

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Accelerating astronomical and cosmological inference with preconditioned Monte Carlo

### Permalink

<https://escholarship.org/uc/item/183453j4>

### Journal

Monthly Notices of the Royal Astronomical Society, 516(2)

### ISSN

0035-8711

### Authors

Karamanis, Minas

Beutler, Florian

Peacock, John A

et al.

### Publication Date

2022-09-06

### DOI






10.1093/mnras/stac2272

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

# Accelerating astronomical and cosmological inference with Preconditioned Monte Carlo

Minas Karamanis <sup>1</sup>,<sup>\*</sup> Florian Beutler <sup>1</sup>, John A. Peacock <sup>1</sup>, David Nabergoj <sup>2</sup> and Uroš Seljak <sup>3</sup>

<sup>1</sup>*Institute for Astronomy, University of Edinburgh, Royal Observatory, Blackford Hill, Edinburgh EH9 3HJ, UK*

<sup>2</sup>*Faculty of Computer and Information Science, University of Ljubljana, Večna pot 113, 1000 Ljubljana, Slovenia*

<sup>3</sup>*Physics Department, University of California and Lawrence Berkeley National Laboratory Berkeley, CA 94720, USA*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

We introduce *Preconditioned Monte Carlo* (PMC), a novel Monte Carlo method for Bayesian inference that facilitates efficient sampling of probability distributions with non-trivial geometry. PMC utilises a *Normalising Flow* (NF) in order to decorrelate the parameters of the distribution and then proceeds by sampling from the preconditioned target distribution using an adaptive *Sequential Monte Carlo* (SMC) scheme. The results produced by PMC include samples from the posterior distribution and an estimate of the model evidence that can be used for parameter inference and model comparison respectively. The aforementioned framework has been thoroughly tested in a variety of challenging target distributions achieving state-of-the-art sampling performance. In the cases of *primordial feature analysis* and *gravitational wave* inference, PMC is approximately 50 and 25 times faster respectively than *Nested Sampling* (NS). We found that in higher dimensional applications the acceleration is even greater. Finally, PMC is directly parallelisable, manifesting linear scaling up to thousands of CPUs. An open-source Python implementation of PMC, called pocomC, is publicly available at <https://github.com/minaskar/pocomc>.

**Key words:** methods: statistical – methods: data analysis – cosmology: large-scale structure of Universe

## 1 INTRODUCTION

Modern astronomical and cosmological analyses have largely adopted the framework of *Bayesian probability* for tasks of parameter inference and model comparison. In the Bayesian context, the *posterior probability distribution*  $\mathcal{P}(\theta) = P(\theta|\mathcal{D}, \mathcal{M})$ , meaning the probability distribution of the parameters  $\theta$  of a model  $\mathcal{M}$ , given some data  $\mathcal{D}$  and the model  $\mathcal{M}$  is given by Bayes' theorem:

$$\mathcal{P}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{\mathcal{Z}}, \quad (1)$$

where  $\mathcal{L}(\theta) = P(\mathcal{D}|\theta, \mathcal{M})$  is the *likelihood function*,  $\pi(\theta) = P(\theta|\mathcal{M})$  is the *prior probability distribution*, and  $\mathcal{Z} = P(\mathcal{D}|\mathcal{M})$  is the *model evidence* or *marginal likelihood* that acts as a normalisation constant for the posterior probability distribution. For a detailed introduction to Bayesian probability theory we refer the reader to Jaynes (2003); Gregory (2005); MacKay et al. (2003) and the reviews Trotta (2017); Sharma (2017) for its use in astronomy and cosmology.

In tasks of parameter inference, the goal is to infer the values of physical and nuisance parameters from the data along with the respective uncertainties. Mathematically, this is formulated as the problem of estimating expectation values (e.g. mean values, standard deviations, 1-D and 2-D marginal posterior distributions, etc.) that correspond to high-dimensional integrals over the posterior probability density. During the past two decades, *Markov chain Monte Carlo* (MCMC) has been established as the standard computational

tool for the calculation of such integrals (see e.g. (Speagle 2019) for a review). MCMC methods generate a sequence of correlated samples, called a Markov chain, that are distributed according to the posterior probability distribution. Those samples can then be used in order to numerically estimate expectation values. Examples of MCMC software implementations in the astronomical and cosmological community are *emcee* (Foreman-Mackey et al. 2013) and *zeus* (Karamanis et al. 2021).

Most modern MCMC methods are based upon the *Metropolis–Hastings* (MH) paradigm that consists of two steps (Metropolis et al. 1953; Hastings 1970). In the first step, known as the *proposal step*, a new sample is drawn from a known proposal distribution that depends only on the position of the current sample. The validity of the new sample, and thus the decision on whether to add it or not to the Markov chain, is determined in the second step, known as the *acceptance step*, which takes into account the new sample, the old sample (i.e. current state) and the proposal distribution that was used in order to generate it. Arguably, the most important element of an efficient MCMC method is the choice of proposal distribution. The degree to which the proposal distribution characterises the local geometry of the target distribution determines the sampling efficiency (i.e. the rate of effectively independent samples) of the method. Unfortunately, choosing or tuning the optimal proposal distribution for a given target distribution is not an easy task. However, certain optimal proposal distributions are known for specific classes of target distributions. For instance, in the case of a normal or Gaussian target distribution, using a normal proposal distribution of the form  $\mathcal{N}(\theta, 2.38^2\Sigma/D)$ , where  $\Sigma$  is the covariance matrix of the target density,  $\theta$  is the current state of the chain, and  $D$  is the number of dimensions yields the

\* E-mail: minas.karamanis@ed.ac.uk

maximum sampling efficiency scheme with acceptance rate of 23.4% in the acceptance step of MH (Gelman et al. 1997). Alternatively, one can use a simpler proposal distribution of the form  $\mathcal{N}(u, 1)$  where  $u = f(\theta)$  and  $f$  is a suitable transformation. In this case,  $f(\theta)$  is proportional to  $L^{-1}\theta$  where  $L$  is the lower triangular matrix of the *Cholesky decomposition* of the covariance matrix  $\Sigma = LL^T$ . In other words, assuming that a suitable transformation can be found, one can increase the sampling efficiency of an MCMC method. This notion of preconditioning is central for the discussion that will follow in the next section.

In recent years, the need for higher sampling efficiency when the correlations between parameters are strong enough or the posterior exhibits multiple modes, as well as the required computation of the model evidence  $\mathcal{Z}$  for model comparison tasks, motivated the development of more advanced sampling methodologies and algorithms. One very popular approach is the *Sequential Monte Carlo* (SMC) algorithm (Del Moral et al. 2006), which evolves a set of particles through a series of intermediate steps that bridge the gap between the prior distribution and the posterior distribution by geometrically interpolating between them. Another class of algorithms called *Nested Sampling* (NS) (Skilling 2004) attempts to approach the problem of Bayesian computation from a slightly different perspective. Instead of evolving a set of particles through a series of geometrically-interpolated steps between prior and posterior distribution, NS splits the posterior distribution into many slices and attempts to sample each slice individually with an appropriate weighting scheme. Many popular versions and implementations of NS exist in the astronomical literature (Speagle 2020; Buchner 2021; Handley et al. 2015; Feroz et al. 2009). Whereas both SMC and NS largely addressed the problem of multimodality, the performance of both methods is still very sensitive to the geometry of the target distribution, meaning the presence of strong non-linear correlations.

