

---

# Transport Elliptical Slice Sampling

---

**Alberto Cabezas**  
Lancaster University

**Christopher Nemeth**  
Lancaster University

## Abstract

We propose a new framework for efficiently sampling from complex probability distributions using a combination of normalizing flows and elliptical slice sampling (Murray et al., 2010). The central idea is to learn a diffeomorphism, through normalizing flows, that maps the non-Gaussian structure of the target distribution to an approximately Gaussian distribution. We then use the elliptical slice sampler, an efficient and tuning-free Markov chain Monte Carlo (MCMC) algorithm, to sample from the transformed distribution. The samples are then *pulled back* using the inverse normalizing flow, yielding samples that approximate the stationary target distribution of interest. Our transport elliptical slice sampler (TESS) is optimized for modern computer architectures, where its adaptation mechanism utilizes parallel cores to rapidly run multiple Markov chains for a few iterations. Numerical demonstrations show that TESS produces Monte Carlo samples from the target distribution with lower autocorrelation compared to non-transformed samplers, and demonstrates significant improvements in efficiency when compared to gradient-based proposals designed for parallel computer architectures, given a flexible enough diffeomorphism.

## 1 INTRODUCTION

Markov Chain Monte Carlo (MCMC) algorithms enable scientists to draw samples from complex distributions, which are typically produced by models that aim to represent the intricate details found in real-world datasets. The exploration of these complex and high-dimensional distributions is challenging, and to be efficient, practitioners use the local pointwise information of the target distribution to create a

Markov chain of dependent samples. The ideal outcome would be to have independent samples, but the Markov chain approach generates samples that are correlated sequentially. Therefore, a major focus in MCMC research is to develop algorithms that reduce these correlations and generate samples that approximate independence.

Designing efficient MCMC algorithms usually relies on using local gradient information from the target distribution; by discretizing, for example, Hamiltonian (Duane et al., 1987; Neal et al., 2011) or Langevin (Rosicky et al., 1978; Grenander and Miller, 1994) dynamics of a process stationary on our target distribution. Calculating gradients has been automated (Linnainmaa, 1976) but optimizing these algorithms to efficiently minimize both computations and correlations between sequential samples requires algorithmic parameters to be manually tuned. Much work has been dedicated to developing efficient, black-box methods to tune these parameters, with notable examples including the NUTS (Hoffman et al., 2014) algorithm which is widely available in probabilistic programming languages (Salvatier et al., 2016; Carpenter et al., 2017; Bingham et al., 2019; Phan et al., 2019).

Within the machine learning community, variational inference (VI; Jordan et al., 1999) has grown in popularity as an inexact but comparatively faster approach to solving the same inferential problem. As such, MCMC has lost its preferential status as the default approach for Bayesian inference for prediction and uncertainty quantification in this thriving community. Recent efforts (Hoffman et al., 2021; Hoffman and Sountsov, 2022) have focused on speeding up MCMC by focusing on widening instead of lengthening computations on modern computer architectures, e.g. utilizing GPUs or TPUs, which allow for vast parallel computations. Tuning parallel MCMC chains has proven to be a somewhat different challenge from its sequential counterpart (Radul et al., 2020) and parallel efforts need to consider the lockstep necessity of gradient evaluations of parallel chains on modern vector oriented libraries (Abadi et al., 2016; Bradbury et al., 2018; Paszke et al., 2019).

Figure 1: Illustration of Algorithm 1 using an exact transport map, i.e. equality holds: sampling from the Banana density  $(x_1; x_2) / \exp[-x_1^2/8 + x_2^2/x_1^2 - 2]$  using the transport map  $T(u_1; u_2) = (\sqrt{8}u_1; u_2 + 2u_1^2)$  starts by transforming the target space to the reference space via a change of variables, drawing samples from an ellipsis on the extended reference space (not pictured) and pushing samples back to the target space.

## 2 TRANSPORT ELLIPTICAL SLICE SAMPLER

In this paper we assume that  $X \subset \mathbb{R}^d$  are model parameters and  $D$  represents our data. Our goal is to then approximate the posterior distribution, where by Bayes rule the posterior density is given by  $p(x) \propto L(D|x) \pi_0(x)$ , for  $L(D|x)$  the likelihood function and  $\pi_0(x)$  the prior density. Our goal is to introduce a new MCMC algorithm which leverages the tuning-free nature of elliptical slice sampling with the efficient density transformation tools of normalising flows, thus creating the transport elliptical slice sampler (TESS); an adaptive mechanism that allows scientists to perform fast parallel sampling from unnormalized densities. An intuitive pictorial representation of our TESS algorithm is given in Figure 1.

### 2.1 Elliptical slice sampling

Introduced by Murray et al. (2010) as a simple MCMC algorithm with no tuning parameters, the elliptical slice sampler builds on a Metropolis-Hasting sampler introduced by Neal (1998), which is designed for situations where the prior  $\pi_0(x)$  is Gaussian. Without loss of generality, we can assume that the prior is a standard Gaussian density  $\pi_0(x) = \mathcal{N}(x)$ <sup>1</sup>. The algorithm of Neal (1998) proceeds by first proposing new state of the Markov chain  $x^0 = \mu + \sigma \sqrt{2}x + v$ ; where  $v \in \mathbb{R}^d$  is an independent momentum variable following a standard Gaussian distribution. The proposal moves along the half ellipse which connects the points  $x$  and  $v$  which pass through, for values  $\alpha \in [0, 1]$ .

Elliptical slice sampling, instead, uses the proposal  $x \cos \alpha + v \sin \alpha$  which moves on the full ellipse connecting  $x$ ,  $x + v$  and  $v$  for  $\alpha \in [0, 2\pi]$ . Both proposals leave the prior density invariant and elliptical slice sampling uses the slice sampling algorithm (Neal, 2003) to choose a value

which ensures that the likelihood  $L(D|x)$  is invariant. Overall, this proposal scheme keeps the target posterior invariant (Murray et al., 2010), details of which are presented in the Supplementary material for completeness.

