiss-calibration. But most importantly for the parameters of interest here, $\log_{10}(N_{\text{Ia}})$ and $\alpha_{\text{IMF}}$ we have a well calibrated posterior.

**Tests of Accuracy with Random Points (TARP)** TARP (Lemos et al. 2023) is an alternative calibration check for the joint distribution, for which defining a rank is not straightforward. Given a test set $(\theta^*, x^*)$ and a set of reference points $\theta_r$, TARP calculates statistics for posterior calibration by - drawing posterior samples $\theta$ given each observation $x^*$ and calculating the distance $r$ between $\theta^*$ and $\theta_r$ counting for how many of the posterior samples the distance to $\theta_r$ is smaller than $r$ (see e.g. Fig. 2 in Lemos et al. 2023, for an illustration).

For each given coverage level $\alpha$, one can then calculate the corresponding average counts and check, whether they correspond to the given $\alpha$. The visualization and interpretation of TARP values is therefore similar to that of SBC. However, in contrast to SBC, TARP provides a necessary and sufficient condition for posterior accuracy, i.e., it can also detect inaccurate posterior estimators. In the middle row of Fig. A.2 we show the result for our NPE in blue in comparison to the ideal line shown
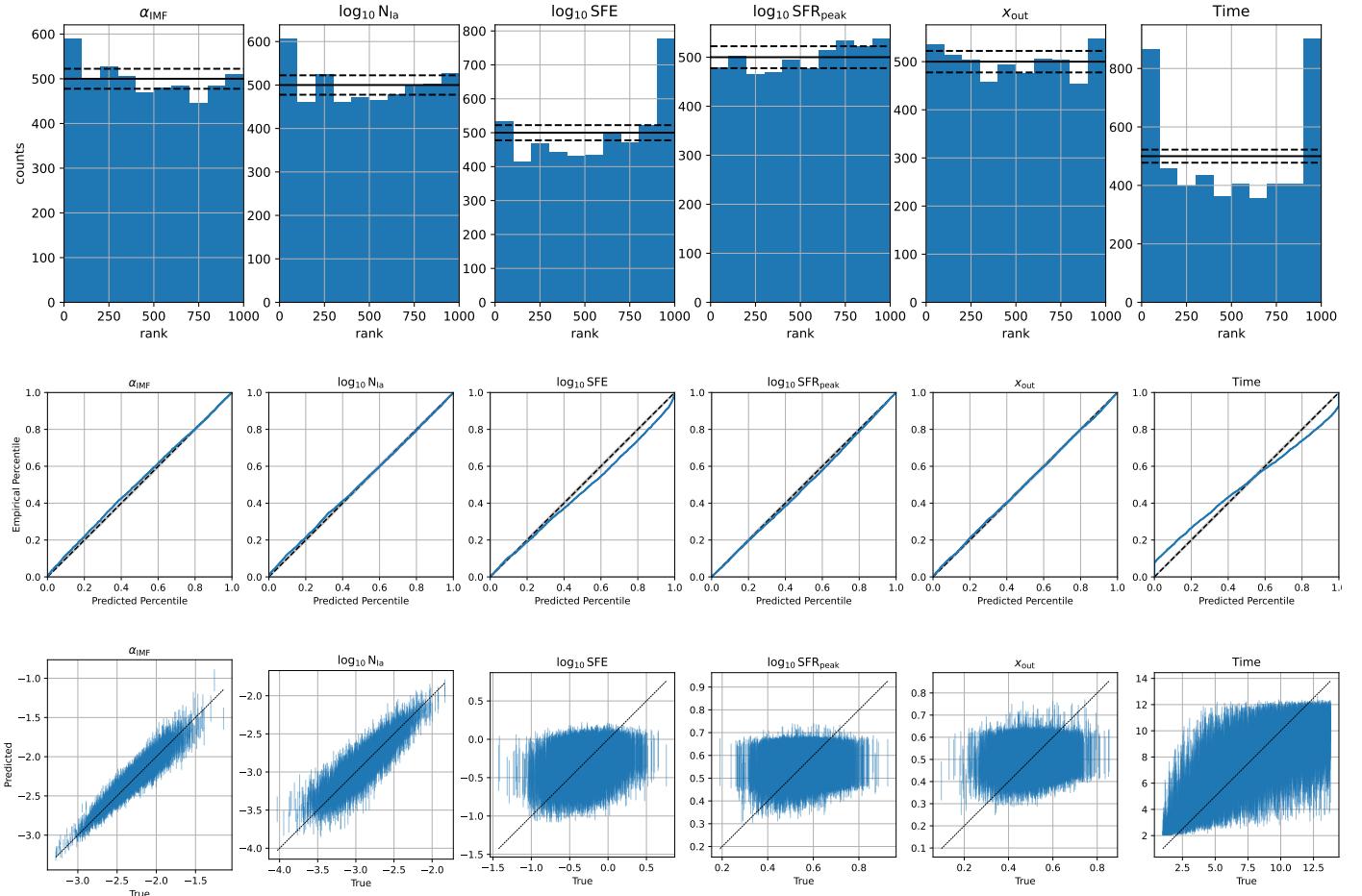
**Fig. A.2.** Posterior calibration diagnostics showing from top to bottom the SBC ranks, the TARP plots and the true vs. predicted parameter plots for each of the six parameters. Top: SBC ranks of ground truth parameters under the inferred posterior samples for each of the six parameters (red bars). The grey area shows the 99% confidence interval of a uniform distribution given the number of samples provided. Middle: TARP plot showing the expected coverage probability vs. the credibility level $\alpha$ for each of the six individual parameters in our inference. The dashed black 1:1-line shows an ideal calibrated posterior and the blue solid line shows the TARP value for our NPE. Bottom: True vs. predicted parameter plots showing the average of the posterior samples and as errorbar the standard deviation of the samples vs their ground truth parameter.