In this paper, we introduce *Preconditioned Monte Carlo* (PMC), a novel Monte Carlo method for Bayesian inference that extends the range of applications of SMC to target distributions with non-trivial geometry, strong non-linear correlations between parameters, and severe multimodality. PMC achieves this by first preconditioning, or transforming the geometry of the target distribution into a more manageable one using a generative model known as a *Normalising Flow* (NF) (Papamakarios et al. 2021), before sampling using a SMC scheme. Hoffman et al. (2019) used a NF to neutralise the bad geometry in Hamiltonian Monte Carlo (HMC) (Betancourt 2017) achieving great results in terms of sampling speed but unreliable estimates for unknown target distributions. Moss (2020) used a NF in order to parameterise efficient MCMC proposals and used it in the context of NS achieving a substantial speedup on several challenging distributions. Both of the aforementioned works used NFs as preconditioning transformations, the first in the context of HMC and the second in NS. In the context of NS and SMC, NFs have also been used as a sampling component of the algorithm (Albergo et al. 2019; Williams et al. 2021; Arbel et al. 2021), albeit not as a preconditioner but as a density from which new samples can be generated independently. The novelty of our work lies in the use of NFs as preconditioning transformations in the context of SMC, thus achieving both robustness and high sampling efficiency.

The structure of the rest of the paper is the following: Section 2 consists of a detailed presentation of the method, Section 3 includes a wide range of empirical tests that act as a demonstration of PMC's sampling performance, and Section 5 is reserved for the conclusions.

We also release a Python implementation of PMC, called *pocomc*, which is publically available at <https://github.com/minaskar/pocomc> and detailed documentation with installation instructions and

examples at <https://pocomc.readthedocs.io>. The code implementation is described in the accompanying paper (Karamanis et al. 2022).

## 2 METHOD

### 2.1 Sequential Monte Carlo

In this subsection, we will present a brief introduction to SMC algorithms. For a more detailed exposition, we refer the reader to Naesseth et al. (2019). We begin by first introducing the concept of *importance sampling*, which is crucial for understanding the function of SMC. Assuming that we have a target probability density  $\pi(\theta)$  that we are able to evaluate up to an unknown multiplicative constant, then if we define another density  $\rho(\theta)$ , called the *importance sampling density*, such that  $\rho(\theta) = 0 \Rightarrow \pi(\theta) = 0$  then the following relation holds for any expectation value:

$$\begin{aligned} E_p[f(\theta)] &= \int f(\theta)w(\theta)\rho(\theta)d\theta \Big/ \int w(\theta)\rho(\theta)d\theta \\ &= E_\rho[f(\theta)w(\theta)]/E_\rho[w(\theta)], \end{aligned} \quad (2)$$

for any function  $f(\theta)$  where  $w(\theta) = p(\theta)/\rho(\theta)$  are called importance weights. We can use samples from the importance density  $\rho(\theta)$  in order to estimate the above expectation value without explicitly sampling from the target density  $p(\theta)$ .

A common measure of the quality of using the importance sampling density  $\rho(\theta)$  to approximate  $p(\theta)$  is the *Effective Sample Size*, defined as:

$$ESS = E_\rho[w(\theta)]^2/E_\rho[w(\theta)^2]. \quad (3)$$

Unfortunately, in high-dimensional scenarios it is difficult to find an appropriate importance sampling density that ensures that the ESS is high enough for the variance of the expectation value to be low. This is exactly the problem that SMC methods address.

SMC samplers extend the importance sampling procedure from the setting of two densities (i.e. importance sampling density and target density) to a sequence of  $T$  probability distributions  $\{p_t\}_{t=1}^T$  in which each individual density  $p_t$  acts as the importance density for the next one in the series. The method proceeds by pushing a collection of  $N$  particles  $\{\theta_t^k\}_{k=1}^N$  through this sequence of densities until the last one is reached. Each iteration of an SMC algorithm consists of three main steps:

(i) **Mutation** – The population of particles is moved from  $\{\theta_{t-1}^k\}_{k=1}^N$  to  $\{\theta_t^k\}_{k=1}^N$  using a *Markov transition kernel*  $K_t(\theta'|\theta)$  that defines the next importance sampling density

$$p_t(\theta') = \int p_{t-1}(\theta)K_t(\theta'|\theta)d\theta. \quad (4)$$

In practice, this step consists of running multiple short MCMC chains (i.e. one for each particle) to get the new states  $\theta'$  starting from the old ones  $\theta$ .

(ii) **Correction** – The particles are reweighted according to the next density in the sequence. This step consists of multiplying the current weight  $W_t^k$  of each particle by the appropriate importance weight:

$$w_t(\theta_t) = p_t(\theta_{t-1})/p_{t-1}(\theta_{t-1}). \quad (5)$$

(iii) **Selection** – The particles are resampled according to their weights which are then set to  $1/N$ . This can be done using *multinomial resampling* or more advanced schemes. The purpose of this step is to eliminate particles with low weight and multiply the ones with high weights.

An important feature of SMC is that it allows for the unbiased estimation of the ratios of normalising constants

$$\mathcal{Z}_t / \mathcal{Z}_{t-1} = \sum_{k=1}^N W_{t-1}^k w_t(\theta_{t-1}^k), \quad (6)$$

between subsequent densities. This is of paramount importance in cases in which the first density in the series corresponds to the prior distribution (i.e. with  $\mathcal{Z} = 1$ ) and the last to the posterior distribution. Then, SMC methods can be used in order to compute the model evidence  $\mathcal{Z}$  for tasks of model comparison.

In principle, there are arbitrary many ways to construct the sequence of densities  $\{p_t\}_{t=1}^T$ . A very common way to do so is to geometrically interpolate between two densities  $\rho(\theta)$  and  $p(\theta)$ :

$$p_t(\theta) \propto \rho(\theta)^{1-\beta_t} p(\theta)^{\beta_t}, \quad t = 1, \dots, T \quad (7)$$

parameterised by a *temperature annealing ladder*:

$$\beta_1 = 0 < \beta_2 < \dots < \beta_T = 1. \quad (8)$$

In the Bayesian context, a natural choice of geometric interpolation is from the prior  $\pi(\theta)$  to the posterior:

$$p_t(\theta) \propto \pi(\theta) \mathcal{L}(\theta)^{\beta_t}, \quad t = 1, \dots, T \quad (9)$$

where  $\mathcal{L}(\theta)$  is the likelihood function. In practice, it can still be difficult to choose a good temperature schedule. However, this can be done adaptively by selecting the next value of  $\beta_t$  such that the ESS is a constant  $\alpha$  fraction of the number of particles  $N$ . Numerically, this can be done by solving

$$\left( \sum_{k=1}^N w_{t+1}^k(\beta_{t+1}) \right)^2 / \sum_{k=1}^N w_{t+1}^k(\beta_{t+1})^2 = \alpha N, \quad (10)$$

the next  $\beta_{t+1}$  such that  $\beta_t < \beta_{t+1} \leq 1$  using, for instance, the *bisection method*.

## 2.2 Normalising Flows

Normalising flows (NF) are generative models, which can facilitate efficient and exact density estimation (Papamakarios et al. 2021). They are based on the formula of change-of-variables  $\theta = f(u)$  where  $u$  is sampled from a base distribution  $u \sim p_u(u)$  (i.e. usually a normal distribution). The NF is a bijective mapping between the base distribution  $p_u(u)$  and the often more complex target distribution  $p_\theta(\theta)$  that can be evaluated exactly using

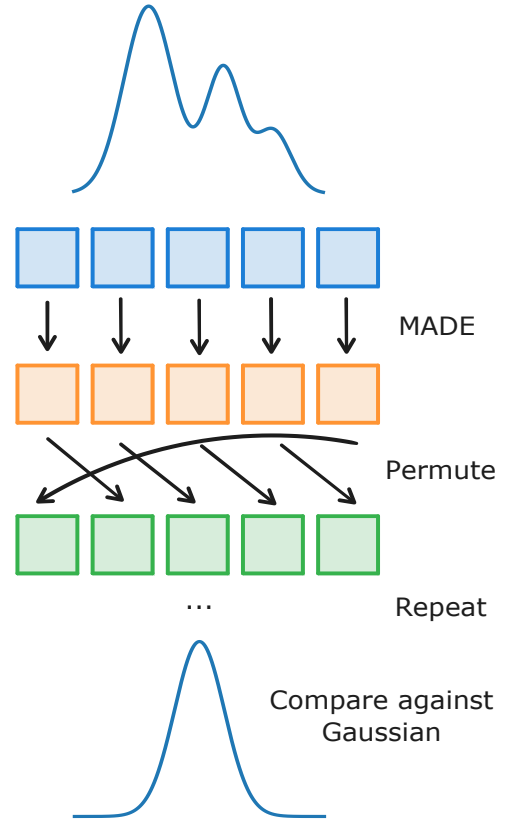
$$p_\theta(\theta) = p_u(f^{-1}(\theta)) \left| \det \left( \frac{\partial f^{-1}}{\partial \theta} \right) \right|, \quad (11)$$

where the Jacobian determinant is tractable.