### 2.2 Normalizing flows

Normalizing flows (NF; Rezende and Mohamed, 2015) are a flexible class of transformations produced by the sequential composition of invertible and differentiable mappings. Using NF involves choosing a simple reference density  $q(u)$ , for example a standard Gaussian distribution  $\mathcal{N}(u)$ , and a parameterized diffeomorphism  $T$ , with optimized parameters, to transform the reference density to our target  $p(x)$  via a change of variables. In other words, we want to find a map  $T$  such that for  $u \in \mathbb{R}^d$  and  $x = T(u)$  we have  $p(x) = q(u) |J_T(u)|$ . Assuming this function exists, applying a change of variable yields the following identities

$$p(x) = q(T^{-1}(x)) |J_{T^{-1}}(x)| = \hat{p}(x) \quad (1)$$

$$q(u) = p(T(u)) |J_T(u)| = \hat{q}(u); \quad (2)$$

where  $J_T$  and  $J_{T^{-1}}$  are the Jacobian matrices of  $T$  and its inverse, respectively. In the context of VI, an approximation of  $\hat{p}(x)$  would serve as an approximation to our target density when carrying-out inference, since this approximation is both normalized and trivial to sample from.

### 2.3 Fixed transport maps with elliptical slice sampling

To fulfill our requirement for a simple and cost-effective MCMC proposal, we begin by generalizing the dimension-independent, gradient-free, and tuning-free elliptical slice sampler. We will add tuning parameters to our generalized elliptical slice sampler using NF. The diffeomorphism will be responsible for efficiently exploring the posterior target density by transforming the proposal's dynamics and tracing the contours of a standard Gaussian density to follow

<sup>1</sup>shift and scale  $\sigma$  if non-standard.

the contours of an approximation of the target. Following previous works in the transport Monte Carlo (as described in Section 3), we present TESS as a two-step procedure. Firstly, we learn the transport map between the target and reference densities. Secondly, we utilize the transport map within the elliptical slice sampler to generate samples from the target density.

1. Map optimization To estimate the parameters of our NF map we minimize a divergence between our target density  $\hat{\mu}(x)$  and the push-forward reference density  $\hat{\nu}(x)$  (1). For our intended purpose, by the law of the unconscious statistician, this is equivalent to minimizing the divergence between the pull-back target density  $\hat{\mu}(u)$  (2) and the reference density  $\hat{\nu}(u)$ . The Kullback-Leibler divergence (Kullback and Leibler, 1951) is arguably the most widely used and studied divergence, here presented in the context of approximate Bayesian inference but also studied in other branches of statistics and information theory (Joyce, 2011). It not only has a tractable Monte Carlo estimate, it is directly related to the foundation of VI and provides intuition into the connection between maximizing the likelihood of observational data and minimizing the distance between target and reference densities (Blei et al., 2017),

$$KL(\hat{\mu} \parallel \hat{\nu}) = \int \log \frac{\hat{\mu}(x)}{\hat{\nu}(x)} \hat{\mu}(x) dx \quad (3)$$

The optimal transport map is found by optimizing the parameters of the diffeomorphism  $T$  such that the Kullback-Leibler divergence between the target and reference densities is minimised, i.e.

$$T = \arg \min_T KL(\hat{\mu} \parallel \hat{\nu} \circ T) \quad (4)$$

2. Sampling from the target Our proposed sampling method generalizes the elliptical slice sampler by targeting the extended state space  $(x, v)$  for any posterior density  $\hat{\mu}(x)$ , regardless of the choice of prior distribution. The target density is preconditioned using a transform via normalizing flow to map to a standard Gaussian distribution. That is, given a map  $T$ , with fixed parameters  $\theta$ , such that  $\hat{\mu}(u) = \hat{\mu}(T(u))$  we proceed as follows: i) from an initial state  $(x; v) = (T(u); v)$ , ii) move around an ellipse connecting  $u$  and  $v$  and iii) accept the new state according to a slice variable chosen uniformly on the interval  $[0, |u - v|]$ . One iteration of this method is detailed in Algorithm 1.

Proposition 1 The transition kernel of the Markov chain derived from Algorithm 1 leaves the target density  $\hat{\mu}(x)$  invariant.

The TESS algorithm is likely to be geometrically ergodic under certain transformations, if those transformations lead to nice tail properties on the pulled back target. A sketch of this argument follows from three key components: (i)

### Algorithm 1 Transport Elliptical Slice Sampler

```

Require:  $\mu; T(\cdot); \hat{\mu}(\cdot)$ 
1:  $v \sim N(0; I_d)$ 
2:  $w \sim \text{Uniform}(0; 1)$ 
3:  $\log s = \log \hat{\mu}(u) + \log \hat{\nu}(v) + \log w$ 
4:  $u \sim \text{Uniform}(0; 2\pi)$ 
5:  $[u^{\min}; u^{\max}] = [u - 2\pi; u]$ 
6:  $u^0 = u \cos u + v \sin u$ 
7:  $v^0 = v \cos u - u \sin u$ 
8: if  $\log \hat{\mu}(u^0) + \log \hat{\nu}(v^0) > \log s$  then
9:    $x^0 = T(u^0)$ 
10:  Return  $x^0, u^0$ 
11: else
12:   if  $u < 0$  then
13:      $u = \min(u, 0)$ 
14:   else
15:      $u = \max(u, 2\pi)$ 
16:   end if
17:    $u \sim \text{Uniform}(u^{\min}; u^{\max})$ 
18:   Go to 6.
19: end if

```

Natarovskii et al. (2021) show that the standard elliptical slice sampler is geometrically ergodic if the target density has tails which are rotationally invariant and monotonically decreasing, e.g.  $\exp(-c|x|)$  for  $c > 0$ . (ii) This implies geometric ergodicity for TESS if for a target density and Markov transition kernel  $P(x; \cdot)$ , with  $C > 0$  and  $\beta \in (0; 1)$ , geometric ergodicity of the elliptical slice sampler holds when

$$\|P^n(x; \cdot) - \pi\|_{TV} \leq C(1 + \|x\|)^n; \quad \forall n \geq 2; \forall x \in \mathbb{R}^d:$$

Then for  $P(x; \cdot)$  the transition kernel of TESS we have,

$$\|P^n(u; \cdot) - \pi\|_{TV} = \|P^n(T(u); \cdot) - \pi\|_{TV} \leq C(1 + \|T(u)\|)^n;$$

which holds only if the transformation  $T$  leads to nice tail properties for  $\hat{\mu}$ . (iii) Following from Theorems 2 and 3 of Johnson and Geyer (2012), if there exists a diffeomorphism  $\tilde{T}$  which ensures that  $\tilde{T}$  pulls in the tails of the distribution enough, then geometric ergodicity holds on the transformed distribution. An open question is to determine the necessary conditions on  $T$  for this result to hold beyond simple transformations.

