

---

# Local-Global MCMC kernels: the best of both worlds

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Recent works leveraging learning to enhance sampling have shown promising results, in particular by designing effective non-local moves and global proposals. However, learning accuracy is inevitably limited in regions where little data is available such as in the tails of distributions as well as in high-dimensional problems. In the present paper we study an Explore-Exploit Markov chain Monte Carlo strategy (Ex<sup>2</sup>MCMC) that combines local and global samplers showing that it enjoys the advantages of both approaches. We prove  $V$ -uniform geometric ergodicity of Ex<sup>2</sup>MCMC without requiring a uniform adaptation of the global sampler to the target distribution. We also compute explicit bounds on the mixing rate of the Explore-Exploit strategy under realistic conditions. Moreover, we propose an adaptive version of the strategy (FlEx<sup>2</sup>MCMC) where a normalizing flow is trained while sampling to serve as a proposal for global moves. We illustrate the efficiency of Ex<sup>2</sup>MCMC and its adaptive version on classical sampling benchmarks as well as in sampling high-dimensional distributions defined by Generative Adversarial Networks seen as Energy Based Models.

## 1 Introduction

We consider the setting where a target distribution  $\pi$  on a measurable space  $(\mathbb{X}, \mathcal{X})$  is known up to a normalizing constant and one tries to estimate the expectations of some function  $f : \mathbb{X} \rightarrow \mathbb{R}$  with respect to  $\pi$ . Examples include the extraction of Bayesian statistics from posterior distributions derived from observations as well as the computation of observables of a physical system  $x \in \mathbb{X}$  under the Boltzmann distribution with non-normalized density  $\pi(x) = e^{-\beta U(x)}$  for the energy function  $U$  at the inverse temperature  $\beta$ .

A common strategy to tackle this estimation is to resort to Markov chain Monte Carlo algorithms (MCMCs). The MCMC approach aims to simulate a realization of a time-homogeneous Markov chain  $\{Y_n, n \in \mathbb{N}\}$ , such that the distribution of the  $n$ -th iterate  $Y_n$  with  $n \rightarrow \infty$  is arbitrarily close to  $\pi$ , regardless of the initial distribution of  $Y_0$ . In particular, the Metropolis-Hastings kernel (MH) is the cornerstone of MCMC simulations, with a number of successful variants following the process of a *proposal* step followed by an *accept/reject* step (see e.g. [60]). In large dimensions, proposal distributions are typically chosen to generate local moves that depend on the last state of the chain in order to guarantee an admissible acceptance rate. However, local samplers suffer from long mixing times as exploration is inherently slow, and mode switching, when there is more than one, can be extremely infrequent.

On the other hand, independent proposals are able to generate more global updates, but they are difficult to design. Developments in deep generative modelling, in particular versatile autoregressive and normalising flows [38, 36, 20, 54], spurred efforts to use learned probabilistic models to improve the exploration ability of MCMC kernels. Among a rapidly growing body of work, references include [35, 2, 52, 25, 33]. While these works show that global moves in a number of practical problems can be successfully informed by machine learning models, it remains the case that the acceptance rate of

39 independent proposals decreases dramatically with dimensions – except in the unrealistic case that  
 40 they perfectly reproduce the target. This is a well-known problem in the MCMC literature [11, 68, 1],  
 41 and it was recently noted that deep learning-based suggestions are no exception in works focusing on  
 42 physical systems [19, 45].

43 In this paper we focus on the benefits of combining local and global samplers. Intuitively, local  
 44 steps interleaved between global updates from an independent proposal (learned or not) increase  
 45 accuracy by allowing accurate sampling in tails that are not usually well handled by the independent  
 46 proposal. Also, mixing time is usually improved by the local-global combination, which prevents  
 47 long chains of consecutive rejections. Here we focus on a global kernel of type iterative-sampling  
 48 importance resampling (i-SIR) [70, 4, 5]. This kernel uses multiple proposals in each iteration to take  
 49 full advantage of modern parallel computing architectures. For local samplers, we consider common  
 50 techniques such as Metropolis Adjusted Langevin (MALA) and Hamiltonian Monte Carlo (HMC).  
 51 We call this combination strategy Explore-Exploit MCMC (Ex<sup>2</sup>MCMC) in the following.

52 **Contributions** The main contributions of the paper are as follows:

- 53 • We provide theoretical bounds on the accuracy and convergence speed of Ex<sup>2</sup>MCMC strategies.  
 54 In particular, we prove  $V$ -uniform geometric convergence of Ex<sup>2</sup>MCMC under assumptions much  
 55 milder than those required to prove uniform geometric ergodicity of the global sampler i-SIR alone.
- 56 • We propose an adaptive version of the strategy, called FlEx<sup>2</sup>MCMC, which involves learning an  
 57 efficient proposal while sampling, as in adaptive MCMC.
- 58 • We perform a numerical evaluation of Ex<sup>2</sup>MCMC and FlEx<sup>2</sup>MCMC for various sampling prob-  
 59 lems, including sampling GANs as energy-based models. The results clearly show the advantages  
 60 of the proposed approaches compared to purely local or purely global MCMC methods.

61 **Notations** Denote  $\mathbb{N}^* = \mathbb{N} \setminus \{0\}$ . For a measurable function  $f : \mathbb{X} \mapsto \mathbb{R}$ , we define  $|f|_\infty =$   
 62  $\sup_{x \in \mathbb{X}} |f(x)|$  and  $\pi(f) := \int_{\mathbb{X}} f(x)\pi(dx)$ . For a function  $V : \mathbb{X} \mapsto [1, \infty)$  we introduce the  
 63  $V$ -norm of two probability measures  $\xi$  and  $\xi'$  on  $(\mathbb{X}, \mathcal{X})$ ,  $\|\xi - \xi'\|_V := \sup_{|f(x)| \leq V(x)} |\xi(f) - \xi'(f)|$ .  
 64 If  $V \equiv 1$ ,  $\|\cdot\|_1$  is equal to the total variation distance (denoted  $\|\cdot\|_{TV}$ ).

## 65 2 Explore-Exploit Samplers

66 Suppose we are given a target distribution  $\pi$  on a measurable space  $(\mathbb{X}, \mathcal{X})$  that is known only up to  
 67 a normalizing constant. We will often assume that  $\mathbb{X} = \mathbb{R}^d$  or a subset thereof. Two related problems  
 68 are sampling from  $\pi$  and estimating integrals of a function  $f : \mathbb{X} \mapsto \mathbb{R}$  w.r.t.  $\pi$ , i.e.,  $\pi(f)$ . Among  
 69 the many methods devoted to solving these problems, there is a popular family of techniques based  
 70 on *Importance Sampling* (IS) and relying on independent proposals, see e.g. [1, 71]. We first give a  
 71 brief overview of IS, to describe the global sampler i-SIR. We recall ergodicity results for the latter  
 72 before investigating the Explore-Exploit sampling strategy which couples the global sampler with a  
 73 local kernel. Then we present the main theoretical result of the paper on the ergodicity of the coupled  
 74 strategy.

### 75 2.1 From Importance Sampling to i-SIR

76 The primary purpose of IS is to approximate integrals of the form  $\pi(f)$ . Its main instrument is a  
 77 (known) *proposal distribution*, which we denote by  $\lambda(dx)$ . To describe the algorithm, we assume  
 78 that  $\pi(dx) = w(x)\lambda(dx)/\lambda(w)$ . In this formula,  $w(x)$  is the *importance weight* function assumed  
 79 to be known and positive, i.e.,  $w(x) > 0$  for all  $x \in \mathbb{X}$ , and  $\lambda(w)$  is the *normalizing constant* of  
 80 the distribution  $\pi$ . Typically  $\lambda(w)$  is unknown. If we assume that  $\pi$  and  $\lambda$  have positive densities  
 81 w.r.t. a common dominant measure, denoted also by  $\pi$  and  $\lambda$  respectively, then the *self-normalized*  
 82 *importance sampling* (SNIS, see [59]) estimator of  $\pi(f)$  is given by

$$\hat{\pi}_N(f) = \sum_{i=1}^N \omega_N^i f(X^i), \quad (1)$$

83 where  $X^{1:N} \stackrel{\text{i.i.d.}}{\sim} \lambda$ , and  $\omega_N^i = w(X^i)/\sum_{j=1}^N w(X^j)$  are the self-normalized importance weights.  
 84 Note that computing  $\omega_N^i$  does not require the knowledge of  $\lambda(w)$ . The main problem in the  
 85 practical applications of IS is the choice of the proposal distribution  $\lambda$ . The representation  
 86  $\pi(dx) = w(x)\lambda(dx)/\lambda(w)$  implies that the support of  $\lambda$  covers the support of  $\pi$ . At the same  
 87 time, too large variance of  $\lambda$  is obviously detrimental to the quality of (1). This suggests *adaptive im-*  
 88 *portance sampling* techniques (discussed in [16]), which involve learning the proposal  $\lambda$  to improve  
 89 the quality of (1). We return to this idea in section 3.

---

**Algorithm 1:** Single stage of i-SIR algorithm with independent proposals

---

```

1 Procedure i-SIR ( $Y_k, \lambda$ ):
2   Input : Previous state  $Y_k$ ; proposal distribution  $\lambda$ ;
3   Output: New state  $Y_{k+1}$ ; pool of proposals  $X_{k+1}^{2:N} \sim \lambda$ ;
4   Set  $X_{k+1}^1 = Y_k$ , draw  $X_{k+1}^{2:N} \sim \lambda$ ; for  $i \in [N]$  do
5      $\lfloor$  compute the normalized weights  $\omega_{i,k+1} = w(X_{k+1}^i) / \sum_{\ell=1}^N w(X_{k+1}^\ell)$ ;
6   Draw the proposal index  $I_{k+1} \sim \text{Cat}(\omega_{1,k+1}, \dots, \omega_{N,k+1})$ ;
7   Set  $Y_{k+1} := X_{k+1}^{I_{k+1}}$ .

```

---

90 IS-based techniques can also be used to draw an (approximate) sample from  $\pi$ . For instance,  
91 Sampling Importance Resampling (SIR, [66]) follows the steps:

- 92 1. Draw  $X^{1:N} \stackrel{\text{i.i.d.}}{\sim} \lambda$ ;
- 93 2. Compute the self-normalized importance weights  $\omega_N^i = w(X^i) / \sum_{\ell=1}^N w(X^\ell)$ ,  $i \in \{1, \dots, N\}$ ;
- 94 3. Select  $M$  samples  $Y^{1:M}$  from the set  $X^{1:N}$  choosing  $X^i$  with probability  $\omega_N^i$  with replacement.