NFs are usually parameterised by neural networks. However, neural networks are not invertible in general, and the Jacobian is not generally tractable. Special care needs to be taken when choosing the architecture of the neural network to ensure the invertibility of the transformation and the tractability of the Jacobian. For instance, if the forward transformation is  $\theta_i = u_i \exp(\alpha_i) + \mu_i$  and inverse transformation is  $u_i = (\theta_i - \mu_i) \exp(-\alpha_i)$ , where  $\mu_i$  and  $\alpha_i$  are constants, then it is straightforward to show that the Jacobian satisfies

$$\left| \det \left( \frac{\partial f^{-1}}{\partial \theta} \right) \right| = \exp \left( - \sum_i \alpha_i \right). \quad (12)$$

To this end, we chose to use the *Masked Autoregressive Flow* (MAF), which has been used many times successfully for density

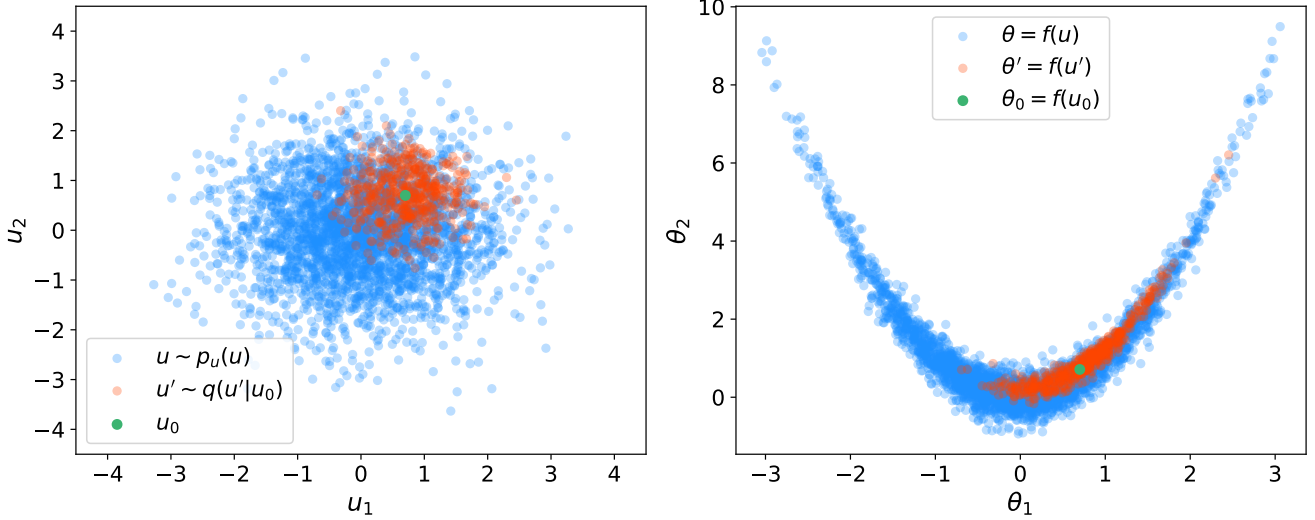


**Figure 1.** Illustration of the inference scheme of a *Masked Autoregressive Flow* (MAF). The arrows show the conditional dependence of the variables as well as the action of the *Masked Autoregressive Density Estimation* (MADE) layer. The input target probability density (top) is mapped into a multivariate normal distribution (bottom). A sequence of MADE layers and permutations is repeated multiple times in order to increase the flexibility of the flow.

estimation tasks due to its superior performance and high flexibility compared to alternative models (Papamakarios et al. 2017). A MAF consists of many stacked layers of a simpler generative model, called *Masked Autoregressive Density Estimator* (MADE) (Germain et al. 2015), with subsequent permutations of its outputs as shown in Figure 1. A MADE model decomposes a joint density  $p(\theta) = \prod_i p(\theta_i | \theta_{1:i-1})$  that ensures that any given value  $\theta_i$  is only a function of the previous values thus maintaining the *autoregressive property*. When the MADE is based on an *autoencoder*, then *masking* is required in order to remove connections between different units in different layers, so as to preserve the aforementioned autoregressive property.

## 2.3 Preconditioning

Most MCMC methods struggle to sample efficiently from highly correlated or skewed target distributions. Often, transforming the parameters of the distribution before sampling, a process also known as *preconditioning*, using appropriate change-of-variable transformations, can help ameliorate this effect by disentangling the dependence between parameters. This is equivalent to choosing an appropriate proposal distribution in the context of *Metropolis–Hastings* (MH) methods. However, finding a valid transformation and selecting an appropriate proposal distribution is often difficult a priori; and there is no obvious way of making this joint choice in an optimal way. For instance, a linear transformation  $\theta \leftarrow L^{-1}\theta$  where  $L$  is the lower



**Figure 2.** The figure illustrates the effect of preconditioning on the *Rosenbrock* distribution. The right panel shows samples (blue) from the true correlated distribution and the left panel shows samples (blue) from the preconditioned/transformed one. The orange samples in the left panel are drawn from a symmetric normal proposal distribution centred around the green point  $u_0$  and they correspond to the respective orange points in the right panel. In other words, the transformed samples from the simple proposal in the left panel correspond to samples that capture the local geometry of the true target distribution in the right panel.

triangular matrix of the *Cholesky decomposition* of the *sample covariance matrix*  $\Sigma = LL^T$  can remove only linear correlations and is not effective against non-linear ones. More sophisticated transformations, such as the use of the *chirp mass* and *mass ratio* instead of the individual black-hole masses in gravitational wave astronomy requires expert knowledge that is problem-specific.

The *Metropolis acceptance criterion* employed by MH methods in order to maintain detailed balance is

$$\alpha = \min \left( 1, \frac{p_\theta(\theta')q(\theta|\theta')}{p_\theta(\theta)q(\theta'|\theta)} \right), \quad (13)$$

where  $p_\theta(\theta)$  is the target distribution and  $q(\theta'|\theta)$  is the proposal distribution. For a general transformation  $\theta = f(u)$  and its inverse  $u = f^{-1}(\theta)$  the modified *Metropolis acceptance criterion* takes the following form

$$\alpha = \min \left( 1, \frac{p_\theta(f^{-1}(u'))q(u|u') \left| \det \frac{\partial f^{-1}(u')}{\partial u'} \right|}{p_\theta(f^{-1}(u))q(u'|u) \left| \det \frac{\partial f^{-1}(u)}{\partial u} \right|} \right), \quad (14)$$

where the Jacobian determinant also appears. In this formulation of MH, the sampler samples the distribution in the transformed space and then samples are pushed through the  $\theta = f(u)$  transformation to the original space. Assuming that the transformation  $\theta = f(u)$  induces a simpler geometry onto the transformed space, sampling using the above acceptance criterion can be substantially more efficient.