### 2.4 Adaptive transport maps

There are two key components to TESS, the MCMC sampling phase and the transformation function  $T$ , which so far we have treated as two independent procedures. However, the function  $T$  is parameterised by  $\theta$  and these parameters must be learnt using samples from the target  $\hat{\mu}$ . We therefore propose an adaptive version of TESS which alternates between optimizing and sampling  $T$  to produce an

accurate map between the reference measure and the target distribution.

The parameters are optimized by first running the TESS sampling procedure (Alg. 1) using parallel Monte Carlo chains with an initial value of  $u$ , resulting in  $k$  approximate samples from our target  $\pi(x)$ . We then run  $m$  iterations of a stochastic gradient descent algorithm on the loss function

$$KL(\pi(x) \parallel \hat{\pi}(x)) = \frac{1}{k} \sum_{i=1}^k \log \frac{\pi(x_i)}{\hat{\pi}(x_i)}; \quad (5)$$

The warm-up stage of the sampler repeats this process for  $h$  epochs with batches of size  $k$ , adjusting the inherited parameters from the previous epoch and finally fixing the parameters to then iterate  $m$  times Algorithm 1, generating samples from our extended target space  $\pi(v)$ . This adaptive sampling algorithm is detailed in Algorithm 2.

An important property of the Kullback-Leibler divergence is that it is an asymmetric divergence, i.e.  $KL(\hat{\pi} \parallel \pi) \neq KL(\pi \parallel \hat{\pi})$ . Minimizing  $KL(\hat{\pi} \parallel \pi)$  forces  $\hat{\pi}(x)$  to cover the mass of  $\pi(x)$  thus produces a poor approximation of the tails of the posterior target density. Alternatively, minimizing  $KL(\pi \parallel \hat{\pi})$  forces  $\hat{\pi}(x)$  to cover the mass of  $\pi(x)$ , providing an overconfident approximation to the target density that can be corrected using a sampling method that leaves the target distribution invariant.

**Algorithm 2 Adaptive TESS**

Require:  $u_{1:k}^{(0)}; h; m; N$ ; TESS. TESS applies Algorithm 1

- 1: Set initial parameters  $\alpha$  and  $\wedge$ .
- 2: for  $t = 1; \dots; h$  do . Warm-up
- 3: for  $i = 1; \dots; k$  do
- 4:  $x_i^{(t)}; u_i^{(t)} \leftarrow \text{TESS}(u_i^{(t-1)}; T; \wedge)$
- 5: end for
- 6: Update  $\alpha$  in  $T$  by running  $m$  iterations of gradient descent on (5) using samples  $\{x_{p,k}^{(t)}\}$ .
- 7: end for
- 8:  $u_{1:k}^{(0)} \leftarrow u_{1:k}^{(h)}$
- 9: for  $t = 1; \dots; N$  do . Sampling
- 10: for  $i = 1; \dots; k$  do
- 11:  $x_i^{(t)}; u_i^{(t)} \leftarrow \text{TESS}(u_i^{(t-1)}; T; \wedge)$
- 12: end for
- 13: end for
- 14: Return  $x_{1:k}^{(1)}; \dots; x_{1:k}^{(N)}$

We follow the approach of Hoffman et al. (2019) and initialize the parameters of the NF using an approximation of the parameters that minimize  $KL(\hat{\pi} \parallel \pi)$  via a stochastic gradient descent scheme. In other words, minimizing the Monte Carlo approximation

$$KL(\pi(u) \parallel \hat{\pi}(u)) = \frac{1}{M} \sum_{i=1}^M \log \frac{\pi(u_i)}{\hat{\pi}(u_i)}; \quad u_i \text{ iid } \pi; \quad (6)$$

**2.5 Choice of transport map**

There is a wide class of linear and nonlinear functions which can be used within our normalizing flow map. In this paper, we focus on the coupling architecture first introduced by Dinh et al. (2014). Consider the disjoint partition  $\mathbb{R}^p = (x^A; x^B) \in \mathbb{R}^p \times \mathbb{R}^{d-p}$  and a coupling function  $t(\cdot; \cdot) : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^p \times \mathbb{R}^p$  parameterized by some set of parameters  $\theta$ . Then, one can define a transformation  $\mathcal{G} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  by the formula

$$x^A = t(u^A; (u^B)) := e^{-1} u^A + \alpha_2 \quad (7)$$

$$x^B = u^B; \quad (8)$$

given parameters:  $\mathbb{R}^{d-p} \times \mathbb{R}^p \rightarrow \mathbb{R}^p \times \mathbb{R}^p$  learned only from the extended input. Here we assume an affine bijection, defined in (7), and make a dense feedforward neural network, for further generalizations and variations see Kobayzev et al. (2020). The main practical advantages of the coupling architecture with affine transformations are that it is easily inverted through a shift and scale of the transformation with parameters given by the unchanged  $u^B = u^B$ , and that the modulus determinant of its Jacobian matrix can be easily computed as  $|\det \text{tr } G(x)| = \prod_{i=1}^d (e^{-1})_i$ . Furthermore, since the inverse of the transformation is of similar structure, also its constant of volume change can be easily derived as  $|\det \text{tr } G^{-1}(x)| = \prod_{i=1}^d (e^{-1})_i$ . Both of these are using parameters given by  $(u^B) = (x^B)$  and we drop the absolute value from our computations since the values being multiplied are non-negative. We allow for arbitrary complexity of our NF by introducing a transformation  $\mathcal{D} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with the same structure as  $\mathcal{G}$  but with the roles of the random variables reversed, i.e.  $\alpha^A = u^A$  and  $x^B = t(u^B; (u^A))$ . Hence making our final NF a sequential composition of  $n-1$  transformations  $\mathcal{F} = \mathcal{D}_n \circ \mathcal{G}_n \circ \dots \circ \mathcal{D}_1 \circ \mathcal{G}_1$ .