95 The drawback of the procedure is that it is only asymptotically valid with  $N \rightarrow \infty$ . Alternatively,  
96 SIR can be repeated to define a Markov Chain as in *iterated SIR* (i-SIR), proposed in [70] and  
97 also studied in [4, 42, 41, 5]. At each iteration of i-SIR described in Algorithm 1, a candidate pool  
98  $X_{k+1}^{2:N}$  is sampled from the proposal and the next state  $Y_{k+1}$  is chosen among the candidates and the  
99 previous state  $X_{k+1}^1 = Y_k$  according to the importance weights. i-SIR shares similarities with the  
100 Multiple-try Metropolis (MTM) algorithm [43], but is computationally simpler and exhibits more  
101 favorable mixing properties; see Appendix A.1. The Markov chain  $\{Y_k, k \in \mathbb{N}\}$  generated by i-SIR  
102 has the following Markov kernel

$$P_N(x, A) = \int \delta_x(dx^1) \sum_{i=1}^N \frac{w(x^i)}{\sum_{j=1}^N w(x^j)} 1_A(x^i) \prod_{j=2}^N \lambda(dx^j).$$

103 Interpreting i-SIR as a systematic-scan two-stage Gibbs sampler (see Appendix A.2 for more details),  
104 it follows easily that the Markov kernel  $P_N$  is reversible w.r.t. the target  $\pi$ , Harris recurrent and  
105 ergodic (see Theorem 5). Provided also that  $|w|_\infty < \infty$ , it was shown in [5] that the Markov kernel  
106  $P_N$  is uniformly geometrically ergodic. Namely, for any initial distribution  $\xi$  on  $(\mathbb{X}, \mathcal{X})$  and  $k \in \mathbb{N}$ ,

$$\|\xi P_N^k - \pi\|_{\text{TV}} \leq \kappa_N^k \quad \text{with } \kappa_N = \frac{N-1}{2L+N-2}, L = |w|_\infty / \lambda(w), \text{ and } \kappa_N = 1 - \epsilon_N. \quad (2)$$

107 We provide a simple direct proof of (2) in Appendix B.1. Yet, note that the bound (2) relies  
108 significantly on the restrictive condition that weights are uniformly bounded  $|w|_\infty < \infty$ . Moreover,  
109 even when this condition is satisfied, the rate  $\kappa_N$  can be close to 1 when the dimension  $d$  is large.<sup>1</sup>  
110 To illustrate this phenomenon, we consider a simple problem of sampling from the standard normal  
111 distribution  $\mathcal{N}(0, I_d)$  with the proposal  $\mathcal{N}(0, 2I_d)$  in increasing dimensions  $d$  up to 300. Results  
112 visualized in Figure 1 show that the performance of vanilla i-SIR quickly deteriorates as most  
113 proposals get rejected. We propose to overcome this problem using the Explore-Exploit strategy  
114 coupling i-SIR with local MCMC steps to define a new sampler.

## 115 2.2 Coupling with local kernels: Ex<sup>2</sup>MCMC

116 After each i-SIR step, we apply a local MCMC kernel  $R$  (rejuvenation kernel), with an invariant  
117 distribution  $\pi$ . We call this strategy Ex<sup>2</sup>MCMC because it combines steps of exploration by i-SIR  
118 and steps of exploitation by the local MCMC moves. The resulting algorithm, formulated in  
119 Algorithm 2, defines a Markov chain  $\{Y_j, j \in \mathbb{N}\}$  with Markov kernel  $K_N(x, \cdot) = P_N R(x, \cdot) =$   
120  $\int P_N(x, dy) R(y, \cdot)$ . The first simple experiments previously considered already illustrates the  
121 advantageous mixing of Ex<sup>2</sup>MCMC with MALA applied as  $R$  (see Figure 1).

122 We now present the main theoretical result of this paper on the properties of Ex<sup>2</sup>MCMC. Under  
123 rather weak conditions, provided that  $R$  is geometrically regular (see [21, Chapter 14]), it is possible  
124 to establish that Ex<sup>2</sup>MCMC remains  $V$ -uniformly geometrically ergodic even if the weight function  
125  $w(x)$  is unbounded.

<sup>1</sup>Indeed, consider a simple scenario  $\pi(x) = \prod_{i=1}^d p(x_i)$  and  $\lambda(x) = \prod_{i=1}^d q(x_i)$  for some densities  $p(\cdot)$  and  $q(\cdot)$  on  $\mathbb{R}$ . Then it is easy to see that  $L = (\sup_{y \in \mathbb{R}} p(y)/q(y))^d$  grows exponentially with  $d$ .

---

**Algorithm 2:** Single stage of Ex<sup>2</sup>MCMC algorithm with independent proposals

---

1 **Procedure** Ex<sup>2</sup>MCMC ( $Y_k, \lambda, R$ ):  
 Input : Previous state  $Y_k$ ; proposal distribution  $\lambda$ ; rejuvenation kernel  $R$ ;  
 Output : New sample  $Y_{k+1}$ ; pool of proposals  $X_{k+1}^{2:N} \sim \lambda$ ;  
 2    $Z_{k+1}, X_{k+1}^{2:N} = i\text{-SIR}(Y_k, \lambda)$ ;  
 3   Draw  $Y_{k+1} \sim R(Z_{k+1}, \cdot)$ .

---

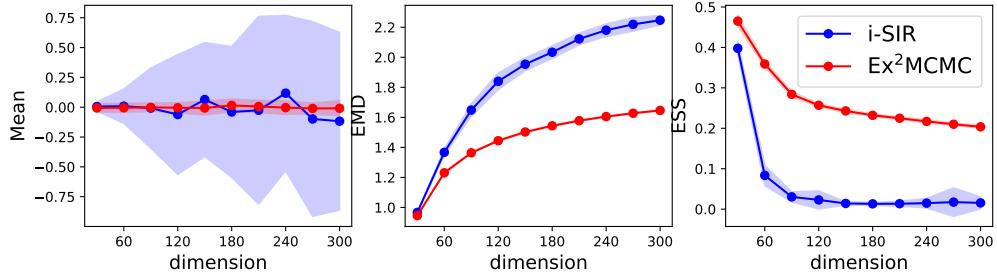


Figure 1: Sampling from  $\mathcal{N}(0, I_d)$  with the proposal  $\mathcal{N}(0, 2 I_d)$ . – See Appendix E.1 for the definitions of EMD and ESS metrics. We display confidence intervals for i-SIR and Ex<sup>2</sup>MCMC obtained from 100 independent runs as blue and red regions, respectively. Ex<sup>2</sup>MCMC helps to achieve efficient sampling even in high dimensions.

126 **Definition 1** ( $V$ -Geometric Ergodicity). A Markov kernel  $Q$  with invariant probability measure  $\pi$  is  
 127  $V$ -geometrically ergodic if there exist constants  $\rho \in (0, 1)$  and  $M < \infty$  such that, for all  $x \in \mathbb{X}$  and  
 128  $k \in \mathbb{N}$ ,  $\|Q^k(x, \cdot) - \pi\|_V \leq M \{\pi(V) + \pi(V)\} \rho^k$ .

129 In particular,  $V$ -geometric ergodicity ensures that the distribution of the  $k$ -th iterate of a Markov chain  
 130 converges geometrically fast to the invariant probability in  $V$ -norm, for all starting points  $x \in \mathbb{X}$ .  
 131 Here the dependence on the initial state  $x$  appears on the right-hand side only in  $V(x)$ . Denote  
 132 by  $\text{Var}_\lambda[w] = \int \{w(x) - \lambda(w)\}^2 \lambda(dx)$  the variance of the importance weight functions under the  
 133 proposal distribution and consider the following assumptions:

134 **A1.** (i)  $R$  has  $\pi$  as its unique invariant distribution; (ii) There exists a function  $V: \mathbb{X} \rightarrow [1, \infty)$ , such  
 135 that for all  $r \geq r_R > 1$  there exist  $\lambda_{R,r} \in [0, 1)$ ,  $b_{R,r} < \infty$ , such that  $RV(x) \leq \lambda_{R,r} V(x) + b_{R,r} 1_{V_r}$ ,  
 136 where  $V_r = \{x: V(x) \leq r\}$ ;

137 **A2.** (i) For all  $r \geq r_R$ ,  $w_{\infty,r} := \sup_{x \in V_r} \{w(x)/\lambda(w)\} < \infty$  and (ii)  $\text{Var}_\lambda[w]/\{\lambda(w)\}^2 < \infty$ .

138 A1-(ii) states that  $R$  satisfies a Foster-Lyapunov drift condition for  $V$ . This condition is fulfilled by  
 139 most classical MCMC kernels - like Metropolis-Adjusted Langevin (MALA) algorithm or Hamil-  
 140 tonian Monte Carlo (HMC), typically under tail conditions for the target distribution; see [61, 22],  
 141 and [21, Chapter 2] with the references therein. A2-(i) states that the (normalized) importance weights  
 142  $w(\cdot)/\lambda(w)$  are upper bounded on level sets of  $V_r$ . This is a mild condition: if  $\mathbb{X} = \mathbb{R}^d$ , and  $V$  is  
 143 norm-like, then the level sets  $V_r$  are compact and  $w(\cdot)$  is bounded on  $V_r$  as soon as  $\pi$  and  $\lambda$  are  
 144 positive and continuous. A2-(ii) states that the variance of the importance weights is bounded; note  
 145 that this variance is also equal to the  $\chi^2$ -distance between the proposal and the target distributions  
 146 which plays a key role in the non-asymptotic analysis of the performance of IS methods [1, 67]. We  
 147 emphasize that these assumptions do not identify the small sets of the rejuvenation kernel  $R$ .

148 **Theorem 2.** Assume A1 and A2. Then, for all  $x \in \mathbb{X}$  and  $k \in \mathbb{N}$ ,

$$\|K_N^k(x, \cdot) - \pi\|_V \leq c_{K_N} \{\pi(V) + V(x)\} \tilde{\kappa}_{K_N}^k, \quad (3)$$

149 where the constant  $c_{K_N}$ ,  $\tilde{\kappa}_{K_N} \in [0, 1]$  are given in the proof. In addition,  $c_{K_\infty} = c_{K_N} + O(N^{-1})$   
 150 and  $\tilde{\kappa}_{K_N} = \tilde{\kappa}_{K_\infty} + O(N^{-1})$  with explicit expressions provided in (13).

151 The proof of Theorem 2 is provided in Appendix B.2. We stress that in many situations, the mixing  
 152 rate  $\tilde{\kappa}_{K_N}$  of the Ex<sup>2</sup>MCMC Markov Kernel  $K_N$  is significantly better than the corresponding mixing  
 153 rate of the local kernel  $R$ , provided  $N$  is large enough. This is illustrated in Appendix C with the  
 154 Metropolis Adjusted Langevin Algorithm (MALA) kernel (see, e.g. [13, 64]).

155 **3 Adaptive version: FlEx<sup>2</sup>MCMC**

156 As already mentioned, the performance of proposal-based samplers depends on the distribution  
 157 of importance weights which is related to the similarity of the proposal and target distributions<sup>2</sup>.  
 158 Therefore, yet another strategy to improve sampling performance is to select the proposal distribution  
 159  $\lambda$  from a family of parameterized distributions  $\{\lambda_\theta\}$  and fit the parameter  $\theta \in \Theta = \mathbb{R}^q$  to the target  
 160  $\pi$ , for example, by minimizing a Kullback-Leibler divergence (KL) [56, 2, 49] or matching moments  
 161 [57]. In *adaptive MCMCs*, parameter adaptation is performed along the MCMC run [6, 9, 62]. In this  
 162 section we propose an adaptive version of Ex<sup>2</sup>MCMC, which we call FlEx<sup>2</sup>MCMC.