Figure 2 shows one such transformation that transforms the banana-shaped *Rosenbrock* distribution into a unit-variance normal distribution and *vice versa*. The same figure also demonstrates the effectiveness of simple proposal distributions  $q(u'|u)$  in the transformed space. A symmetric normal proposal distribution  $q(u'|u_0)$  centred around a point  $u_0$  corresponds to a highly effective proposal distribution in the original space, which captures the local geometry of the target distribution around that point.

## 2.4 Preconditioned Monte Carlo

*Preconditioned Monte Carlo* (PMC) is the result of the amalgamation of SMC, NFs and preconditioning as they were introduced in the previous paragraphs. In particular, we suggest the use of the transformation  $\theta = f(u)$  of a NF in order to precondition the *Mutation* step of SMC. A pseudocode of the algorithm is presented at Algorithm 1. The *Mutation* step in this case consists of  $N$  *Random-Walk Metropolis* (RWM) steps, meaning MH with an isotropic Gaussian proposal distribution centred around the current state of the Markov chain, in which the algorithm targets the preconditioned density. We fix the acceptance rate of MH to its optimal value 23.4% between temperature steps by adapting the proposal scale (Gelman et al. 1997). As the optimal proposal scale of MH for a Gaussian target distribution is

$$\sigma_{\text{opt}} = \frac{2.38}{\sqrt{D}}, \quad (15)$$

where  $D$  is the number of dimensions, we can assess the performance of the NF preconditioner by estimating the ratio of the true scale  $\sigma$  to the optimal one  $\sigma_{\text{opt}}$ . Assuming that the NF preconditioner perfectly the target density and maps it into a unit-variance Gaussian distribution, this ratio should be equal to one. In practice, this ratio can deviate slightly from the optimal value of unity, and one can utilise this ratio as a metric of the preconditioning quality. The number  $N$  of the MCMC steps performed in each iteration is determined adaptively during the run. The process we used is based on the mean correlation coefficient between the initial positions of the particles in the beginning of an iteration and their current positions. In particular, the particles are updated using MCMC until their mean correlation coefficient drops below a prespecified threshold. The lower the threshold, the higher the number  $N$  of MCMC steps. It is important to note that the correlation coefficient is computed in the preconditioned  $u$  space.

**Algorithm 1** Preconditioned Monte Carlo

---

```

1: input Number of particles  $N$ 
2:  $t \leftarrow 1, \beta_1 \leftarrow 0, \mathcal{Z} \leftarrow 1$ 
3: for  $k = 1$  to  $N$  do sample  $\theta_1^k \sim \pi(\theta)$  and set  $W_1^k = 1/N$ 
4: train  $\theta = f(u)$  using  $\{\theta_1^k\}_{k=1}^N$ 
5: while  $\beta_t \neq 1$  do
6:    $t \leftarrow t + 1$ 
7:    $\beta_t \leftarrow$  solution to Eq. 10
8:   for  $k = 1$  to  $N$  do  $w_{t-1}^k \leftarrow W_{t-1}^k \mathcal{L}(\theta) \beta_t - \beta_{t-1}$ 
9:    $\mathcal{Z} \leftarrow \mathcal{Z} \sum_{k=1}^N w_t^k$ 
10:   $\{\tilde{\theta}_{t-1}^k\}_{k=1}^N \leftarrow$  resample  $\{\theta_{t-1}^k\}_{k=1}^N$  according to  $\{w_{t-1}^k\}_{k=1}^N$ 
11:  for  $k = 1$  to  $N$  do  $W_t^k \leftarrow 1/N$ 
12:   $\{\theta_t^k\}_{k=1}^N \leftarrow$  move  $\{\tilde{\theta}_{t-1}^k\}_{k=1}^N$  according to
     $K_t(\{\theta_t^k\}_{k=1}^N \leftarrow \{\tilde{\theta}_{t-1}^k\}_{k=1}^N; f)$ 
13:  train  $\theta = f(u)$  using  $\{\theta_t^k\}_{k=1}^N$ 
14: end while
15: return samples  $\{\theta_t^k\}_{k=1}^N$  and estimate of the marginal likelihood  $\mathcal{Z}$ 

```

---

**Table 1.** The table shows the default values for the hyperparameters of PMC.

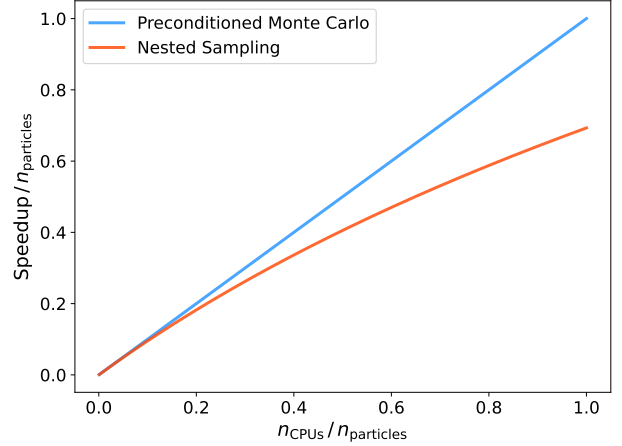
NF hyperparameters		SMC hyperparameters	
blocks	6	particles	1000 – 4000
neurons	$3 \times D$	ESS	95%
batch	1000	threshold	75%
epochs	500		
tolerance	30		
lr	$10^{-2} - 10^{-5}$		
b	0.2		

## 2.5 Hyperparameters

We organise the hyperparameters of PMC into two groups, those related to the normalising flow and those related to the SMC algorithm. The first group consists of structure and training hyperparameters for the NF. The NF structure parameters include the number of MADE layers (blocks), as well as the number of neurons per hidden layer (neurons). The NF training hyperparameters include the learning rate (lr) of the Adam optimiser (Kingma & Ba 2014), the maximum number of epochs (epochs), the training batch size (batch), the tolerance for early stopping (tolerance), and the Laplace prior scale (b) used for regularisation. On the other hand, the SMC hyperparameters include the number of particles (particles), the desired effective sample size (ESS), and the correlation coefficient threshold (threshold). The default values for those hyperparameters are shown in Table 1. We found that this configuration was robust and efficient for a wide range of applications and thus decided to recommend it as the default choice.

## 2.6 Parallelization

An important property of PMC is its ideal scaling with the available number of CPUs. In particular, the mutation step of PMC is exactly parallelisable, meaning that the speedup gained by using more than one CPU scales linearly with the number of CPUs as long as  $n_{\text{CPUs}} \leq n_{\text{particles}}$ . Similar methods that also use a large collection of particles scale less favourably. For instance, *Nested Sampling* (NS) exhibits sub-linear scaling as shown in Figure 3 of Handley et al. (2015). The aforementioned scaling characteristic of PMC renders it ideal

**Figure 3.** Parallelization of PMC compared to nested sampling. PMC (blue) exhibits linear speedup compared to the sub-linear one achieved by NS (orange).

for computationally costly applications that are often encountered in astronomy and cosmology.