**3 RELATED WORK**

Elliptical slice sampling The original elliptical slice sampler paper (Murray et al., 2010) presented a simple algorithm that worked well on scenarios of strong prior (Gaussian) information. Nishihara et al. (2014) were the first to explore the idea of generalizing this algorithm to any target distribution, while trying to maintain a simple kernel. Their proposal used a Student-t distribution to approximate the target, under the premise that this proposal would adequately cover the tails of the target density. Their work also considered an adaptive mechanism using parallel computing architectures, which accelerated the MCMC sampler by utilizing multiple chains with fewer iterations per chain. Fagan et al. (2016) also used a generalized elliptical slice sampling proposal paired with a preconditioning step to alleviate complex geometry on their target, in their case using expectation propagation to learn correlation structures for subsets of the parameter space. The main difference between previous

elliptical slice sampling work and our methodology is the algorithm (Hoffman et al., 2019). We compare the effect of use of normalizing flows to create a transport map between a Gaussian density (for which the sampler works well) and the target density of interest. As shown in Section 4, using the richness of nonlinear transport maps produces a fast and highly efficient MCMC algorithm.

Monte Carlo transport maps Our work draws inspiration and is closely related to several threads of work that approach the problem of simulation by simplifying the structure of the target density through a preconditioning step. For general MCMC proposals, Parno and Marzouk (2018) introduced the idea of learning a diffeomorphism using samples from an MCMC algorithm to approximate the target density. Their work built on El Moselhy and Marzouk (2012)'s proposal for approximate inference, adding an MCMC kernel that corrects the approximation and provides asymptotically exact samples. Their work showed that a relatively simple transformation can provide valuable information about the global structure of the target density, thus improving the efficiency of MCMC algorithms that use local gradient information on certain, especially degenerate, test cases.

MCMC with normalizing flows The NeuTra Hamiltonian Monte Carlo (HMC) algorithm introduced in Hoffman et al. (2019) combines neural transport maps with the HMC algorithm. This builds on the earlier work of Marzouk et al. (2016), who frame the approximate inference problem as solving a two-step process, where first an optimization problem is solved to find a preconditioned diffeomorphism which minimizes the integrated autocorrelation time, and then the preconditioned target is sampled from using a HMC algorithm. The NeuTra algorithm relies on gradient-based proposals to explore the target density. A key difference from the transport elliptical slice sampler is that gradients of the target density are not required, this makes the algorithm faster than gradient-based MCMC algorithms, and as illustrated in Section 4, this is achieved without sacrificing sampling accuracy due to transport mapping. Additionally, TESS can be applied in settings where it is either infeasible to calculate target gradients, or they may be unstable (e.g the Neal's funnel density (Gorinova et al., 2020)).

We compare the experimental results of each algorithm based on their Monte Carlo sample efficiency, as indicated by the maximum integrated autocorrelation time ( $\tau_{\text{max}}(\cdot)$ ) with standard deviation ( $\sigma$ ). Additionally, we present the effective sample size (ESS) in terms of the median worst case integrated autocorrelation time, both for individual chains and when all chains are grouped together. A more efficient algorithm is indicated by lower autocorrelations and higher ESS, as this indicates that samples are closer to being independent. To demonstrate the impact of computational cost on each algorithm, we normalize the ESS by the run time in seconds. ESS/sec considers the time spent adapting and sampling, therefore provides a fair comparison between algorithms. To assess the accuracy of the posterior approximation for each algorithm, we use the kernelized Stein discrepancy with U- and V-statistics, as described in (Gorham and Mackey, 2017). Lower values of U- and V-statistics indicate a better approximation of the target posterior. Further information on these diagnostics can be found in the Supplementary Material.

## 4 EXPERIMENTS

### 4.1 Biochemical oxygen demand model

In this section we compare the performance of the adaptive form of TESS (Alg. 2) with the performance of several state-of-the-art MCMC algorithms designed for parallel computer architectures. Specifically, MEADS (Hoffman and Sountsov, 2022), ChEES-HMC (Hoffman et al., 2021) and the popular NUTS algorithm (Hoffman et al., 2014), where an adaptive step size is tuned such that the average cross-chain harmonic-mean acceptance rate is approximately 0.8. We also precondition the latter NUTS method, using the same NF as in TESS, which leads to the NeuTra

We start with an experiment from (Parno and Marzouk, 2018) designed to undermine gradient methods because of its rapidly changing posterior correlation structure, which is challenging for standard samplers to explore. Gradient methods capture local geometry, but the local geometry in this example is not representative of the global geometry of the target and thus provides insufficient information for efficient sampling. On the other hand, the non-linear transformation of the target space with a NF-based approach captures the global, non-Gaussian structure of the target density.

Algorithm	$\tau_{\max}$		ESS	ESS/chain	ESS/sec	Stein U-stat.	Stein V-stat.
TESS	0.555	1.485	11523	90	1129.199	4.269e+02	4.570e+02
MEADS	9.959	1.468	643	5	208.613	1.476e+15	1.486e+15
ChEES-HMC	6.406	2.228	999	8	224.290	1.505e+16	1.510e+16
NUTS	9.579	1.427	668	5	19.625	1.187e+15	1.192e+15
NeuTra	9.553	1.502	670	5	15.579	1.082e+15	1.087e+15

Table 1: Biochemical oxygen demand model. Algorithm diagnostics where  $\tau_{\max}$  is the maximum integrated autocorrelation time over all dimensions; ESS is the corresponding minimum effective sample size. Results are averaged over multiple chains of each sampler, and  $\sigma_{\max}$  is the empirical standard deviation of  $\tau_{\max}$  over these runs.

rejection probabilities, thus being inefficient at producing samples from the posterior.

#### 4.2 Sparse logistic regression

Next, we consider a sparse logistic regression model with hierarchies. Regression parameters of the logistic likelihood are given a horseshoe prior (Carvalho et al., 2009) which induces sparsity on the regressors, i.e. variable selection. These types of hierarchies on the prior scale of a parameter create funnel geometries that are hard to efficiently explore without local or global structure of the target.

Figure 2: Samples from the target density of the Biochemical oxygen demand model acquired by the TESS algorithm, mapped to  $\hat{u}$  (1), with diffeomorphism learned from the warm-up procedure of Algorithm 2. With an approximation that overestimates the real variance of our target (left) we are able to capture its global, non-Gaussian structure and explore it using a dimension independent and gradient-free method.