163 **Normalizing flow proposal.** A flexible way to parameterize proposal distributions is to combine  
 164 a tractable distribution  $\varphi$  with an invertible parameterized transformation. Let  $T : \mathbb{X} \mapsto \mathbb{X}$  be a  
 165  $C^1$  diffeomorphism. We denote by  $T\#\varphi$  the push-forward of  $\varphi$  under  $T$ , that is, the distribution  
 166 of  $X = T(Z)$  with  $Z \sim \varphi$ . Assuming that  $\varphi$  has a p.d.f. (also denoted  $\varphi$ ), the corresponding  
 167 push-forward density (w.r.t. the Lebesgue measure) is given by  $\lambda_T(y) = \varphi(T^{-1}(y)) J_{T^{-1}}(y)$ , where  
 168  $J_T$  denotes the Jacobian determinant of  $T$ . The parameterized family of diffeomorphisms  $\{T_\theta\}_{\theta \in \Theta}$   
 169 defines a family of distributions  $\{\lambda_{T_\theta}\}_{\theta \in \Theta}$ , denoted for simplicity as  $\{\lambda_\theta\}_{\theta \in \Theta}$ . This construction is  
 170 called a *normalizing flow* (NF) and a great deal of work has been devoted to ways of parameterizing  
 171 invertible flows  $T_\theta$  with neural networks; see [39, 54] for reviews.

172 **Simultaneous learning and sampling.** As with adaptive MCMC methods, we learn the parameters of  
 173 a NF proposal for the global proposal during sampling. We work with  $M$  copies of the Markov chains  
 174  $\{(Y_k[j], X_k^{1:N}[j])\}_{k \in \mathbb{N}^*}$  indexed by  $j \in \{1, \dots, M\}$ . At each step  $k \in \mathbb{N}^*$ , each copy is sampled as  
 175 in Ex<sup>2</sup>MCMC using the NF proposal, independently from the other copies, but conditionally to the  
 176 the current value of the parameters  $\theta_{k-1}$ , i.e.

$$X_k^1[j] = Y_{k-1}[j] \text{ and } X_k^\ell[j] = T_{\theta_{k-1}}(Z_k^\ell[j]), \ell \in \{2, \dots, N\} \text{ where } \{Z_k^{1:N}[j]\} \stackrel{\text{i.i.d.}}{\sim} \varphi.$$

177 We then adapt the parameters by taking steps of gradient descent on a convex combination of  
 178 the *forward* KL,  $\text{KL}(\pi || \lambda_\theta) = \int \pi(x) \log(\pi(x)/\lambda_\theta(x)) dx$  and the *backward* KL  $\text{KL}(\lambda_\theta || \pi) =$   
 179  $\int \lambda_\theta(x) \log(\pi(x)/\lambda_\theta(x)) dx = \int \varphi(z) \log w_\theta \circ T_\theta(z) dz$ . Let  $\{\gamma_k, k \in \mathbb{N}\}$  be a sequence of nonneg-  
 180 ative stepsizes and  $\{\alpha_k, k \in \mathbb{N}\}$  be a nondecreasing sequence in  $[0, 1]$  with  $\alpha_\infty = \lim_{k \rightarrow \infty} \alpha_k$ . The  
 181 update rule is  $\theta_k = \theta_{k-1} + \gamma_k M^{-1} \sum_{j=1}^M H(\theta_{k-1}, X_k^{1:N}[j], Z_k^{2:N}[j])$  where  $H(\theta, x^{1:N}, z^{2:N}) =$   
 182  $\alpha_k H^f(\theta, x^{1:N}) + (1 - \alpha_k) H^b(\theta, z^{2:N})$  with

$$H^f(\theta, x^{1:N}) = \sum_{\ell=1}^N \frac{w_\theta(x^\ell)}{\sum_{i=1}^N w_\theta(x^i)} \nabla_\theta \log \lambda_\theta(x^\ell), \quad w_\theta(x) = \pi(x)/\lambda_\theta(x), \quad (4)$$

$$H^b(\theta, z^{2:N}) = -\frac{1}{N-1} \sum_{\ell=2}^N \{\nabla_\theta \log \pi \circ T_\theta(z^\ell) + \nabla_\theta \log J_{T_\theta}(z^\ell)\}. \quad (5)$$

183 Note that we use a Rao-Blackwellized estimator of the gradient of the forward KL (4) where we  
 184 fully recycle all the  $N$  candidates sampled at each iteration of i-SIR. The quality of this estimator is  
 185 expected to improve along the iterations  $k$  of the algorithm as the variance of importance weights  
 186 decreases as the proposal improves. Note also that using only gradients from the backward KL (5) is  
 187 prone to mode-collapse [56, 53, 49, 25], hence the need for also using gradients from the forward KL  
 188  $H^f(\theta, x^{1:N})$ , which requires the simultaneous sampling from  $\pi$ .

189 Since the parameters of the Markov kernel  $\theta_k$  are updated using samples  $X_k^{1:N}$  from the chain,  
 190  $((Y_k, X_k^{1:N}))_{k \in \mathbb{N}}$  is no longer Markovian. This type of problems has been considered in [47, 12, 30, 7]  
 191 and to prove convergence of the strategy we need to strengthen assumptions compared to the previous  
 192 section.

193 **A3.** *There exists a function  $W : \mathbb{X} \rightarrow \mathbb{R}_+$  such that  $\varphi(W^2) = \int W^2(z) \varphi(dz) < \infty$ , and a*  
 194 *constant  $L < \infty$  such that, for all  $\theta, \theta' \in \Theta$  and  $z \in \mathbb{X}$ ,  $\|\nabla_\theta \log \pi \circ T_\theta(z) - \nabla_\theta \log \pi \circ T_{\theta'}(z)\| \leq$*   
 195  *$L\|\theta - \theta'\|W(z)$  and  $\|\nabla_\theta \log J_{T_\theta}(z) - \nabla_\theta \log J_{T_{\theta'}}(z)\| \leq L\|\theta - \theta'\|W(z)$ .*

196 **A4.** *(i) For all  $d \geq d_R$ ,  $w_{\infty,d} = \sup_{\theta \in \Theta} \sup_{x \in V_d} w_\theta(x)/\lambda_\theta(w_\theta) < \infty$  and (ii)  $\sup_{\theta \in \Theta} \text{Var}_\varphi(w_\theta \circ$   
 197  $T_\theta)/\{\lambda_\theta(w_\theta)\}^2 < \infty$ .*

198 A3 is a continuity condition on the NF push-forward density w.r.t. its parameters  $\theta$ . A4 implies that  
 199 the Markov kernel  $K_{N,\theta} = P_{N,\theta} R$  satisfies a drift and minorization condition uniform in  $\theta$ .

---

<sup>2</sup>more specifically, it depends on the the quantities appearing in A2, namely, the maximum of the importance weight on a level set of the drift function for the local kernel R and the variance of the importance weights under the proposal

200 **Theorem 3** (simplified). Assume **A 1-A 3-A 4** and that  $\sum_{k=0}^{\infty} \gamma_k = \infty$ ,  $\sum_{k=0}^{\infty} \gamma_k^2 < \infty$  and  
 201  $\lim_{k \rightarrow \infty} \alpha_k = \alpha_{\infty}$ . Then, w.p. 1, the sequence  $\{\theta_k, k \in \mathbb{N}\}$  converges to the set  $\{\theta \in \Theta, 0 =$   
 202  $\alpha_{\infty} \nabla \text{KL}(\pi || \lambda_{\theta}) + (1 - \alpha_{\infty}) \nabla \text{KL}(\lambda_{\theta} || \pi)\}$ .

203 Theorem 3 proves the convergence of the learning of parameters  $\theta$  to a stationary point of the  
 204 loss. The proof is postponed to Appendix D. Note that once the proposal learning has converged,  
 205 FlEx<sup>2</sup>MCMC boils back to Ex<sup>2</sup>MCMC with a fixed learned proposal. Our experiments show that  
 206 adaptivity can significantly speed up mixing for i-SIR, especially for distributions with complex  
 207 geometries and that the addition of a rejuvenation kernel further improves samples quality.

## 208 4 Related Work

209 The possibility to parametrize very flexible probabilistic models with neural networks thanks to  
 210 deep learning has rekindled interest in adapting MCMC kernels; see e.g. [69, 35, 2, 52, 33]. While  
 211 significant performance gain were found in problems of moderate dimensions, these learning-based  
 212 methods were found to suffer from increasing dimensions as fitting models accurately becomes more  
 213 difficult [19, 45]. Similarly to FlEx<sup>2</sup>MCMC, a few work proposed adaptive algorithms that alternates  
 214 between global and local MCMC moves to ensure ergodicity without requiring a perfect learning  
 215 of the proposal. [57] focused on multimodal distributions and analysed a mode jumping algorithm  
 216 using proposals parametrized as mixture distributions. [25] introduced a combination of a local and a  
 217 global sampler similar to this work, using a more classical independent Metropolis-Hastings sampler  
 218 instead of i-SIR. However the author of [25] did not provide a clear analysis of the benefit of the  
 219 local kernel, conversely to the present work.

220 Another line of work exploits both normalizing flows and common local MCMC kernels for sampling  
 221 [56, 35, 53, 73], yet in a different way. As in this work, the flow  $T$  is trained to transport a simple  
 222 distribution  $\varphi$  near  $\pi$ , which is equivalent to bringing  $T^{-1} \# \pi$  (the pushforward of the original target  
 223 distribution  $\pi$  by the inverse flow  $T^{-1}$ ) close to  $\varphi$ . If  $\varphi$  is chosen to be efficiently sampled by local  
 224 samplers, the hope is that local samplers can also obtain high-quality samples of  $T^{-1} \# \pi$  – samples  
 225 which can be transported back through  $T$  to obtain samples of  $\pi$ . This method, sometimes referred to  
 226 as “neural transport”, effectively reparametrizes the space to disentangle problematic geometries for  
 227 local kernels. Yet, it is unclear what will happen in the tails of the distribution for which the flow is  
 228 likely poorly learned. [56] derived an ergodicity theory for these transported samplers describing  
 229 substantial constraints on maps (see section 2.2.2.).

## 230 5 Numerical experiments

### 231 5.1 Synthetic examples

232 **Multimodal distributions.** Let us start with a toy example highlighting differences between purely  
 233 global i-SIR, purely local MALA and Ex<sup>2</sup>MCMC combining both. We consider sampling from  
 234 a mixture of 3 equally weighted Gaussians in dimension  $d = 2$ . In Figure 2, we compare single  
 235 chains produced by each algorithms. The global proposal is a wide Gaussian, with pools of  $N = 3$   
 236 candidate. The MALA stepsize is chosen to reach a target acceptance rate of  $\sim 0.67$ . This simple  
 237 experiment illustrates the drawbacks of both approaches: i-SIR samples reach all the modes of the  
 238 target, but the chains often get stuck for several steps hindering variability. MALA allows for better  
 239 local exploration of each particular mode, yet it fails to cover all the target support. Meanwhile,  
 240 Ex<sup>2</sup>MCMC retains the benefits of both methods, combining the i-SIR-based global exploration with  
 241 MALA-based local exploration.