## 3 EMPIRICAL EVALUATION

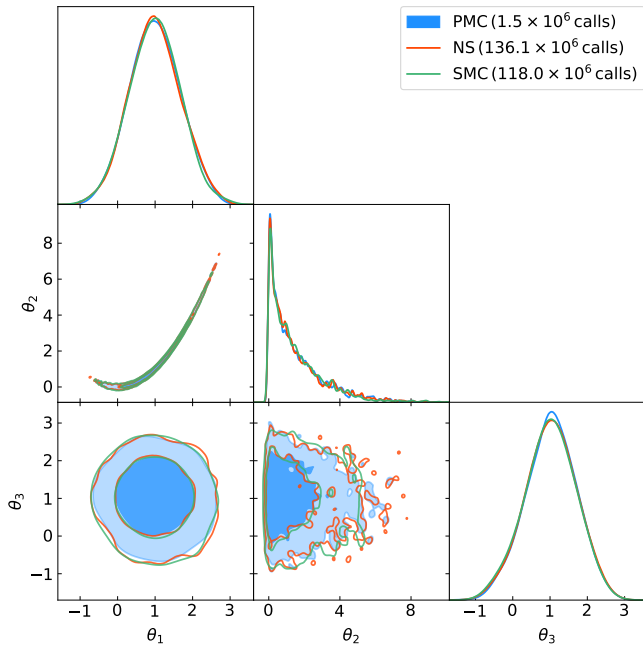
In this section, we present two toy examples and two realistic parameter inference examples that reproduce common astronomical and cosmological analyses. In all cases, the hyperparameters of PMC were set to their default values as shown in Table 1. In both analyses, the performance of PMC is compared to that of SMC using the same settings (e.g. number of particles, ESS, etc.) as PMC but no preconditioning, as well as *Nested Sampling* (NS), a popular particle Monte Carlo alternative<sup>1</sup>. The metric that we use in order to evaluate the performance of each method is the total number of model evaluations performed until convergence. Convergence in all methods is well-defined: in PMC and SMC the algorithm converges when  $\beta = 1$ , whereas in NS the run stops when less than 1% of the model evidence is left unaccounted. All other computational costs are negligible, including the training and evaluation of the normalising flow in the case of PMC that only required a few seconds for the whole inference procedure. All methods used 1000 particles.

### 3.1 Rosenbrock distribution

The first toy example that we used is the *Rosenbrock* distribution, which exhibits strong non-linear correlation between its parameters. For this reason, the *Rosenbrock* distribution has often been used as a benchmark target for optimization and sampling tasks. Here we use a 20-dimensional generalisation defined through the probability density function given by:

$$\log P(\theta) = - \sum_{i=1}^{N/2} \left[ 10 \left( \theta_{2i-1}^2 - \theta_{2i} \right)^2 + (\theta_{2i-1} - 1)^2 \right]. \quad (16)$$

<sup>1</sup> We used the popular Python implementation dynesty (Speagle 2020) for NS.



**Figure 4.** Illustration of the 1–dimensional and 2–dimensional marginal posteriors for the first three out of 20 parameters of the *Rosenbrock* distribution. The figure shows the 1– $\sigma$  and 2– $\sigma$  contours generated by *Preconditioned Monte Carlo* (PMC) in blue, *Nested Sampling* (NS) in orange, and *Sequential Monte Carlo* (SMC) in green. The legend also shows the computational cost of each method in terms of the total number of required model evaluations until convergence is reached.

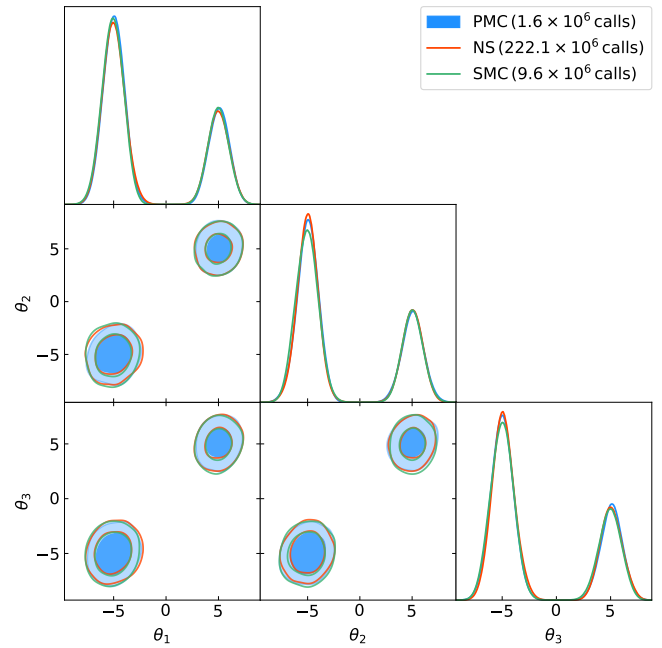
**Table 2.** The table shows a comparison of PMC, NS, and SMC in terms of their computational cost (i.e. total number of model evaluations until convergence).

Distribution	Model evaluations ( $\times 10^6$ )		
	PMC	NS	SMC
Rosenbrock	<b>1.5</b>	136.1	118.0
Gaussian Mixture	<b>1.6</b>	222.1	9.6
Primordial Features	<b>0.4</b>	21.3	19.5
Gravitational Waves	<b>0.4</b>	10.2	4.6

We use flat priors  $\mathcal{U}(-10, 10)$  for all parameters. Figure 4 shows the 2–dimensional marginal posterior for the first two parameters as generated by the three methods. The total computational cost of PMC, NS, and SMC is  $1.5 \times 10^6$ ,  $136.1 \times 10^6$ , and  $118.0 \times 10^6$  model evaluations, respectively. PMC requires approximately 1/91 of the number of model evaluations that NS does and approximately 1/79 of those that SMC does.

### 3.2 Gaussian Mixture

The second toy example that we used is a 50–dimensional Gaussian Mixture with two components, one of them being twice as massive as the other. This is a highly multimodal problem as the target distribution exhibits two distinct modes that are well separated. Just as in the *Rosenbrock* case, we use flat priors  $\mathcal{U}(-10, 10)$  for all parameters. Figure 5 shows the 1–dimensional and 2–dimensional marginal posteriors for the first three parameters as generated by the three methods. The total computational cost of PMC, NS, and SMC is  $1.6 \times 10^6$ ,  $222.1 \times 10^6$ , and  $9.6 \times 10^6$  model evaluations respec-



**Figure 5.** Illustration of the 1–dimensional and 2–dimensional marginal posteriors for the first three out of 50 parameters of the two–component Gaussian mixture distribution. The figure shows the 1– $\sigma$  and 2– $\sigma$  contours generated by *Preconditioned Monte Carlo* (PMC) in blue, *Nested Sampling* (NS) in orange, and *Sequential Monte Carlo* (SMC) in green. The legend also shows the computational cost of each method in terms of the total number of required model evaluations until convergence is reached.