in black dashed style. This figure clearly shows that our NPE is well calibrated.

Note, however, that this property depends on the choice of reference point distribution: to obtain the full diagnostic power of TARP, one would need to sample reference points from a distribution that depends on $x$. Thus, in general, it is recommended using and interpreting TARP like SBC and complementing coverage checks with posterior predictive checks.

Finally, the bottom row of Fig. A.2 shows the predicted vs. true parameter plots where blue dots show the average of the posterior samples and errorbars show the standard deviation of the samples vs their ground truth parameter. We find that for the global parameters we recover the 1:1 relation well while for the other parameters the agreement is tilted.

## Appendix C: Correlation analysis of inference results

In order to characterize the relation between the global parameters $\alpha_{IMF}$ and $\log_{10} N_{Ia}$, we have decided to study the correlation obtained from the inference results on the whole validation set. Figure 4 display a possible positive correlation, and in order to obtain a population level statics of this results, we calculate

for each star in the validation set the 2x2 covariance matrix[10] of the posterior samples, extracting the off diagonal value. In Figure C.1, the histogram of the correlation value is shown, and all the inference results suggest a positive correlation. We kept only the left sided 99 percentile for graphical reasons, but also those removed outliers are in agreement with the conclusion. We checked the independence of this results from the accuracy of the inference with the central and right plots of Fig. C.1.