242 To illustrate further the performance of the proposed method, we keep the same target mixture model  
 243 yet assigning the uneven weights  $(2/3, 1/6, 1/6)$  to the 3 modes. We start  $M$  chains drawing from  
 244 the initial distribution  $\xi \sim \mathcal{N}(0, 4 I_d)$  and use the same hyper-parameters as above. In Figure 3a  
 245 we provide a simple illustration to the statement (2) and Theorem 2, namely we compare the target  
 246 density to the instantaneous distributions for each sampler propagating  $\xi$  during burn-in steps. As  
 247 MALA does not mix easily between modes, the different statistical weights of the different modes  
 248 can hardly be rendered in few iterations and KL and TV distances stalls after a few iterations. i-SIR  
 249 can visit the different modes, yet it does not necessarily move at each step which slows down its  
 250 covering of the modes full support, which again shows in the speed of decrease of the TV and KL.  
 251 Overcoming both of these shortcomings, Ex<sup>2</sup>MCMC instantaneous density comes much closer to  
 252 the target. Finally, Figure 3b evaluates the same metrics yet for the density estimate obtained with  
 253 single chain samples after burn-in. Results demonstrate once again the superiority of Ex<sup>2</sup>MCMC.  
 254 Further details on these experiments can be found in Appendix E.2.

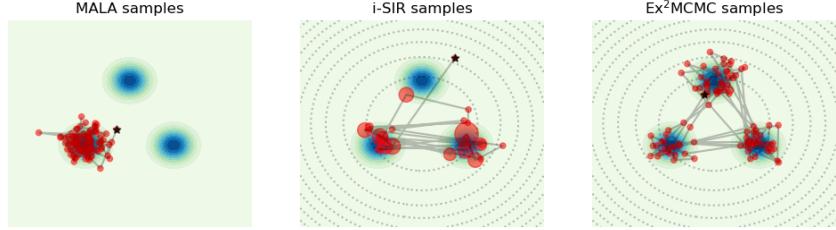


Figure 2: Single chain mixing visualization. – Blue color levels represent the target 2d density. Random chain initialization is noted in black, 100 steps are plotted per sampler: the size of each red dot corresponds to the number of consecutive steps the walkers remains at a given location. For MALA, we generate 300 samples and choose each 3-rd one for comparability. Note that the variance of the global proposal (dotted contour lines) should be relatively large to cover well all the modes. The step size of MALA also can not be increased much to keep reasonable acceptance ratio.

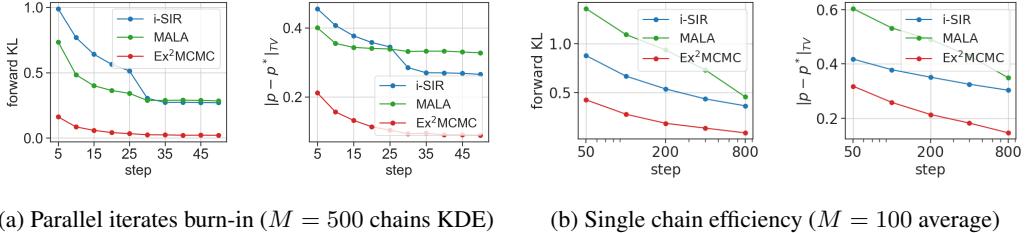


Figure 3: Inhomogeneous 2d Gaussian mixture. – Quantitative analysis during burn-in of parallel chains (a) and for after burn-in for single chains statistics (b).

255 **Distributions with complex geometry.** Next, we turn to highly anisotropic distributions in high  
 256 dimensions. Following [51] and [32], we consider the *funnel* and the *banana-shape* distributions.  
 257 We remind densities in Appendix E.4 along with providing experiments details. For  $d \in [10; 200]$ ,  
 258 we run i-SIR, MALA, Ex<sup>2</sup>MCMC, FIE<sup>2</sup>MCMC, adaptive i-SIR (using the same proposal as  
 259 FIE<sup>2</sup>MCMC, but without interleaved local steps) and the versatile sampler NUTS [34] as a baseline.  
 260 Here the parameter adaptation for FIE<sup>2</sup>MCMC is performed in a pre-run and parameters are frozen  
 261 before sampling. For the adaptive samplers, a simple RealNVP-based normalizing flow [20] is used  
 262 such that total running times, including training, are comparable with NUTS. For Ex<sup>2</sup>MCMC and  
 263 i-SIR the global proposal is a wide Gaussian with a pool of  $N = 2000$  candidates drawn at each  
 264 iteration. For MALA we tune the step size in order to keep acceptance rate approximately at  
 265 0.5. We report the average sliced TV distance and ESS in Figure 4 (see Appendix E.1 for metrics  
 266 definition). In most cases, FIE<sup>2</sup>MCMC is the most reliable algorithm. The only exception is at  
 267 very high dimension for the banana where NUTS performs the best: in this case, tuning the flow to  
 268 learn tails in high-dimension faithfully was costly such that we proceeded to an early stopping to  
 269 maintain comparability with the baseline. Remarkably, FIE<sup>2</sup>MCMC compensates significantly for  
 270 the imperfect flow training, improving over adaptive-i-SIR, but NUTS eventually performs better.  
 271 Conversely, for the funnel, most of the improvement comes from well-trained proposal flow, leading  
 272 to similar behaviors of adaptive i-SIR and FIE<sup>2</sup>MCMC, while both algorithms clearly outperforms  
 273 NUTS in terms of metrics.

## 274 5.2 Sampling from GANs as Energy-based models (EBMs)

275 Generative adversarial networks (GANs [27]) are a class of generative models defined by a pair of a  
 276 generator network  $G$  and a discriminator network  $D$ . The generator  $G$  takes a latent variable  $z$  from a  
 277 prior density  $p_0(z)$ ,  $z \in \mathbb{R}^d$ , and generates an observation  $G(z) \in \mathbb{R}^D$  in the observation space. The  
 278 discriminator takes a sample in the observation space and aims to discriminate between true examples  
 279 and false examples produced by the generator. Recently, it has been advocated to consider GANs as  
 280 Energy-Based Models (EBMs) [72, 17]. Following [17], we consider the EBM model induced by the  
 281 GAN in latent space. Recall that an EBM is defined by a Boltzmann distribution  $p(z) = e^{-E(z)} / Z$ ,  
 282  $z \in \mathbb{R}^d$ , where  $E(z)$  is the energy function and  $Z$  is the normalizing constant. Note that Wasserstein  
 283 GANs also allow for an energy-based interpretation (see [17]), although the interpretation of the  
 284 discriminator in this case is different. The energy function is given by

$$E_{JS}(z) = -\log p_0(z) - \text{logit}(D(G(z))), \quad E_W(z) = -\log p_0(z) - D(G(z)), \quad z \in \mathbb{R}^d, \quad (6)$$

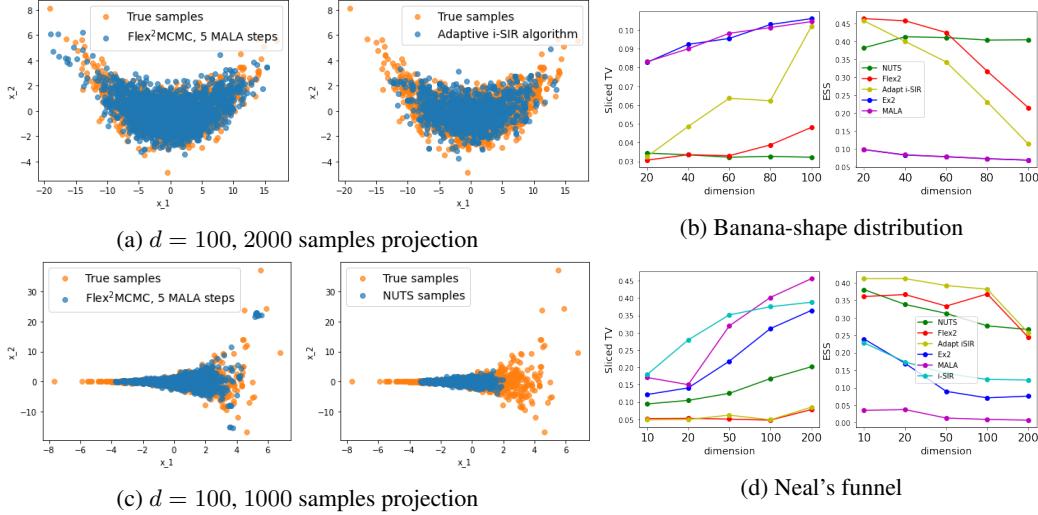


Figure 4: Anisotropic Funnel and Banana-shape distributions – (a) and (b) visualize samples projected onto the first 2 coordinates of tested algorithms (blue) versus true samples obtained by reparametrization (orange). (c) and (d) compare Sliced Total Variation and Effective Sample Size as a function of dimension. i-SIR is removed from (b) as corresponding metrics for  $d > 20$  are significantly worse.

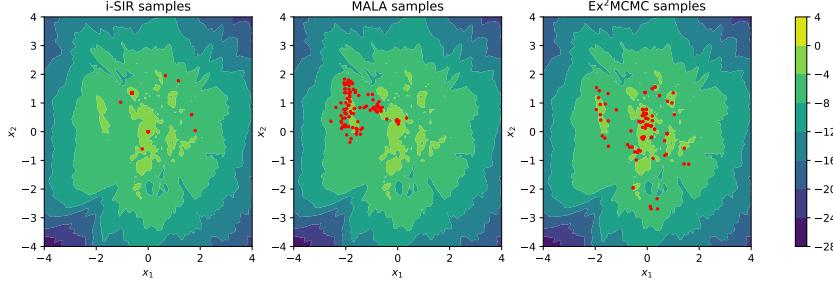


Figure 5: MNIST energy landscape and single chain latent samples visualizations.

for the vanilla Jensen-Shannon and Wasserstein GANs, respectively. Here  $\text{logit}(y)$ ,  $y \in (0, 1)$  is the inverse of the sigmoid function and  $p_0(z) = \mathcal{N}(0, I_d)$ .

**MNIST results.** We consider a simple Jensen-Shannon GAN model trained on the MNIST dataset with latent space dimension  $d = 2$ . We compare samples obtained by i-SIR, MALA, and Ex<sup>2</sup>MCMC from the energy-based model associated with  $E_{JS}(z)$ , see (6). We use a wide normal distribution as the global proposal for i-SIR and Ex<sup>2</sup>MCMC, and pools of candidates at each iteration  $N = 10$ . The step-size of MALA is tuned to keep an acceptance rate  $\sim 0.5$ . We visualize chains of 100 steps in the latent space obtained with each method in Figure 5. Note that the poor agreement between the proposal and the landscape makes it difficult for i-SIR to accept from the proposal and for MALA to explore many modes of the latent distribution, as shown in Figure 5. Ex<sup>2</sup>MCMC combines effectively global and local moves, encouraging better diversity associated with a better mixing time. The images corresponding to the sampled latent space locations are displayed in Figure 6 and reflect the diversity issue of MALA and i-SIR. Further details and experiments are provided in Appendix E.5.1, including similar results for WGAN-GP [31] and the associated EBM  $E_W(z)$ .