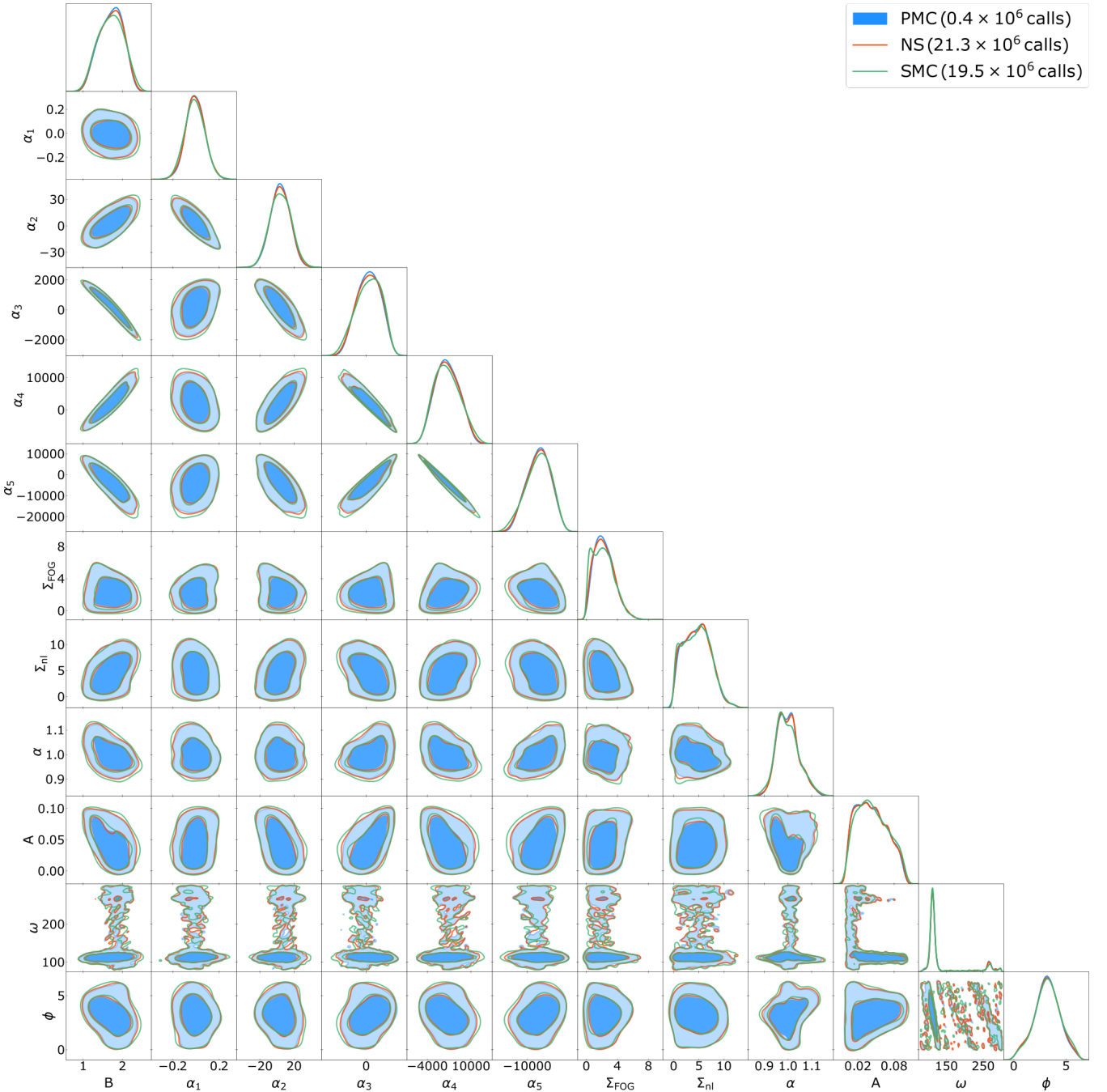
tively. PMC requires approximately 1/139 of the number of model evaluations that NS does and 1/6 of those that SMC does.

### 3.3 Primordial Features

The first realistic application that we study is the search for primordial features along the Baryon Acoustic Oscillation (BAO) signature in the distribution of galaxies observed by the Sloan Digital Sky Survey (SDSS) (Eisenstein et al. 2011). In particular, the data that we analysed come from the 12th data release (DR12) of the high–redshift North Galactic Cap (NGC) sample of the Baryon Oscillation Spectroscopic Survey (BOSS) (Dawson et al. 2013). Our analysis follows closely that of Beutler et al. (2019) for the linear oscillation model. The inference problem includes 12 free parameters with either flat/uniform or normal priors. Figure 6 shows the 1–dimensional and 2–dimensional marginal posteriors of the aforementioned analysis. The posterior distribution exhibits a highly non–Gaussian geometry that can hinder the sampling performance of conventional methods. The total computational cost of PMC, NS, and SMC is  $0.4 \times 10^6$ ,  $21.3 \times 10^6$ , and  $19.5 \times 10^6$  model evaluations respectively. PMC requires approximately 1/53 of the number of model evaluations that NS does, and 1/49 of those that SMC does.

### 3.4 Gravitational Waves

The second realistic application is the simulated gravitational wave analysis of an injected signal. For this, we used the standard CBC (i.e. compact binary coalescence) injected signal configuration provided by BILBY (Ashton et al. 2019). The inference problem includes 13 free parameters with a variety of common priors. Figure 7 shows



**Figure 6.** Illustration of the 1-dimensional and 2-dimensional marginal posteriors for the 12 parameters of the primordial features posterior. The figure shows the 1- $\sigma$  and 2- $\sigma$  contours generated by *Preconditioned Monte Carlo* (PMC) in blue, *Nested Sampling* (NS) in orange, and *Sequential Monte Carlo* (SMC) in green. The legend also shows the computational cost of each method in terms of the total number of required model evaluations until convergence is reached.

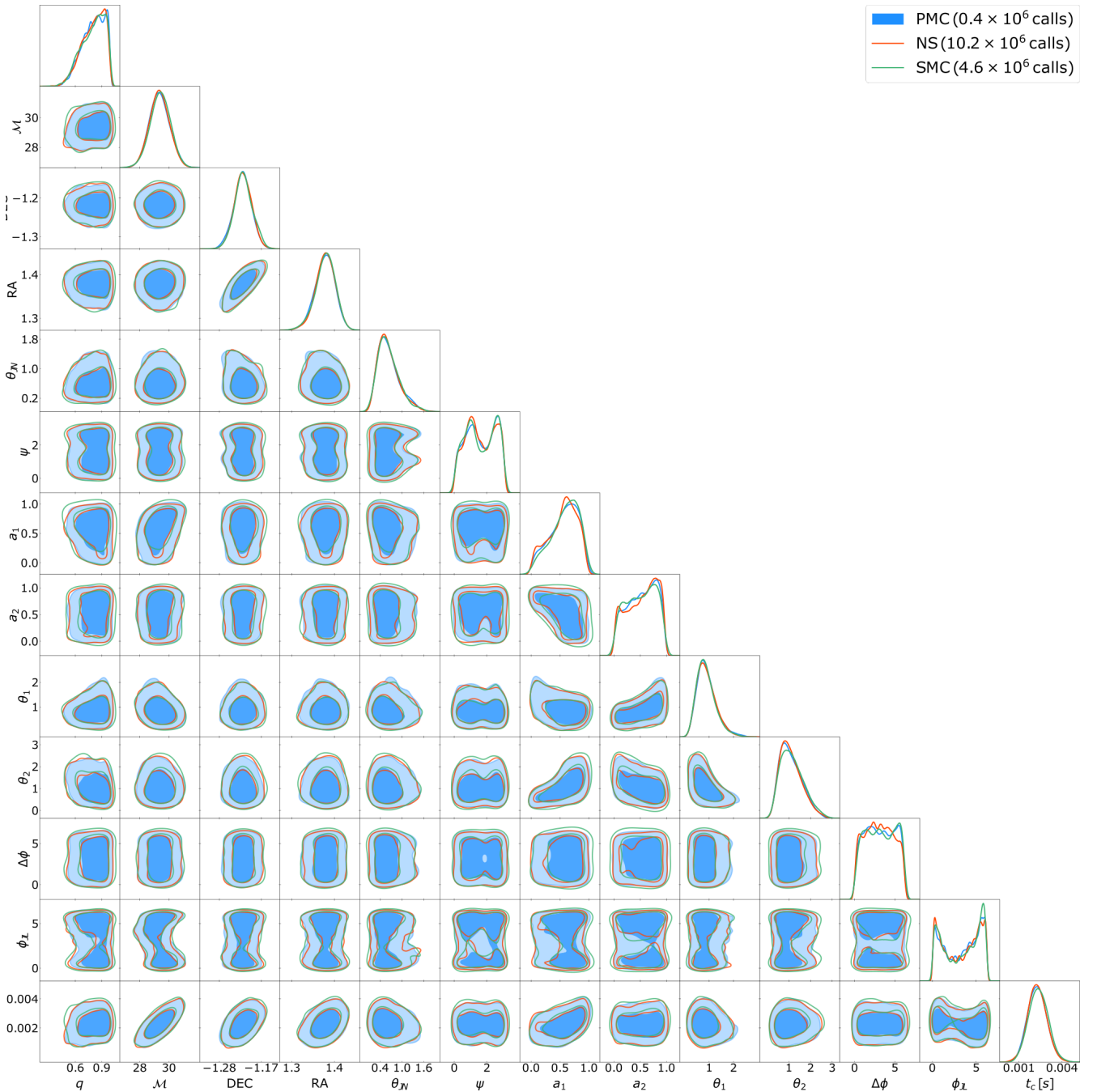
the 1-dimensional and 2-dimensional marginal posteriors of the aforementioned analysis. The posterior distribution exhibits a highly non-Gaussian geometry that can hinder the sampling performance of conventional methods. The total computational cost of PMC, NS, and SMC is  $0.4 \times 10^6$ ,  $10.2 \times 10^6$ , and  $4.6 \times 10^6$  model evaluations respectively. PMC requires approximately 1/25 of the number of model evaluations that NS does and 1/11 of those that SMC does.