## Appendix D: Additional inference results

---

[10] We have used only the global parameters components of the sample to calculate the covariance matrix.
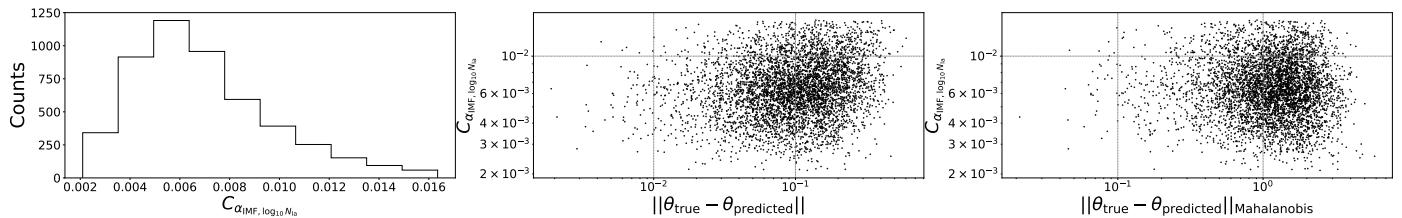
**Fig. C.1.** The left Figure shows the distribution of the correlation between the global parameters $\alpha_{\mathrm{IMF}}$ and $\log_{10} N_{\mathrm{Ia}}$, obtained by the covariance matrix of the sample of each star in the validation set. The central and right Figures show the correlation as a function of the Euclidean and Mahalanobis distance of the true value $\theta_{\mathrm{true}}$ and the sample average $\theta_{\mathrm{predicted}}$. This results shows that the correlation is independent on the accuracy of the inference.
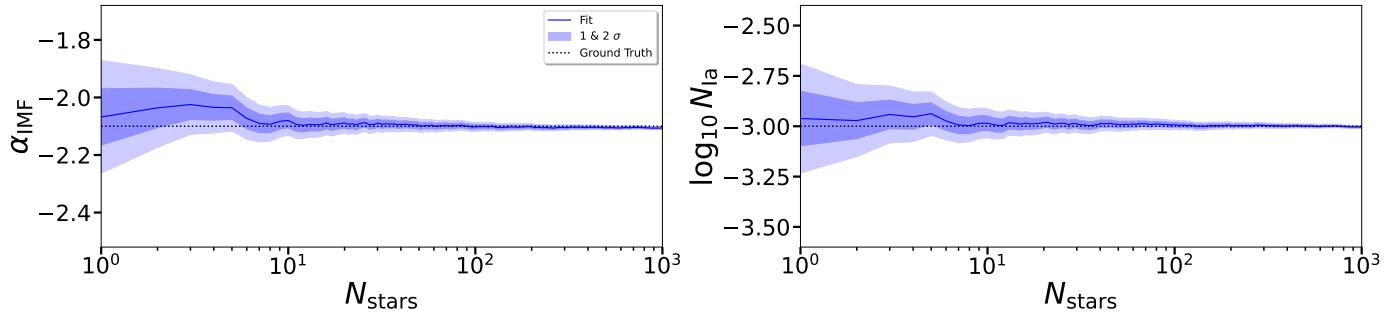


**Fig. D.1.** Same as Fig. 5 but for mock data created with different parameters for $\alpha_{IMF} = -2.1$ and $\log_{10}(N_{Ia}) = -3.0$.
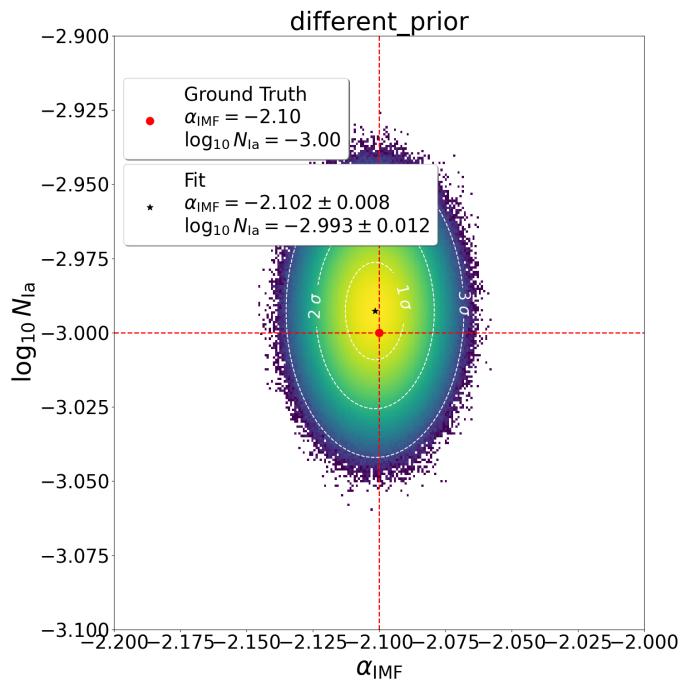


**Fig. D.2.** Same as Fig. 10 but for the mock data taken created with a different parameter combination for $\alpha_{IMF} = -2.1$ and $\log_{10}(N_{Ia}) = -3.0$.