**Cifar-10 results.** We consider two popular architectures trained on Cifar-10, DC-GAN [58] and SN-GAN [48]. In both cases the dimension of the latent space equals  $d = 128$ . Together with the non-trivial geometry of the corresponding energy landscapes, the large dimension makes sampling with NUTS unfeasible in terms of computational time. We perform sampling from mentioned GANs as energy-based models using i-SIR, MALA, Ex<sup>2</sup>MCMC, and FLEX<sup>2</sup>MCMC. In i-SIR and Ex<sup>2</sup>MCMC we use the prior  $p_0(z)$  as a global proposal with a pool of  $N = 10$  candidates. For FLEX<sup>2</sup>MCMC we perform training and sampling simultaneously. Additional implementation details are provided in Appendix E.5.2. To evaluate sampling quality, we report the values of the energy function  $E(z)$ , averaged over 500 independent runs of each sampler. We present the results in Figure 7 together with the images produced by each sampler. Note that we observe that Ex<sup>2</sup>MCMC

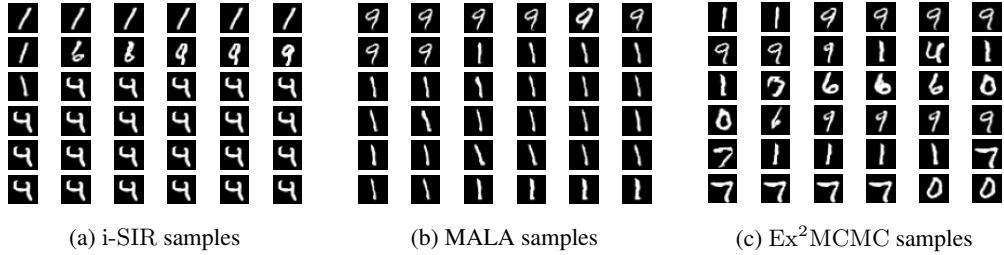


Figure 6: MNIST samples visualization. – Single chains run, sequential steps.

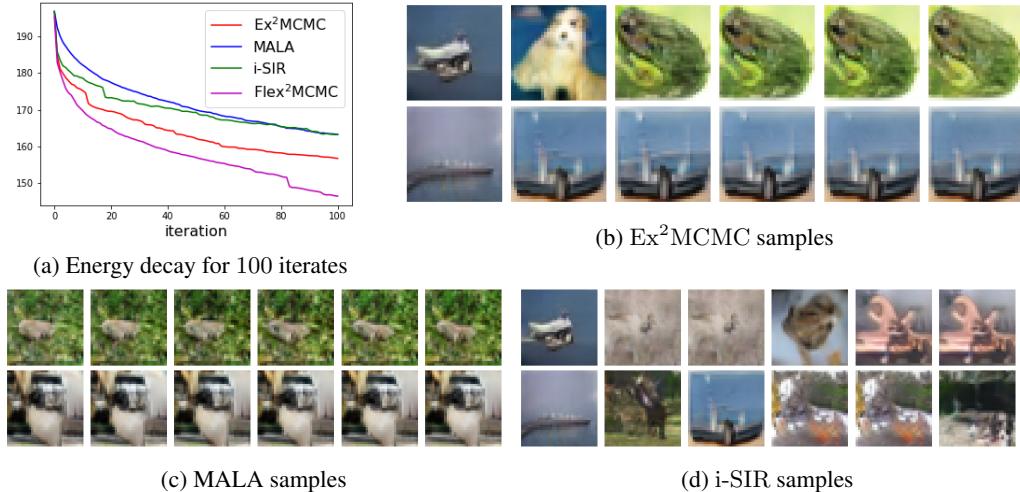


Figure 7: Cifar-10 energy and sampling results with DC-GAN architecture. Along the horizontal lines we visualize each 10th sample from a single trajectory.

and FlEx<sup>2</sup>MCMC reach low level of energies faster than other methods. Visualizations indicate that MALA is unlikely to escape the mode of the distribution  $p(z)$  it started from, while i-SIR and Ex<sup>2</sup>MCMC/FlEx<sup>2</sup>MCMC better explores the target support. However, global move appear to become more rare after some number of iterations for Ex<sup>2</sup>MCMC/FlEx<sup>2</sup>MCMC, which then exploit a particular mode with MALA steps. We here hit the following limitation: i-SIR remains at relatively high-energies, failing to explore well modes basins but still accepting global moves, while Ex<sup>2</sup>MCMC/FlEx<sup>2</sup>MCMC explores well modes basins but eventually remains trapped. We predict that improving further the quality of the FlEx<sup>2</sup>MCMC proposal by scaling the normalizing flow architecture and increasing learning time would allow for more global moves. We provide additional experiments (including ones with SN-GAN) and visualizations in Appendix E.5.2.

## 319 6 Conclusions and further research directions

320 The present paper examines the benefits of combining local and global samplers. From a theoretical  
 321 point of view, we show that global samplers are more robust when coupled with local samplers.  
 322 Namely, a  $V$ -geometric ergodicity is obtained for the Ex<sup>2</sup>MCMC kernel under minimal assumptions.  
 323 Meanwhile, the global samplers drives exploration when properly adjusted. Therefore, we also  
 324 describe the adaptive version FlEx<sup>2</sup>MCMC of the strategy involving the learning of a global proposal  
 325 parametrized by a normalizing flow. We also check for the learning convergence along the adaptive  
 326 MCMC run. Finally, a series of numerical experiments confirms the superiority of the strategy,  
 327 including the high-dimensional examples. While the startegy was described and analyzed for the  
 328 i-SIR global kernel, we note that it would be possible to extend the theory to other independent  
 329 global samplers. We expect that the benefit of the combination would remain. Further studies of  
 330 FlEx<sup>2</sup>MCMC, in particular the derivation of its mixing rate, is an interesting direction for future  
 331 work.

332 **References**

- 333 [1] S. Agapiou, O. Papaspiliopoulos, D. Sanz-Alonso, and A. M. Stuart. Importance sampling:  
334 Intrinsic dimension and computational cost. *Statistical Science*, 32(3):405–431, 2017.
- 335 [2] M. Albergo, G. Kanwar, and P. Shanahan. Flow-based generative models for Markov chain  
336 Monte Carlo in lattice field theory. *Physical Review D*, 100(3):034515, 2019.
- 337 [3] C. Andrieu. On random-and systematic-scan samplers. *Biometrika*, 103(3):719–726, 2016.
- 338 [4] C. Andrieu, A. Doucet, and R. Holenstein. Particle Markov chain Monte Carlo methods. *Journal*  
339 *of the Royal Statistical Society: Series B*, 72(3):269–342, 2010.
- 340 [5] C. Andrieu, A. Lee, M. Vihola, et al. Uniform ergodicity of the iterated conditional SMC and  
341 geometric ergodicity of particle Gibbs samplers. *Bernoulli*, 24(2):842–872, 2018.
- 342 [6] C. Andrieu and É. Moulines. On the ergodicity properties of some adaptive mcmc algorithms.  
343 *The Annals of Applied Probability*, 16(3):1462–1505, 2006.
- 344 [7] C. Andrieu, É. Moulines, and P. Priouret. Stability of stochastic approximation under verifiable  
345 conditions. *SIAM Journal on control and optimization*, 44(1):283–312, 2005.
- 346 [8] C. Andrieu, V. B. Tadić, and M. Vihola. On the stability of some controlled markov chains and  
347 its applications to stochastic approximation with markovian dynamic. *The Annals of Applied*  
348 *Probability*, 25(1):1–45, 2015.
- 349 [9] C. Andrieu and J. Thoms. A tutorial on adaptive MCMC. *Statistics and computing*, 18(4):343–  
350 373, 2008.
- 351 [10] C. Andrieu and M. Vihola. Markovian stochastic approximation with expanding projections.  
352 *Bernoulli*, 20(2):545–585, 2014.
- 353 [11] T. Bengtsson, P. J. Bickel, and B. Li. Curse-of-dimensionality revisited: Collapse of the particle  
354 filter in very large scale systems. *arXiv: Statistics Theory*, pages 316–334, 2008.
- 355 [12] A. Benveniste, M. Métivier, and P. Priouret. *Adaptive algorithms and stochastic approximations*,  
356 volume 22 of *Applications of Mathematics (New York)*. Springer-Verlag, Berlin, 1990. Translated  
357 from the French by Stephen S. Wilson.
- 358 [13] J. Besag. Comments on “Representations of knowledge in complex systems” by U. Grenander  
359 and M. Miller. *J. Roy. Statist. Soc. Ser. B*, 56:591–592, 1994.
- 360 [14] N. Bonneel, M. Van De Panne, S. Paris, and W. Heidrich. Displacement interpolation using  
361 lagrangian mass transport. In *Proceedings of the 2011 SIGGRAPH Asia Conference*, pages  
362 1–12, 2011.
- 363 [15] V. S. Borkar. *Stochastic approximation: a dynamical systems viewpoint*, volume 48. Springer,  
364 2009.
- 365 [16] M. F. Bugallo, V. Elvira, L. Martino, D. Luengo, J. Miguez, and P. M. Djuric. Adaptive  
366 importance sampling: The past, the present, and the future. *IEEE Signal Processing Magazine*,  
367 34(4):60–79, 2017.
- 368 [17] T. Che, R. Zhang, J. Sohl-Dickstein, H. Larochelle, L. Paull, Y. Cao, and Y. Bengio. Your GAN  
369 is Secretly an Energy-based Model and You Should Use Discriminator Driven Latent Sampling.  
370 In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in*  
371 *Neural Information Processing Systems*, volume 33, pages 12275–12287. Curran Associates,  
372 Inc., 2020.
- 373 [18] M.-F. Chen and F.-Y. Wang. Estimation of spectral gap for elliptic operators. *Trans. Amer. Math.*  
374 *Soc.*, 349(3):1239–1267, 1997.
- 375 [19] L. Del Debbio, J. Marsh Rossney, and M. Wilson. Efficient modeling of trivializing maps for  
376 lattice  $\phi^4$  theory using normalizing flows: A first look at scalability. *Physical Review D*, 104(9),  
377 2021.

- 378 [20] L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using real NVP. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- 381 [21] R. Douc, E. Moulines, P. Priouret, and P. Soulier. *Markov chains*. Springer Series in Operations  
382 Research and Financial Engineering. Springer, Cham, 2018.
- 383 [22] A. Durmus and E. Moulines. On the geometric convergence for MALA under verifiable  
384 conditions. 2022.
- 385 [23] A. Eberle. Reflection couplings and contraction rates for diffusions. *Probab. Theory Related  
386 Fields*, pages 1–36, 2015.
- 387 [24] D. L. Ermak. A computer simulation of charged particles in solution. i. technique and equilib-  
388 rium properties. *The Journal of Chemical Physics*, 62(10):4189–4196, 1975.
- 389 [25] M. Gabrié, G. M. Rotskoff, and E. Vanden-Eijnden. Adaptive Monte Carlo augmented with  
390 normalizing flows. *Proceedings of the National Academy of Sciences*, 119(10), mar 2022.
- 391 [26] M. Girolami and B. Calderhead. Riemann manifold langevin and hamiltonian monte carlo  
392 methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(2):123–  
393 214, 2011.
- 394 [27] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville,  
395 and Y. Bengio. Generative adversarial nets. In *Proceedings of the 27th International Conference  
396 on Neural Information Processing Systems - Volume 2*, NIPS’14, page 2672–2680, Cambridge,  
397 MA, USA, 2014. MIT Press.
- 398 [28] U. Grenander. Tutorial in pattern theory. Division of Applied Mathematics, Brown University,  
399 Providence, 1983.
- 400 [29] U. Grenander and M. I. Miller. Representations of knowledge in complex systems. *J. Roy.  
401 Statist. Soc. Ser. B*, 56(4):549–603, 1994. With discussion and a reply by the authors.
- 402 [30] M. G. Gu and F. H. Kong. A stochastic approximation algorithm with markov chain monte-carlo  
403 method for incomplete data estimation problems. *Proceedings of the National Academy of  
404 Sciences*, 95(13):7270–7274, 1998.
- 405 [31] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved train-  
406 ing of wasserstein gans. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,  
407 S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,  
408 volume 30. Curran Associates, Inc., 2017.
- 409 [32] H. Haario, E. Saksman, and J. Tamminen. Adaptive proposal distribution for random walk  
410 metropolis algorithm. *Computational Statistics*, 14(3):375–395, 1999.
- 411 [33] D. C. Hackett, C.-C. Hsieh, M. S. Albergo, D. Boyd, J.-W. Chen, K.-F. Chen, K. Cranmer,  
412 G. Kanwar, and P. E. Shanahan. Flow-based sampling for multimodal distributions in lattice  
413 field theory. *arXiv preprint*, 2107.00734, 2021.
- 414 [34] M. D. Hoffman, A. Gelman, et al. The no-U-turn sampler: adaptively setting path lengths in  
415 Hamiltonian Monte Carlo. *J. Mach. Learn. Res.*, 15(1):1593–1623, 2014.
- 416 [35] M. D. Hoffman, P. Sountsov, J. V. Dillon, I. Langmore, D. Tran, and S. Vasudevan. NeuTra-  
417 lizing Bad Geometry in Hamiltonian Monte Carlo Using Neural Transport. In *1st Symposium  
418 on Advances in Approximate Bayesian Inference, 2018 1–5*, 2019.
- 419 [36] C.-W. Huang, D. Krueger, A. Lacoste, and A. Courville. Neural autoregressive flows. In J. Dy  
420 and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*,  
421 volume 80 of *Proceedings of Machine Learning Research*, pages 2078–2087. PMLR, 10–15 Jul  
422 2018.
- 423 [37] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR 2015*, 2015.