#### 4 DISCUSSION

While we have demonstrated PMC's superior sampling performance for a number of target distributions, including two real-world applications, the real test is based on researchers applying the method to their analyses. Different applications pose different computational challenges and there is no one single sampler to rule them all. Sometimes, certain kinds of distributions will be better handled by other, perhaps simpler, approaches.

In general, we expect PMC to be a useful tool when dealing with computationally expensive likelihood functions and highly corre-





**Figure 7.** Illustration of the 1–dimensional and 2–dimensional marginal posteriors for the 13 parameters of the gravitational waves posterior. The figure shows the 1– $\sigma$  and 2– $\sigma$  contours generated by *Preconditioned Monte Carlo* (PMC) in blue, *Nested Sampling* (NS) in orange, and *Sequential Monte Carlo* (SMC) in green. The legend also shows the computational cost of each method in terms of the total number of required model evaluations until convergence is reached.

lated or multimodal posteriors. There are two main reasons for this. First, training of the normalising flow takes about  $O(1\text{ s})$  per iteration on a laptop computer, whereas the actual vectorised evaluation of the bijective mapping takes almost  $O(10\text{ ms})$  per MCMC step for the whole population of particles. This means that if the cost of evaluating the likelihood is low enough to be comparable to that of the normalising flow, as discussed above, the chances are that there are simpler methods (e.g. MCMC) that can obtain the results more quickly. The second reason has to do with the geometry of the posterior distribution. If the latter is trivial enough (e.g. approximately Gaussian with no non–linear correlation or multiple modes), then

the use of the normalising flow as a preconditioner would offer no benefit and instead only help delay the run.

On the other hand, if both of these conditions are met, that is, the likelihood function is computationally expensive, as is often the case in cosmology, and the posterior is non–Gaussian, then PMC can be a valuable asset in the astronomer’s toolkit. Furthermore, when the cost of evaluating the likelihood function is large enough to dominate both the normalising flow evaluation and any potential *MPI* communication overhead, one can capitalise on the availability of multiple CPUs in order to accelerate PMC. In particular, if the evaluation of the likelihood function takes  $O(1\text{ s})$ , one should be

able to use up to thousands of CPUs, potentially parallelising all or a substantial fraction of the particles.

## 5 CONCLUSIONS

We introduced PMC, a preconditioned generalisation of the standard SMC algorithm. PMC is a novel sampling method that can accelerate Bayesian inference and model comparison in computationally challenging astronomical and cosmological analyses.

After introducing the method in Section 2, we presented a thorough demonstration of *Preconditioned Monte Carlo*'s sampling capabilities by comparing its sampling performance to that of *Nested Sampling* and *Sequential Monte Carlo* in a range of target distributions characterised by non-trivial geometry. The results are presented in Table 2. We found that *Preconditioned Monte Carlo* is one to two orders of magnitude faster than either *Nested Sampling* or *Sequential Monte Carlo*, both of which performed similarly to each other. Furthermore, in the realistic analyses of primordial features and gravitational waves, *Preconditioned Monte Carlo* required approximately 50 and 25 times fewer model evaluations compared to NS in order to converge. The reduced computational cost, combined with the superior parallelisation scaling, renders *Preconditioned Monte Carlo* ideal for astronomical and cosmological Bayesian analyses with computationally expensive, strongly correlated, multimodal and high-dimensional posteriors.

We hope that *Preconditioned Monte Carlo* will prove useful to the astronomical community by facilitating challenging Bayesian data analyses and enabling the investigation of complex models and sparse datasets. We also released a Python implementation of *Preconditioned Monte Carlo*, called *pocomc*, which is publicly available at <https://github.com/minaskar/pocomc> and detailed documentation with installation instructions and examples at <https://pocomc.readthedocs.io>.

## ACKNOWLEDGEMENTS

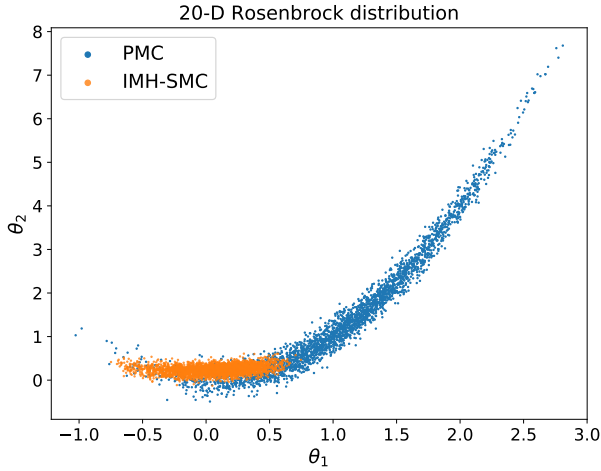
The authors extend their gratitude to Jamie Donald-McCann, Richard Grumitt, Biwei Dai, and James Sullivan for providing useful comments. MK would also like to thank George Vretinaris for providing valuable feedback on an early version of the code. This work has benefited from a variety of Python packages including *numpy* (Van Der Walt et al. 2011), *scipy* (Virtanen et al. 2020), *torch* (Paszke et al. 2019), *matplotlib* (Hunter 2007), *seaborn* (Waskom 2021), *getdist* (Lewis 2019), *sklearn* (Pedregosa et al. 2011), *tqdm* (da Costa-Luis 2019), *dynesty* (Speagle 2020), and *mpi4py* (Dalcin et al. 2011). This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement 853291), and by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research under Contract No. DE-AC02-05CH11231 at Lawrence Berkeley National Laboratory to enable research for Data-intensive Machine Learning and Analysis. FB is a University Research Fellow.

## DATA AVAILABILITY

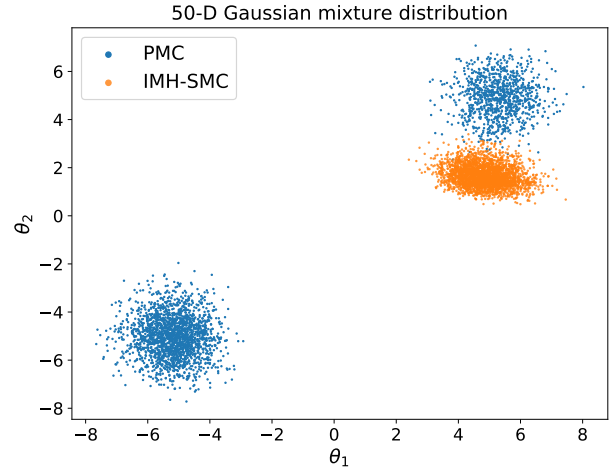
All data used in this work are publicly available. Power spectrum estimates, covariance matrices and window functions used in the cosmological inference example are available at [http://www.sdss3.org/science/boss\\_publications.php](http://www.sdss3.org/science/boss_publications.php).