Algorithms are run on the non-centered parametrization (Papaspiliopoulos et al., 2007) of our model using the numerical version of the German credit dataset. The target posterior is defined by the likelihood  $L(y; \beta; \gamma) = \prod_i \text{Bernoulli}(y_i; (\sigma(\beta^T X_i)))$ , with sigmoid function  $\sigma(\cdot)$ , and prior  $\pi_0(\beta; \gamma) = \text{Gamma}(\beta; 1=2; 1=2) \prod_j N(\gamma_j; 0; 1) \text{Gamma}(\gamma; 1=2; 1=2)$ . Numerical results for each MCMC algorithm are shown in Table 2. Notice how NUTS and NeuTra provide the best results but long sampling times reflect their inefficiency when running in parallel: every iteration takes as long as the longest chain takes to iterate. Waiting for all chains to catch up severely slows down sampling time, the same effect can be observed in all experiments.

The simple biochemical oxygen demand model is given by  $B(t) = \sigma(1 - \exp(-\gamma t))$  for time  $t < 5$ . In this synthetic data experiment, we set the parameters  $\beta = 1$  and  $\gamma = 0:1$  and simulate  $y(t_i)$  observations at times  $t_i$  evenly spaced in  $[0; 5)$  for  $i = 1; \dots; 20$  such that  $y(t_i) = \sigma(1 - \exp(-\gamma t_i)) + \epsilon_i$ , where  $\epsilon_i \sim N(0; \frac{2}{y})$  and  $\text{ced} \frac{2}{y} = 2 \cdot 10^{-4}$ . The target posterior density is given by the likelihood  $L(y; \beta; \gamma) = \prod_i N(y(t_i); B(t_i; \beta; \gamma); \frac{2}{y})$  and at prior  $\pi_0(\beta; \gamma) / 1$ . The numerical results are shown in Table 1 and Figure 2 plots the Monte Carlo approximation of the posterior for the original and transformed densities.

As the dimension of the parameter space grows (51 in this example), TESS will require more samples, i.e. more chains, for a low variance estimate (6). In addition, a more complicated NF is required to capture the non-Gaussian structure of the high-dimensional target space. When either of these fail, and the diffeomorphism is unable to capture the structure of the target space, the simple sampling procedure inherited from the elliptical slice sampler will struggle to sample from the target space, even if producing uncorrelated samples. We purposely illustrate the effect of a deficient transformation on a high dimensional problem in order for the practitioner to understand the caveats of our method. Studying ways to lower the variance of (6), using control variates (Lemieux, 2014) and similar methods (Botev and Ridder, 2017), as well as alternative

It is clear from the results that local gradient information is insufficient to efficiently sample from the rapidly changing local correlation structure of the target density. On the other hand, the learned transport map from the warm-up procedure of TESS provides a mass-covering approximation of the global structure of the target, demonstrated in Figure 2 by  $\hat{u}$ , which allows the algorithm to move farther away from its initial position, exploring efficiently the entire target space, and yielding not only shorter autocorrelation times, but also the correct posterior estimates of the parameter space. In this specific case, gradient-based algorithms are forced to take very small steps while still encountering largest

<https://archive.ics.uci.edu/ml/datasets/german+credit+data>

Algorithm	$\tau_{max}$		ESS	ESS/chain	ESS/sec	Stein U-stat.	Stein V-stat.
TESS	5.182	0.352	1235	10	34.744	1.591e+00	1.693e+00
MEADS	7.105	0.413	901	7	49.453	9.408e-01	1.079e+00
ChEES-HMC	5.666	0.380	1130	9	81.588	1.193e+00	1.312e+00
NUTS	4.734	0.833	1352	11	0.379	1.004e+00	1.138e+00
NeuTra	2.482	1.949	2579	20	0.401	3.618e-01	4.971e-01

Table 2: Sparse logistic regression. Algorithm diagnostics where  $\tau_{max}$  is the maximum integrated autocorrelation time over all dimensions; ESS is the corresponding minimum effective sample size. Results are averaged over multiple chains of each sampler, and  $\sigma_{\tau_{max}}$  is the empirical standard deviation of  $\tau_{max}$  over these runs.

NF schemes is left to future work.

### 4.3 Regime switching Hidden Markov model

A important use of inference and uncertainty quanti ca- tion is on time series data. In this example, we analyze financial time series, speci cally the daily difference in log price data of Google's stock, referred to as returns for  $t = 1; :::; 431$ . We shall assume that at any given time  $t$  the stock's returns will follow one of two regimes: an independent random walk regime  $N(r_t; \sigma^2)$ , or an autoregressive regime  $N(r_t + r_{t-1}; \sigma^2)$ . We define the two regimes as  $s_t \in \{0, 1\}$  and the probability of switching between, or remaining within a regime at time  $t$  will depend on the regime at  $t-1$ , i.e.  $p_{s_t-1; s_t}$  for  $s_t \in \{0, 1\}$ . The transition probabilities  $p_{1;1}$  and  $p_{2;2}$ , and their complementary probabilities  $p_{1;2} = 1 - p_{1;1}$  and  $p_{2;1} = 1 - p_{2;2}$  are treated as model parameters. Since the regime at any time is unobserved, we instead carry over time the probability of belonging to either regime as  $s_t + s_{t-1} = 1$ . Finally, we define the initial values, both for returns  $r_0$  and the probability of belonging to one of the two regimes  $s_0$ .

The regime switching model is defined by the likelihood

$$L(r; \theta; \sigma^2; p; r_0; s_0) = \prod_{t=1}^T \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} (r_t - s_t r_{t-1})^2\right) \quad (9)$$

where  $s_t = \frac{1 + s_{t-1} r_{t-1}}{1 + s_{t-1} r_{t-1} + (1 - s_{t-1}) r_{t-1}}$ ;

and  $j_t = p_{j;1} N(r_t; r_{t-1}; \sigma^2) + p_{j;2} N(r_t; r_{t-1} + r_{t-1}; \sigma^2)$  for  $j \in \{0, 1\}$ . The prior distributions for the parameters are

$$r_0; \sigma; r_0 \sim N(0; 1); \quad N^0(1; 0; 1); \quad (10)$$

$$p_{j;2} \sim C^+(1); \quad (11)$$

$$p_{1;1}; p_{2;2} \sim \text{Beta}(10; 2); \quad s_0 \sim \text{Beta}(2; 2); \quad (12)$$

where  $N^0$  indicates a Gaussian distribution which is truncated at zero and  $C^+$  is the half-Cauchy distribution. Numerical results are shown in Table 3.