- 424 [38] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling. Improving  
 425 variational inference with inverse autoregressive flow, 2016.
- 426 [39] I. Kobyzev, S. Prince, and M. Brubaker. Normalizing flows: An introduction and review of  
 427 current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- 428 [40] H. J. Kushner and G. G. Yin. *Stochastic approximation and recursive algorithms and applica-*  
 429 *tions*, volume 35 of *Applications of Mathematics (New York)*. Springer-Verlag, New York,  
 430 second edition, 2003. Stochastic Modelling and Applied Probability.
- 431 [41] A. Lee. *On auxiliary variables and many-core architectures in computational statistics*. PhD  
 432 thesis, University of Oxford, 2011.
- 433 [42] A. Lee, C. Yau, M. B. Giles, A. Doucet, and C. C. Holmes. On the utility of graphics  
 434 cards to perform massively parallel simulation of advanced Monte Carlo methods. *Journal of*  
 435 *computational and graphical statistics*, 19(4):769–789, 2010.
- 436 [43] J. S. Liu, F. Liang, and W. H. Wong. The multiple-try method and local optimization in  
 437 Metropolis sampling. *Journal of the American Statistical Association*, 95(449):121–134, 2000.
- 438 [44] J. S. Liu, W. H. Wong, and A. Kong. Covariance structure of the gibbs sampler with applications  
 439 to the comparisons of estimators and augmentation schemes. *Biometrika*, 81(1):27–40, 1994.
- 440 [45] A. H. Mahmoud, M. Masters, S. J. Lee, and M. A. Lill. Accurate Sampling of Macromolecular  
 441 Conformations Using Adaptive Deep Learning and Coarse-Grained Representation. *Journal of*  
 442 *Chemical Information and Modeling*, 62(7):1602–1617, apr 2022.
- 443 [46] J. Mattingly, A. Stuart, and D. Higham. Ergodicity for {SDEs} and approximations: locally  
 444 lipschitz vector fields and degenerate noise. *Stochastic Processes and their Applications*,  
 445 101(2):185 – 232, 2002.
- 446 [47] M. Métivier and P. Priouret. Théorèmes de convergence presque sûre pour une classe  
 447 d’algorithmes stochastiques à pas décroissant. *Probability Theory and related fields*, 74(3):403–  
 448 428, 1987.
- 449 [48] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative  
 450 adversarial networks. *arXiv:1802.05957*, 2018.
- 451 [49] C. A. Naesseth, F. Lindsten, and D. Blei. Markovian score climbing: Variational inference with  
 452  $\text{KL}(\text{pliq})$ . *Advances in Neural Information Processing Systems*, 2020-Decem(MCMC), 2020.
- 453 [50] R. M. Neal. Bayesian learning via stochastic dynamics. In *Advances in Neural Information*  
 454 *Processing Systems 5, [NIPS Conference]*, pages 475–482, San Francisco, CA, USA, 1993.  
 455 Morgan Kaufmann Publishers Inc.
- 456 [51] R. M. Neal. Slice sampling. *The Annals of Statistics*, 31(3):705 – 767, 2003.
- 457 [52] K. A. Nicoli, S. Nakajima, N. Strodthoff, W. Samek, K. R. Müller, and P. Kessel. Asymptotically  
 458 unbiased estimation of physical observables with neural samplers. *Physical Review E*, 101(2),  
 459 2020.
- 460 [53] F. Noé, S. Olsson, J. Köhler, and H. Wu. Boltzmann generators: Sampling equilibrium states of  
 461 many-body systems with deep learning. *Science*, 365(6457), 2019.
- 462 [54] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan. Nor-  
 463 malizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*,  
 464 22(57):1–64, 2021.
- 465 [55] G. Parisi. Correlation functions and computer simulations. *Nuclear Physics B*, 180:378–384,  
 466 1981.
- 467 [56] M. D. Parno and Y. M. Marzouk. Transport map accelerated markov chain monte carlo.  
 468 *SIAM-ASA Journal on Uncertainty Quantification*, 6(2):645–682, 2018.

- 469 [57] E. Pompe, C. Holmes, and K. Łatuszyński. A framework for adaptive mcmc targeting multi-  
 470 modal distributions. *Annals of Statistics*, 48(5):2930–2952, 2020.
- 471 [58] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional  
 472 generative adversarial networks. *arXiv:1511.06434*, 2016.
- 473 [59] C. Robert and G. Casella. *Monte Carlo statistical methods*. Springer Science & Business Media,  
 474 2013.
- 475 [60] C. P. Robert. *The Metropolis–Hastings Algorithm*, pages 1–15. John Wiley & Sons, Ltd, 2015.
- 476 [61] G. O. Roberts and J. S. Rosenthal. General state space markov chains and mcmc algorithms.  
 477 *Probability surveys*, 1:20–71, 2004.
- 478 [62] G. O. Roberts and J. S. Rosenthal. Examples of adaptive mcmc. *Journal of computational and*  
 479 *graphical statistics*, 18(2):349–367, 2009.
- 480 [63] G. O. Roberts and R. L. Tweedie. Exponential convergence of Langevin distributions and their  
 481 discrete approximations. *Bernoulli*, 2(4):341–363, 1996.
- 482 [64] G. O. Roberts and R. L. Tweedie. Geometric convergence and central limit theorems for  
 483 multidimensional Hastings and Metropolis algorithms. *Biometrika*, 83(1):95–110, 03 1996.
- 484 [65] P. J. Rossky, J. D. Doll, and H. L. Friedman. Brownian dynamics as smart Monte Carlo  
 485 simulation. *The Journal of Chemical Physics*, 69(10):4628–4633, 1978.
- 486 [66] D. B. Rubin. Comment: A noniterative Sampling/Importance Resampling alternative to the data  
 487 augmentation algorithm for creating a few imputations when fractions of missing information are  
 488 modest: The SIR algorithm. *Journal of the American Statistical Association*, 82(398):542–543,  
 489 1987.
- 490 [67] D. Sanz-Alonso. Importance sampling and necessary sample size: an information theory  
 491 approach. *SIAM/ASA Journal on Uncertainty Quantification*, 6(2):867–879, 2018.
- 492 [68] C. Snyder, T. Bengtsson, P. Bickel, and J. Anderson. Obstacles to high-dimensional particle  
 493 filtering. *Monthly Weather Review*, 136(12):4629 – 4640, 2008.
- 494 [69] J. Song, S. Zhao, and S. Ermon. A-NICE-MC: Adversarial training for MCMC. In *Advances in*  
 495 *Neural Information Processing Systems*, pages 5140–5150, 2017.
- 496 [70] H. Tjelmeland. Using all Metropolis–Hastings proposals to estimate mean values. Technical  
 497 report, 2004.
- 498 [71] S. T. Tokdar and R. E. Kass. Importance sampling: a review. *WIREs Computational Statistics*,  
 499 2(1):54–60, 2010.
- 500 [72] R. Turner, J. Hung, E. Frank, Y. Saatchi, and J. Yosinski. Metropolis-Hastings generative  
 501 adversarial networks. In *International Conference on Machine Learning*, pages 6345–6353.  
 502 PMLR, 2019.
- 503 [73] L. Zhang, C. A. Naesseth, and D. M. Blei. Transport Score Climbing: Variational Inference  
 504 Using Forward KL and Adaptive Neural Transport. *arXiv preprint*, 2202.01841, 2022.

505 **Checklist**

- 506 1. For all authors...
- 507 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
508 contributions and scope? **[Yes]**
- 509 (b) Did you describe the limitations of your work? **[Yes]**
- 510 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 511 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
512 them? **[Yes]** The paper suggests novel MCMC technique and is validated on artificial  
513 and standard datasets.
- 514 2. If you are including theoretical results...
- 515 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** See Section 2  
516 and Section 3 in the main text.
- 517 (b) Did you include complete proofs of all theoretical results? **[Yes]** Yes, the proofs of  
518 Section 2 and Section 3 are provided in Appendix B and Appendix D.
- 519 3. If you ran experiments...
- 520 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
521 perimental results (either in the supplemental material or as a URL)? **[Yes]** Code to  
522 reproduce experiments is attached to the supplement. Due to size constraints, we are  
523 not available to present all the pre-trained GANs models for the section Section 5, but  
524 we intend to do so when possible.
- 525 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
526 were chosen)? **[Yes]** The hyperparameters are provided in the supplement paper.
- 527 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
528 ments multiple times)? **[Yes]** Partially yes, but not for all experiments.
- 529 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
530 of GPUs, internal cluster, or cloud provider)? **[Yes]** We provide this information in the  
531 supplement paper.
- 532 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 533 (a) If your work uses existing assets, did you cite the creators? **[N/A]** We use only the  
534 common knowledge datasets.
- 535 (b) Did you mention the license of the assets? **[N/A]**
- 536 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
- 537 (d) Did you discuss whether and how consent was obtained from people whose data you're  
538 using/curating? **[N/A]**
- 539 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
540 information or offensive content? **[N/A]**
- 541 5. If you used crowdsourcing or conducted research with human subjects...
- 542 (a) Did you include the full text of instructions given to participants and screenshots, if  
543 applicable? **[N/A]**
- 544 (b) Did you describe any potential participant risks, with links to Institutional Review  
545 Board (IRB) approvals, if applicable? **[N/A]**
- 546 (c) Did you include the estimated hourly wage paid to participants and the total amount  
547 spent on participant compensation? **[N/A]**

548 **A i-SIR Algorithm**

549 **A.1 i-SIR and Multiple-try Metropolis (MTM) algorithm**

550 In the MTM algorithm,  $N$  i.i.d.sample proposals  $\{X_{k+1}^i\}_{i=1}^N$  are drawn from a kernel  $\mathsf{T}(y, \cdot)$  in each  
 551 iteration. In a second step, a sample  $Y_{k+1}^*$  is selected with probability proportional to the weights (the  
 552 exact expression of the weighting weights differs from ours, but this does not change the complexity  
 553 of the algorithm). In a third step,  $N - 1$  i.i.d. proposals are drawn from the kernel  $\mathsf{T}(Y_{k+1}^*, \cdot)$  and it  
 554 is assumed that the move is  $Y_{k+1} = Y_{k+1}^*$  with an *generalised M-H* ratio, see [43, eq. 3]. This step is  
 555 bypassed in i-SIR, reducing the computational complexity by a factor of 2.