## REFERENCES

- Albergo M., Kanwar G., Shanahan P., 2019, *Physical Review D*, 100, 034515
- Arbel M., Matthews A., Doucet A., 2021, in *International Conference on Machine Learning*. pp 318–330
- Ashton G., et al., 2019, *The Astrophysical Journal Supplement Series*, 241, 27
- Betancourt M., 2017, arXiv preprint arXiv:1701.02434
- Beutler F., Biagetti M., Green D., Slosar A., Wallisch B., 2019, *Physical Review Research*, 1, 033209
- Buchner J., 2021, arXiv preprint arXiv:2101.09604
- Dalcin L. D., Paz R. R., Kler P. A., Cosimo A., 2011, *Adv. Water Resour.*, 34, 1124
- Dawson K. S., et al., 2013, *AJ*, 145, 10
- Del Moral P., Doucet A., Jasra A., 2006, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68, 411
- Eisenstein D. J., et al., 2011, *AJ*, 142, 72
- Feroz F., Hobson M., Bridges M., 2009, *Monthly Notices of the Royal Astronomical Society*, 398, 1601
- Foreman-Mackey D., Hogg D. W., Lang D., Goodman J., 2013, *Publications of the Astronomical Society of the Pacific*, 125, 306
- Gelman A., Gilks W. R., Roberts G. O., 1997, *The annals of applied probability*, 7, 110
- Germain M., Gregor K., Murray I., Larochelle H., 2015, in *International Conference on Machine Learning*. pp 881–889
- Gregory P., 2005, *Bayesian logical data analysis for the physical sciences: a comparative approach with mathematica® support*. Cambridge University Press
- Handley W., Hobson M., Lasenby A., 2015, *Monthly Notices of the Royal Astronomical Society*, 453, 4384
- Hastings W. K., 1970, *Biometrika*, 57, 97
- Hoffman M., Sountsov P., Dillon J. V., Langmore I., Tran D., Vasudevan S., 2019, arXiv preprint arXiv:1903.03704
- Hunter J. D., 2007, *IEEE Ann. Hist. Comput.*, 9, 90
- Jaynes E. T., 2003, *Probability theory: The logic of science*. Cambridge university press
- Karamanis M., Beutler F., Peacock J. A., 2021, *Monthly Notices of the Royal Astronomical Society*, 508, 3589
- Karamanis M., Nabergoj D., Beutler F., Peacock J. A., Seljak U., 2022, in prep
- Kingma D. P., Ba J., 2014, arXiv preprint arXiv:1412.6980
- Lewis A., 2019, preprint (arXiv:1910.13970)
- MacKay D. J., Mac Kay D. J., et al., 2003, *Information theory, inference and learning algorithms*. Cambridge university press
- Metropolis N., Rosenbluth A. W., Rosenbluth M. N., Teller A. H., Teller E., 1953, *J. Chem. Phys.*, 21, 1087
- Moss A., 2020, *Monthly Notices of the Royal Astronomical Society*, 496, 328
- Naesseth C. A., Lindsten F., Schön T. B., 2019, arXiv preprint arXiv:1903.04797
- Papamakarios G., Pavlakou T., Murray I., 2017, *Advances in neural information processing systems*, 30
- Papamakarios G., Nalisnick E., Rezende D. J., Mohamed S., Lakshminarayanan B., 2021, *Journal of Machine Learning Research*, 22, 1
- Paszke A., et al., 2019, *Advances in neural information processing systems*, 32
- Pedregosa F., et al., 2011, *J. Mach. Learn. Res.*, 12, 2825
- Sharma S., 2017, *Annual Review of Astronomy and Astrophysics*, 55, 213
- Skilling J., 2004, in *AIP Conference Proceedings*. pp 395–405
- Speagle J. S., 2019, preprint (arXiv:1909.12313)
- Speagle J. S., 2020, *Monthly Notices of the Royal Astronomical Society*, 493, 3132
- Trotta R., 2017, arXiv preprint arXiv:1701.01467
- Van Der Walt S., Colbert S. C., Varoquaux G., 2011, *Comput. Sci. Eng.*, 13, 22
- Virtanen P., et al., 2020, *Nat. Methods*, 17, 261
- Waskom M. L., 2021, *J. Open Source Softw.*, 6, 3021
- Williams M. J., Veitch J., Messenger C., 2021, *Physical Review D*, 103, 103006



**Figure A1.** Comparison of the first two parameters of samples generated using PMC (blue) and IMH-SMC (orange) for the 20-D Rosenbrock target distribution. PMC produces representative samples, whereas IMH-SMC does not.



**Figure A2.** Comparison of the first two parameters of samples generated using PMC (blue) and IMH-SMC (orange) for the 50-D two-component Gaussian mixture target distribution. PMC produces representative samples, whereas IMH-SMC does not.

da Costa-Luis C. O., 2019, *J. Open Source Softw.*, 4, 1277

## APPENDIX A: COMPARISON TO INDEPENDENT METROPOLIS-HASTINGS SEQUENTIAL MONTE CARLO

Recent practice in the literature (Albergo et al. 2019; Williams et al. 2021; Arbel et al. 2021) is to use normalising flows as auxiliary densities for *Importance Sampling (IS)* and *Independent Metropolis-Hastings (IMH)* estimators. The latter approach can also be accommodated in the context of *Sequential Monte Carlo (SMC)* as an alternative to PMC. For this reason we will offer an experimental comparison of PMC to IMH-SMC.

The IMH-SMC algorithm is identical to Algorithm 1 with the exception that the mutation step of line 12 takes place using the modified *Metropolis acceptance criterion*

$$\alpha = \min \left( 1, \frac{p_{\theta}(f^{-1}(u'))q(u) \left| \det \frac{\partial f^{-1}(u')}{\partial u'} \right|}{p_{\theta}(f^{-1}(u))q(u') \left| \det \frac{\partial f^{-1}(u)}{\partial u} \right|} \right), \quad (\text{A1})$$

instead of equation 14. The difference between the two criteria is that the proposal distribution  $q(u) = \mathcal{N}(u|0, 1)$  is no longer conditional on the previous state of the Markov chain.

The number  $M$  of IMH steps performed in each iteration of IMH-SMC is determined adaptively during the run, based on the observed acceptance rate  $\alpha$ , as

$$M = \frac{\log(1-p)}{\log(1-\alpha)}, \quad (\text{A2})$$

where  $p$  is the target probability of generating a new independent sample. In our examples below, the value of  $p$  is chosen such that the computational cost of IMH-SMC is similar to that of PMC for the same example. This results in  $p > 0.99$  which corresponds to very conservative sampling.

Despite this, as shown in Figures A1 and A2, for the 20-dimensional Rosenbrock and the 50-dimensional two-component Gaussian mixture studied in the main text respectively, IMH-SMC does not manage to produce typical samples from the posterior distribution. It is important to note here that the acceptance rate of

IMH-SMC was high throughout both runs, and as such offered no indication on its own that NF is not correct.

The origin of this discrepancy between IMH-SMC and PMC in both cases, and the ultimate inability of IMH-SMC to compete with PMC, originates in the substantial mismatch between the NF distribution and target distribution in high dimensions and the subsequent over-fitting of the NF to the particle distribution leading to a narrower distribution. The high acceptance rate does not imply high quality of NF solution, and other tests of the quality of solution are needed, such as comparing expectation of  $\log p$  between samples from NF and true MCMC samples. On the other hand, PMC does not suffer from this pathology as the local exploration offered by MCMC helps diversify the particles in order to avoid over-fitting. Furthermore, local MCMC methods generally scale better with the number of dimensions compared to IMH and IS.

This paper has been typeset from a  $\text{\TeX}/\text{\LaTeX}$  file prepared by the author.