The marginal unimodality and somewhat independent correlation structure of the parameters makes this posterior

distribution easy to sample from, diagnostic results show the best performance for all algorithms with respect to other models. TESS's learned reversible transformation of the target density, allowing it to propose uncorrelated sequential samples, is fundamental for its superior diagnostics. ChEES-HMC outputs the samples with the lowest Stein discrepancy, but since it uses the same step size for all target dimensions it struggles to mix well on the worst-case dimension. On the other hand, a reversible transport map is able to capture the covariance structure of the target, allowing fast mixing even on the worst-case dimension. Pair density plots can be found in the Supplementary Material.

### 4.4 Predator-prey system

We consider a likelihood defined as a solution of an ODE system, specifically, the predator-prey system defined by the Lotka-Volterra equations (Goel et al., 1971),

$$\frac{dp}{dt} = p - pq; \quad \text{and} \quad \frac{dq}{dt} = -q + pq; \quad (13)$$

where  $p$  and  $q$  are the prey and predator populations, respectively. We can solve the ODE system of equations numerically and account for measurement error by modelling the observations as  $\log p_t \sim N(\log p(t); \frac{\sigma_p^2}{p(t)})$  and  $\log q_t \sim N(\log q(t); \frac{\sigma_q^2}{q(t)})$  for all  $t > 0$ . Furthermore  $p(0)$  and  $q(0)$  are the initial values. Since we cannot analytically solve the system of equations, we approximate its solution using the Runge-Kutta method, adding an approximation error to our likelihood function. Data for the Hudson's Bay historical lynx-hare population are used as observations in the model. The likelihood is defined as

$$L(p; q; \theta) = \prod_{t=1}^T N\left(\begin{matrix} \log p_t; & \log p(t); & \frac{\sigma_p^2}{p} & 0 \\ \log q_t; & \log q(t); & 0 & \frac{\sigma_q^2}{q} \end{matrix}\right)$$

where  $\theta = (p; q; \sigma_p^2; \sigma_q^2; p(0); q(0))$  and  $f(p(t); q(t); \theta)$  are approximate solutions to the Lotka-Volterra system of equations initialized at  $(p(0); q(0))$ .

The marginal unimodality and somewhat independent correlation structure of the parameters makes this posterior <http://people.whitman.edu/~hundredr/courses/M250F03/LynxHare.txt>

Algorithm	$\tau_{max}$		ESS	ESS/chain	ESS/sec	Stein U-stat.	Stein V-stat.
TESS	0.267	0.893	23985	187	985.969	5.120e-02	1.301e-01
MEADS	1.382	1.197	4631	36	319.949	3.066e-01	3.867e-01
ChEES-HMC	3.451	1.825	1855	14	121.756	-8.203e-03	7.073e-02
NUTS	0.282	0.403	22672	177	182.255	2.222e-02	1.009e-01
NeuTra	0.441	1.020	14530	114	209.069	1.092e-01	1.880e-01

Table 3: Regime switching Hidden Markov model. Algorithm diagnostics where  $\tau_{max}$  is the maximum integrated auto-correlation time over all dimensions; ESS is the corresponding minimum effective sample size. Results are averaged over multiple chains of each sampler, and  $\sigma_{max}$  is the empirical standard deviation of  $\tau_{max}$  over these runs.

Prior distributions for parameters are

$$\mu \sim N^0(1; 1=2); \quad \sigma \sim N^0(1=20; 1=20); \quad (14)$$

$$\log \rho; \log q \sim N(\cdot; 1); \quad (15)$$

$$\log p(0); \log q(0) \sim N(\log 10; 1); \quad (16)$$

where  $N^0$  is a Gaussian distribution truncated at zero.

This experiment exhibits a situation similar to Section 4.1: gradient methods, without global information on the structure of our target, lack enough information to move efficiently around its rapidly changing correlation structure. On the other hand, TESS captures the global structure of the target using a NF and is able to move purposely around it when sampling. Figure 3 illustrates the contrast: MEADS is unable to converge towards a sensible solution, exploring a region of the target space with large error variance and insignificant initial positions, both for the predator and the prey populations; on the other hand, TESS is able to converge towards reasonable initial populations and concentrate sampling around small error variance. Samples from the other gradient methods give similar results to MEADS. Gradient methods need a learned correlation matrix that captures the global correlation structure of the target and use gradient information to propose large steps locally, while TESS is able to capture both the global correlation and local structure by learning an overconfident transport map, then using this information on a cheap and gradient-free method for sampling.

TESS MEADS

Figure 3: Density plots of the approximate posterior distribution for the initial values and scale parameters from the predator-prey system model, drawn with transport elliptical slice sampling on the left and MEADS on the right.

## 5 DISCUSSION

In this paper we proposed TESS, an MCMC algorithm that performs dimension independent and gradient-free sampling from any unnormalized target density. We also proposed an adaptive version of our algorithm that learns a non-Gaussian approximation to the target, helping the algorithm explore complex geometries efficiently. TESS is also able to utilize parallel computer architectures to accelerate sampling from posterior distributions. We believe that this will allow practitioners to perform uncertainty quantification of their models with parallel computational resources and little time.

We found that our algorithm is able to outperform gradient-based competitors in a variety of models. However, it is im-

portant to develop flexible transport maps and low-variance Monte Carlo approximations of the KL divergence, specially for high-dimensional models. Future work will explore the role of the transport map on the algorithm's efficiency, its efficacy in capturing issues in Bayesian posterior geometries, and develop flexible transport maps for high-dimensional models.

### Acknowledgements

The authors would like to thank the anonymous reviewers for their helpful feedback which has significantly improved the quality of the paper. CN gratefully acknowledges the support of EPSRC grants EP/V022636/1, EP/S00159X/1 and EP/R01860X/1.