556 **A.2 i-SIR as a systematic scan two-stage Gibbs sampler**

557 We analyze a slightly modified version of the i-SIR algorithm, with an extra randomization of the  
 558 state position. The  $k$ -th iteration is defined as follows. Given a state  $Y_k \in \mathbb{X}$ ,

- 559 (i) draw  $I_{k+1} \in \{1, \dots, N\}$  uniformly at random and set  $X_{k+1}^{I_{k+1}} = Y_k$ ;
- 560 (ii) draw  $X_{k+1}^{1:N \setminus \{I_{k+1}\}}$  independently from the proposal distribution  $\lambda$ ;
- 561 (iii) compute, for  $i \in \{1, \dots, N\}$ , the normalized importance weights

$$\omega_{N,k+1}^i = w(X_{k+1}^i) / \sum_{\ell=1}^N w(X_{k+1}^\ell);$$

- 562 (iv) select  $Y_{k+1}$  from the set  $X_{k+1}^{1:N}$  by choosing  $X_{k+1}^i$  with probability  $\omega_{N,k+1}^i$ .

563 Thus, compared to the simplified i-SIR algorithm given in the introduction, the state is inserted  
 564 uniformly at random into the list of candidates instead of being inserted at the first position. Of course,  
 565 this change has no impact as long as we are interested in integrating functions that are permutation  
 566 invariant with respect to candidates, which is the case throughout our work. Still, this randomization  
 567 makes the analysis much more transparent.

568 In what follows, we show that i-SIR can be interpreted as a systematic-scan two-stage Gibbs sampler  
 569 sampling, which alternately samples from the full conditionals of the extended target  $\varphi_N$ , which  
 570 is carefully defined below in terms of the state and candidate pool. Here we essentially follow the  
 571 work of [70, 4, 5]. This is formalized by a dual representation of  $\varphi_N$ , presented below in Theorem 4,  
 572 which provides the two complete conditionals in question. We introduce the Markov kernel

$$\Lambda_N(y, dx^{1:N}) = \frac{1}{N} \sum_{i=1}^N \delta_y(dx^i) \prod_{j \neq i} \lambda(dx^j)$$

on  $\mathbb{X} \times \mathcal{X}^{\otimes N}$ , which probabilistically describes the candidate selection operation in i-SIR. Note that  
 by construction, for each  $y \in \mathbb{X}$ ,  $\ell \in \{1, \dots, N\}$  and nonnegative measurable function  $h : \mathbb{X} \rightarrow \mathbb{R}^+$ ,

$$\Lambda_N h(y) = \int \Lambda_N(y, dx^{1:N}) h(x^\ell) = \left(1 - \frac{1}{N}\right) \lambda(h) + \frac{1}{N} h(y).$$

573 Using the kernel  $\Lambda_N$  we may now define properly the extended target  $\varphi_N$  as the probability law

$$\varphi_N(dy, x^{1:N}) = \pi(dy) \Lambda_N(y, dx^{1:N}) = \frac{1}{N} \sum_{i=1}^N \pi(dy) \delta_y(dx^i) \prod_{j \neq i} \lambda(dx^j)$$

574 on  $(\mathbb{X}^{N+1}, \mathcal{X}^{\otimes(N+1)})$ . Note that since for every  $A \in \mathcal{X}$ ,  $\varphi_N(1_{A \times \mathbb{X}}) = \pi(A)$ , the target  $\pi$  coincides  
 575 with the marginal of  $\varphi_N$  with respect to the state. Moreover, it is easily seen that  $\Lambda_N$  provides the  
 576 conditional distribution, under  $\varphi_N$ , of the candidate pool given the state.

577 On the other hand, using that  $\pi(dy) \delta_y(dx^i) = w(x^i) \lambda(dx^i) \delta_{x^i}(dy) / \lambda(w)$ , the marginal distribution  
 578  $\pi_N$  of  $\varphi_N$  with respect to  $x^{1:N}$  is given by

$$\pi_N(dx^{1:N}) = \frac{1}{\lambda(w)} \Gamma_N 1_{\mathbb{X}}(x^{1:N}) \prod_{j=1}^N \lambda(dx^j), \quad (7)$$

579 where we have set

$$\Gamma_N(x^{1:N}, dy) = \sum_{i=1}^N w(x^i) \delta_{x^i}(dy)/N, \quad \Pi_N(x^{1:N}, dy) = \Gamma_N(x^{1:N}, dy)/\Gamma_N 1_{\mathbb{X}}(x^{1:N})$$

580 It is interesting to note that the marginal  $\pi_N$  has a probability density function, proportional to  
581  $\Gamma_N 1_{\mathbb{X}}(x^{1:N}) = \sum_{i=1}^N w(x^i)/N$ , with respect to the product measure  $\lambda^{\otimes N}$ . Using (7), we immedi-  
582 ately obtain the following result.

583 **Theorem 4** (duality of extended target). *For every  $N \in \mathbb{N}^*$ ,*

$$\varphi_N(dy, x^{1:N}) = \pi(dy) \Lambda_N(y, dx^{1:N}) = \pi_N(dx^{1:N}) \Pi_N(x^{1:N}, dy).$$

584 Using this dual representation of  $\varphi_N$ , i-SIR can be interpreted as a two-stage Gibbs sampler. Given  
585 the state  $Y_k$ ,  $N$  candidates  $X_{k+1}^{1:N}$  are sampled from  $\Lambda_N(Y_k, \cdot)$ . In a second step, the next state  $Y_{k+1}$   
586 is sampled given the current candidates from  $\Pi_N(X_{k+1}^{1:N}, \cdot)$ . The two-stages Gibbs sampler generates  
587 a Markov chain  $((Y_k, X_k^{1:N}))_{k \in \mathbb{N}}$  with Markov kernel

$$\mathbf{P}_N((y, x^{1:N}), C) = \int \Lambda_N(y, d\tilde{x}^{1:N}) \Pi_N(\tilde{x}^{1:N}, d\tilde{y}) 1_C(dy, \tilde{x}^{1:N}), \quad C \in \mathcal{X}^{\otimes(N+1)}.$$

588 Note that the Markov kernel  $\mathbf{P}_N(y, x^{1:N}, \cdot)$  does not depend on  $x^{1:N}$ , which means that only the  
589 state  $Y_k$  needs to be stored from one iteration to another. Given a distribution  $\xi$  on  $(\mathbb{X}^{n+1}, \mathcal{X}^{\otimes(n+1)})$ ,  
590 we denote by  $\mathbb{P}_{\xi}$  the distribution of the canonical Markov chain  $((Y_k, X_k^{1:N}))_{k \in \mathbb{N}}$  with kernel  $\mathbf{P}_N$ .  
591 With these notations, for any nonnegative measurable function  $f : \mathbb{X}^{n+1} \rightarrow \mathbb{R}$ , we get, for  $k \in \mathbb{N}^*$ ,

$$\mathbb{E}_{\xi} [f(Y_k, X_k^{1:N}) | \mathcal{F}_{k-1}] = \int \mathbf{P}_N((Y_{k-1}, X_{k-1}^{1:N}), d(y, x^{1:N})) f(x^{1:N}) = \mathbf{P}_N f(Y_{k-1}, X_{k-1}^{1:N}).$$

592 The systematic scan two-stages Gibbs sampler is one of the MCMC algorithm structures that has  
593 given rise to many works. We summarize in the theorem below the important properties of this  
594 sampler; see [44], [59, Chapter 9], [3] and the references therein.

595 **Theorem 5.** *Assume that for any  $y \in \mathbb{X}$ ,  $w(y) > 0$ . Then,*

- 596 • *The Markov kernel  $\mathbf{P}_N$  is Harris recurrent and ergodic with unique invariant distribution  $\varphi_N$ .*
- 597 • *The Markov kernel  $\mathbf{P}_N$  is reversible w.r.t.  $\pi$ , Harris recurrent and ergodic.*

598 The proof follows from [59, Theorem 9.6, Lemma 9.11]. The following theorem establishes the  
599 unbiasedness of the estimator  $\Pi_N f(X^{1:N})$  under  $\varphi_N$ .

600 **Theorem 6.** *For every  $N \in \mathbb{N}^*$  and  $\pi$ -integrable function  $f$ ,*

$$\pi(f) := \int \Pi_N f(x^{1:N}) \pi_N(dx^{1:N}) = \int \Pi_N f(x^{1:N}) \pi(dx^1) \prod_{j=2}^N \lambda(dx^j).$$

600 *Proof.* Using (7) we get

$$\begin{aligned} \int \pi_N(dx^{1:N}) \Pi_N f(x^{1:N}) &= \int \frac{1}{N \lambda(w)} \sum_{\ell=1}^N w(x^\ell) \Pi_N f(x^{1:N}) \prod_{j=1}^N \lambda(dx^j) \\ &= \frac{1}{N \lambda(w)} \int \sum_{i=1}^N w(x^i) f(x^i) \prod_{j=1}^N \lambda(dx^j) = \pi(f), \end{aligned}$$

601 and the first identity follows. The second identity stems from the fact that the function  $\Pi_N f(x^{1:N})$  is  
602 invariant under permutation.  $\square$

603 **B Proofs of Section 2**

604 **B.1 Uniform geometric ergodicity of the i-SIR Markov kernel**

605 Here we provide a simple direct proof of the bound (2). We preface the proof by a technical lemma.

606 **Lemma 7.** *Let  $Y^{1:M}$  be  $M$  independent random variables, satisfying  $\mathbb{E}[Y_i] = 1$ , and  $\text{Var}[Y_i] < \infty$   
607 for  $i \in \{1, \dots, M\}$ . Then, for  $S_M = \sum_{i=1}^M Y_i$  and  $a, b > 0$*

$$\mathbb{E}\left[(a + bS_M)^{-1}\right] \leq (a + bM/2)^{-1} + (4/a) \text{Var}[S_M]/M^2.$$

608 *Proof.* Let  $K \geq 0$ . Then we get

$$\frac{1}{a + bS_M} = \frac{1}{a + bS_M} \mathbf{1}\{S_M < K\} + \frac{1}{a + bS_M} \mathbf{1}\{S_M \geq K\} \leq \frac{1}{a + bK} + \frac{1}{a} \mathbf{1}\{S_M < K\}$$

609 and in particular,  $\mathbb{E}[(a + bS_M)^{-1}] \leq (a + bK)^{-1} + a^{-1}\mathbb{P}(S_M < K)$ . By Markov's inequality,

$$\mathbb{P}(S_M < K) = \mathbb{P}(S_M - M < -(M - K)) \leq \frac{\text{Var}[S_M]}{(M - K)^2}$$