### References

Iain Murray, Ryan Adams, and David MacKay. Elliptical slice sampling. In Proceedings of the thirteenth international conference on artificial intelligence and statistics



- pages 541–548. JMLR Workshop and Conference Proceedings, 2010.
- Simon Duane, Anthony D Kennedy, Brian J Pendleton, and Duncan Roweth. Hybrid monte carlo. *Physics letters B* 195(2):216–222, 1987.
- Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo* 10(11):2, 2011.
- Peter J Rosicky, Jimmie D Doll, and Harold L Friedman. Brownian dynamics as smart monte carlo simulation. *The Journal of Chemical Physics* 69(10):4628–4633, 1978.
- Ulf Grenander and Michael I Miller. Representations of knowledge in complex systems. *Journal of the Royal Statistical Society: Series B (Methodological)* 56(4):549–581, 1994.
- Seppo Linnainmaa. Taylor expansion of the accumulated rounding error. *BIT Numerical Mathematics* 16(2):146–160, 1976.
- Matthew D Hoffman, Andrew Gelman, et al. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *J. Mach. Learn. Res.* 15(1):1593–1623, 2014.
- John Salvatier, Thomas V Wiecki, and Christopher Fonnesbeck. Probabilistic programming in python using pymc3. *PeerJ Computer Science* 2:e55, 2016.
- Bob Carpenter, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of statistical software* 76(1), 2017.
- Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal probabilistic programming. *Mach. Learn. Res.* 20:28:1–28:6, 2019.
- Du Phan, Neeraj Pradhan, and Martin Jankowiak. Composable effects for flexible and accelerated probabilistic programming in numpyro. *arXiv preprint arXiv:1912.11554* 2019.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning* 37(2): 183–233, 1999.
- Matthew Hoffman, Alexey Radul, and Pavel Soutsov. An adaptive-mcmc scheme for setting trajectory lengths in hamiltonian monte carlo. *International Conference on Artificial Intelligence and Statistics* pages 3907–3915. PMLR, 2021.
- Matthew D Hoffman and Pavel Soutsov. Tuning-free generalized hamiltonian monte carlo. *International Conference on Artificial Intelligence and Statistics* pages 7799–7813. PMLR, 2022.
- Alexey Radul, Brian Patton, Dougal Maclaurin, Matthew Hoffman, and Rif A Saurous. Automatically batching control-intensive programs for modern accelerated. *Proceedings of Machine Learning and Systems* 2:390–399, 2020.
- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensor flow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* 2016.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. Alché-Buc, E. Fox, and R. Garnett, editors *Advances in Neural Information Processing Systems* 32, pages 8024–8035. Curran Associates, Inc., 2019.
- R. Neal. Regression and classification using gaussian process priors. In *Bayesian statistics* volume 6, page 475. 1998.
- Radford M Neal. Slice sampling. *The annals of statistics* 31(3):705–767, 2003.
- Daniilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning* pages 1530–1538. PMLR, 2015.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics* 22(1):79–86, 1951.
- James M. Joyce. Kullback-Leibler Divergence. pages 720–722. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011. ISBN 978-3-642-04898-2. doi: 10.1007/978-3-642-04898-2\_327.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association* 112(518):859–877, 2017.
- Viacheslav Natarovskii, Daniel Rudolf, and Björn Sprungk. Geometric convergence of elliptical slice sampling. In *International Conference on Machine Learning* pages 7969–7978. PMLR, 2021.

- Leif T Johnson and Charles J Geyer. Variable transformation to obtain geometric ergodicity in the random-walk metropolis algorithm. *The Annals of Statistics*, pages 3050–3076, 2012.
- Matthew Hoffman, Pavel Sountsov, Joshua V Dillon, Ian Langmore, Dustin Tran, and Srinivas Vasudevan. Neutralizing bad geometry in hamiltonian monte carlo using neural transport. *arXiv preprint arXiv:1903.03704*, 2019.
- Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.
- Robert Nishihara, Iain Murray, and Ryan P Adams. Parallel mcmc with generalized elliptical slice sampling. *The Journal of Machine Learning Research*, 15(1):2087–2112, 2014.
- Francois Fagan, Jalaj Bhandari, and John P Cunningham. Elliptical slice sampling with expectation propagation. In *UAI*, 2016.
- Matthew D Parno and Youssef M Marzouk. Transport map accelerated markov chain monte carlo. *SIAM/ASA Journal on Uncertainty Quantification*, 6(2):645–682, 2018.
- Tarek A El Moselhy and Youssef M Marzouk. Bayesian inference with optimal maps. *Journal of Computational Physics*, 231(23):7815–7850, 2012.
- Youssef Marzouk, Tarek Moselhy, Matthew Parno, and Alessio Spantini. An introduction to sampling via measure transport. *arXiv preprint arXiv:1602.05023*, 2016.
- Maria Gorinova, Dave Moore, and Matthew Hoffman. Automatic reparameterisation of probabilistic programs. In *International Conference on Machine Learning*, pages 3648–3657. PMLR, 2020.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jackson Gorham and Lester Mackey. Measuring sample quality with kernels. In *International Conference on Machine Learning*, pages 1292–1301. PMLR, 2017.
- Carlos M Carvalho, Nicholas G Polson, and James G Scott. Handling sparsity via the horseshoe. In *Artificial Intelligence and Statistics*, pages 73–80. PMLR, 2009.
- Omiros Papaspiliopoulos, Gareth O Roberts, and Martin Sköld. A general framework for the parametrization of hierarchical models. *Statistical Science*, pages 59–73, 2007.
- Christiane Lemieux. Control variates. *Wiley StatsRef: Statistics Reference Online*, pages 1–8, 2014.
- Zdravko Botev and Ad Ridder. Variance reduction. *Wiley statsRef: Statistics reference online*, pages 1–6, 2017.
- Narendra S Goel, Samaresh C Maitra, and Elliott W Montroll. On the volterra and other nonlinear models of interacting populations. *Reviews of modern physics*, 43(2):231, 1971.
- Ulli Wolff, Alpha Collaboration, et al. Monte carlo errors with less errors. *Computer Physics Communications*, 156(2):143–153, 2004.
- Qiang Liu, Jason Lee, and Michael Jordan. A kernelized stein discrepancy for goodness-of-fit tests. In *International conference on machine learning*, pages 276–284. PMLR, 2016.
- Charles Stein, Persi Diaconis, Susan Holmes, and Gesine Reinert. Use of exchangeable pairs in the analysis of simulations. *Lecture Notes-Monograph Series*, pages 1–26, 2004.

## A ELLIPTICAL SLICE SAMPLING ALGORITHM

---

**Algorithm 3** Elliptical slice sampler (Murray et al., 2010)

---

**Require:**  $x; L(Dj)$

- 1:  $v \sim N(0; I_d)$
- 2:  $w \sim \text{Uniform}(0; 1)$
- 3:  $\log s \sim \log L(Djx) + \log w$
- 4:  $\theta \sim \text{Uniform}(0; 2\pi)$
- 5:  $[min; max] \sim [ -2; 2 ]$
- 6:  $x^0 = x \cos \theta + v \sin \theta$
- 7: **if**  $\log L(Djx^0) > \log s$  **then**
- 8:     Return  $x^0$
- 9: **else**
- 10:    **if**  $\theta < 0$  **then**
- 11:      $min \leftarrow min - \theta$
- 12:    **else**
- 13:      $max \leftarrow max + \theta$
- 14:    **end if**
- 15:     $\theta \sim \text{Uniform}(min; max)$
- 16:    Go to 6.
- 17: **end if**

---

## B PROOF

### B.1 Proof of Proposition 1

As established in Murray et al. (2010) and Nishihara et al. (2014), the elliptical slice sampler and generalized elliptical sampler target the correct stationary distribution as the algorithm is reversible and produces an irreducible, aperiodic Markov chain.

The same result holds for the TESS algorithm from initial state  $u = T^{-1}(x)$  and where  $(u; v)$  and  $(u^0; v^0)$  represent the initial and accepted transformed parameters of the sampler (steps 1 and 6-7), with  $s$  the slice variable (step 3) and  $f_{k=1}^K g_{k=1}^K$  the parameters representing points in the slice expressed in radians until acceptance at  $K$  (steps 4 and 17). Let

$$f_{k=1}^K = \begin{cases} f_{k=1}^K & \text{if } k < K \\ f_{k=1}^K & \text{if } k = K \end{cases}; \quad (17)$$

then by the properties of the elliptical slice sampler, the transformation  $(u; v; s; f_{k=1}^K g_{k=1}^K) \mapsto (u^0; v^0; s; f_{k=1}^0 g_{k=1}^0)$  is bijective, preserves volume and  $p(f_{k=1}^K g_{k=1}^K; u; v; s) = p(f_{k=1}^0 g_{k=1}^0; u^0; v^0; s)$ . Using the uniform density of the slice variable  $s$  it is easy to see that  $p(u^0; v^0; f_{k=1}^0 g_{k=1}^0; s; j u^0; v^0) \wedge (u^0) \wedge (v^0) = p(u; v; f_{k=1}^K g_{k=1}^K; s; j u; v) \wedge (u) \wedge (v)$ , and so if  $(u; v) \wedge (u) \wedge (v)$  then  $(u^0; v^0) \wedge (u^0) \wedge (v^0)$ . Finally, as  $x = T(u)$  we have  $(x^0; v^0) \wedge (T(u^0)) j \det T(u^0) \wedge (v^0) = (x^0; v^0) \wedge (v^0)$ .

## C DIAGNOSTIC TOOLS CALCULATION DETAILS

Here we describe the calculation of the maximum integrated autocorrelation time  $\tau_{\max}$  and Kernelized Stein discrepancy U- and V- statistics used throughout our results. Assume we have as output from chain  $c$  a sequence of  $N$  samples from our target  $\mathbf{X}_{c,1}; \dots; \mathbf{X}_{c,N}$ , where each sample is on a  $d$  dimensional parameter space. Then, compute the integrated

autocorrelation time for dimension  $j = 1, \dots, d$  on chain  $c$  as

$$c_j = \frac{1}{2} + 2 \sum_{t=1}^{N-1} \frac{1}{N} \frac{\hat{C}(t)}{2\hat{C}(0)} \quad (18)$$

$$\hat{C}(t) = \frac{1}{N-t} \sum_{0 < i < N-t} (x_{c;i;j} - x_{c;j})(x_{c;i+t;j} - x_{c;j}) \quad (19)$$

$$x_{c;j} = \frac{1}{N} \sum_{i=1}^N x_{c;i;j} \quad (20)$$

The value  $\hat{C}(t)$  for all  $t = 1, \dots, N-1$  is computed by applying the Fourier transform method from Wolff et al. (2004). We then define  $c_{\max}$  and ESS as

$$c_{\max} = \max_{j=1, \dots, d} \text{median}_{c=1, \dots, C} c_j \quad (21)$$

$$\text{ESS} = \min_{j=1, \dots, d} \text{median}_{c=1, \dots, C} \frac{N}{2 c_j} \quad (22)$$

The Kernelized Stein discrepancy's U- and V-statistics are calculated using the inverse multi-quadratic kernel  $k(x; x^0) = (1 + (x - x^0)^T(x - x^0))^{-1/2}$  with  $\beta = 1/2$  on all experiments as

$$\text{U-stat} = \frac{1}{C^2 N(N-1)} \sum_{c,i} \sum_{c^0 \neq c; i^0 \neq i} A A^0 K(\mathbf{x}_{c;i}; \mathbf{x}_{c^0;i^0}) \quad (23)$$

$$\text{V-stat} = \frac{1}{C^2 N^2} \sum_{c,i} \sum_{c^0; i^0} A A^0 K(\mathbf{x}_{c;i}; \mathbf{x}_{c^0;i^0}) \quad (24)$$

$$A A^0 K(x; x^0) = r_x r_{x^0} k(x; x^0) + r_x k(x; x^0) r_{x^0} \log(x^0) + r_{x^0} k(x; x^0) r_x \log(x) + k(x; x^0) r_x \log(x) r_{x^0} \log(x) \quad (25)$$

It can be shown that the U-statistic is an unbiased estimate of  $E_{x; x^0} [A A^0 K(x; x^0)]$  for process generating the samples, while the V-statistic is biased but always non-negative (Liu et al., 2016). If  $\beta = 1/2$  then  $E_{x; x^0} [A A^0 K(x; x^0)] = 0$  by Stein's identity (Stein et al., 2004).