610 In particular, for  $K = M/2$ , we have  $\mathbb{P}(S_M < K) \leq 4 \text{Var}[S_M]/M^2$ .  $\square$

611 *Proof of (2).* For  $(x, \mathsf{A}) \in \mathbb{X} \times \mathcal{X}$ , we get

$$\begin{aligned} \mathsf{P}_N(x, \mathsf{A}) &= \int \delta_x(dx^1) \sum_{i=1}^N \frac{w(x^i)}{\sum_{j=1}^N w(x^j)} \mathbf{1}_{\mathsf{A}}(x^i) \prod_{j=2}^N \lambda(dx^j) \\ &= \int \frac{w(x)}{w(x) + \sum_{j=2}^N w(x^j)} \mathbf{1}_{\mathsf{A}}(x) \prod_{j=2}^N \lambda(dx^j) + \int \sum_{i=2}^N \frac{w(x^i)}{w(x) + \sum_{j=2}^N w(x^j)} \mathbf{1}_{\mathsf{A}}(x^i) \prod_{j=2}^N \lambda(dx^j) \\ &\geq \sum_{i=2}^N \int \frac{w(x^i)}{w(x) + w(x^i) + \sum_{j=2, j \neq i}^N w(x^j)} \mathbf{1}_{\mathsf{A}}(x^i) \prod_{j=2}^N \lambda(dx^j) \\ &\stackrel{(a)}{\geq} \sum_{i=2}^N \int \pi(dx^i) \mathbf{1}_{\mathsf{A}}(x^i) \int \frac{\lambda(w)}{w(x) + w(x^i) + \sum_{j=2, j \neq i}^N w(x^j)} \prod_{j=2, j \neq i}^N \lambda(dx^j). \end{aligned} \quad (8)$$

612 Here in (a) we used Fubini's theorem together with  $w(x)\lambda(dx) = \pi(dx)\lambda(w)$ . Finally, since the  
613 function  $f: z \mapsto (z + a)^{-1}$  is convex on  $\mathbb{R}_+$  and  $a > 0$ , we get for  $i \in \{2, \dots, N\}$ ,

$$\begin{aligned} &\int \frac{\lambda(w)}{w(x) + w(x^i) + \sum_{j=2, j \neq i}^N w(x^j)} \prod_{j=2, j \neq i}^N \lambda(dx^j) \\ &\geq \frac{\lambda(w)}{\int w(x) + w(x^i) + \sum_{j=2, j \neq i}^N w(x^j) \prod_{j=2, j \neq i}^N \lambda(dx^j)} \\ &\geq \frac{\lambda(w)}{w(x) + w(x^i) + (N - 2)\lambda(w)} \geq \frac{1}{2L + N - 2}. \end{aligned}$$

614 With the bound above we obtain the inequality

$$\mathsf{P}_N(x, \mathsf{A}) \geq \pi(\mathsf{A}) \times \frac{N - 1}{2L + N - 2} = \epsilon_N \pi(\mathsf{A}). \quad (9)$$

615 This means that the whole space  $\mathbb{X}$  is  $(1, \epsilon_N \pi)$ -small (see [21, Definition 9.3.5]). Since  $\mathsf{P}_N(x, \cdot)$  and  
616  $\pi$  are probability measures, (9) implies

$$\|\mathsf{P}_N(x, \cdot) - \pi\|_{\text{TV}} = \sup_{\mathsf{A} \in \mathcal{X}} |\mathsf{P}_N(x, \mathsf{A}) - \pi(\mathsf{A})| \leq 1 - \epsilon_N = \kappa_N.$$

617 The statement follows from [21, Theorem 18.2.4] applied with  $m = 1$ .  $\square$

618 **B.2 Proof of Theorem 2**

619 We preface the proof with some preparatory lemmas.

620 **Lemma 8.** Let  $K \subset \mathbb{X}$ , such that  $w_{\infty, K} := \sup_{x \in K} \{w(x)/\lambda(w)\} < \infty$  and  $\pi(K) > 0$ . Then, for all  
621  $(x, A) \in K \times \mathcal{X}$ , we get that

$$P_N(x, A) \geq \epsilon_{N, K} \pi_K(A),$$

622 with  $\epsilon_{N, K} = (N - 1)\pi(K)/[2w_{\infty, K} + N - 2]$  and  $\pi_K(A) = \pi(A \cap K)/\pi(K)$ .

623 Note that if the weight function  $w$  is upper semi-continuous, then for any compact  $K$ ,  $w_{\infty, K} =$   
624  $\sup_{x \in K} w(x) < \infty$ . Moreover,  $\lim_{N \rightarrow \infty} \epsilon_{N, K} = \pi(K)$ .

625 *Proof.* Let  $(x, A) \in \mathbb{X} \times \mathcal{X}$ . Then, using the lower bound (8), we obtain

$$\begin{aligned} P_N(x, A) &\geq \sum_{i=2}^N \int \pi(dx^i) 1_A(x^i) \int \frac{\lambda(w)}{w(x) + w(x^i) + \sum_{j=2, j \neq i}^N w(x^j)} \prod_{j=2, j \neq i}^N \lambda(dx^j) \\ &\geq (N - 1) \int \pi(dy) 1_A(y) \frac{1}{w(x)/\lambda(w) + w(y)/\lambda(w) + N - 2}, \end{aligned}$$

626 where the last inequality follows from Jensen's inequality and the convexity of the function  $z \mapsto$   
627  $(z + a)^{-1}$  on  $\mathbb{R}_+$ . We conclude by noting that

$$\begin{aligned} P_N(x, A) &\geq (N - 1) \int \pi(dy) 1_{A \cap K}(y) \frac{1}{w(x)/\lambda(w) + w(y)/\lambda(w) + N - 2} \\ &\geq \frac{N - 1}{2w_{\infty, K} + N - 2} \int \pi(dy) 1_{A \cap K}(y) = \frac{(N - 1)\pi(K)}{2w_{\infty, K} + N - 2} \pi_K(A). \end{aligned}$$

628  $\square$

629 **Lemma 9.** Assume A1. Then for all  $x \in \mathbb{X}$ , any function  $V : \mathbb{X} \rightarrow [1, \infty)$  with  $\pi(V) < \infty$ ,  
630  $\lambda(V) < \infty$ , and  $N \geq 3$ , it holds that

$$P_N V(x) \leq V(x) + b_{P_N}, \quad (10)$$

631 where  $b_{P_N}$  is given in (12).

632 Note that

$$b_{P_\infty} := \lim_{N \rightarrow \infty} b_{P_N} = 2\pi(V) + 4 \operatorname{Var}_\lambda[w]/\lambda(V). \quad (11)$$

633 *Proof.* Note first that

$$\begin{aligned} P_N V(x) &= V(x) \int \frac{w(x)}{w(x) + \sum_{j=2}^N w(x^j)} \prod_{j=2}^N \lambda(dx^j) + \int \sum_{i=2}^N \frac{w(x^i)}{w(x) + \sum_{j=2}^N w(x^j)} V(x^i) \prod_{j=2}^N \lambda(dx^j) \\ &\leq V(x) + (N - 1)U_N \end{aligned}$$

634 where we have set

$$U_N = \int \frac{w(x^2)V(x^2)\lambda(dx^2)}{w(x^2) + \sum_{j=3}^N w(x^j)} \prod_{j=3}^N \lambda(dx^j).$$

635 Since the function  $z \mapsto z/(z + a)$  is concave on  $\mathbb{R}_+$  for  $a > 0$ , we have

$$\begin{aligned} \int \frac{w(x^2)}{w(x^2) + \sum_{j=3}^N w(x^j)} V(x^2)\lambda(dx^2) &= \lambda(V) \int \frac{w(x^2)}{w(x^2) + \sum_{j=3}^N w(x^j)} \frac{V(x^2)\lambda(dx^2)}{\lambda(V)} \\ &\leq \lambda(V) \frac{\int w(x^2)V(x^2)\lambda(dx^2)/\lambda(V)}{\int w(x^2)V(x^2)\lambda(dx^2)/\lambda(V) + \sum_{j=3}^N w(x^j)} \leq \frac{\pi(V)\lambda(w)}{\pi(V)\lambda(w)/\lambda(V) + \sum_{j=3}^N w(x^j)}. \end{aligned}$$

636 The bound above implies that, with renormalization,

$$U_N \leq \int \frac{\pi(V)}{\pi(V)/\lambda(V) + \sum_{j=3}^N w(x^j)/\lambda(w)} \prod_{j=3}^N \lambda(dx^j)$$

637 Applying now Lemma 7 with  $a = \pi(V)/\lambda(V)$ ,  $b = 1$ ,  $M = N - 2$ , and  $Y_j = w(x^j)/\lambda(w)$ , we  
638 obtain that

$$U_N \leq \frac{\pi(V)}{\pi(V)/\lambda(V) + (N-2)/2} + \frac{4 \operatorname{Var}_\lambda[w]}{(N-2)\lambda(V)}.$$

639 Combining the bounds above yields (10) with

$$b_{P_N} = \frac{(N-1)\pi(V)}{\pi(V)/\lambda(V) + (N-2)/2} + \frac{4(N-1) \operatorname{Var}_\lambda[w]}{(N-2)\lambda(V)}. \quad (12)$$

640  $\square$

641 **Lemma 10.** Let  $P$  be a Markov kernel on  $(\mathbb{X}, \mathcal{X})$ ,  $\gamma$  be a probability measure on  $(\mathbb{X}, \mathcal{X})$ , and  $\epsilon > 0$ .  
642 Let  $C \in \mathcal{X}$  be an  $(1, \epsilon\gamma)$ -small set for  $P$ . Then for arbitrary Markov kernel  $Q$  on  $(\mathbb{X}, \mathcal{X})$ , the set  $C$  is  
643 an  $(1, \epsilon\gamma_Q)$ -small set for  $PQ$ , where  $\gamma_Q(A) = \int \gamma(dy)Q(y, A)$  for  $A \in \mathcal{X}$ .

644 *Proof.* Let  $(x, A) \in C \times \mathcal{X}$ . Then it holds

$$PQ(x, A) = \int P(x, dy)Q(y, A) \geq \epsilon \int \gamma(dy)Q(y, A) = \epsilon\gamma_Q(A).$$

645  $\square$

646 **Lemma 11.** Let  $P$  and  $Q$  be two irreducible Markov kernels with  $\pi$  as their unique invariant  
647 distribution. Let  $V : \mathbb{X} \rightarrow [1, \infty)$  be a measurable function. Suppose that there exist  $\lambda_Q \in [0, 1)$   
648 and  $b_P, b_Q \in \mathbb{R}_+$  such, that  $PV(x) \leq V(x) + b_P$  and  $QV(x) \leq \lambda_Q V(x) + b_Q$ . Let  $r_0 > 1$ .  
649 Also assume that for all  $r \geq r_0$ , there exist  $\epsilon_r > 0$  and a probability measure  $\gamma_r$  such that for all  
650  $(x, A) \in V_r \times \mathcal{X}$ ,  $P(x, A) \geq \epsilon_r \gamma_r(A)$ , where  $V_r = \{x \in \mathbb{X} : V(x) \leq r\}$ . Define  $K = PQ$  and  
651  $\lambda_K = \lambda_Q, b_K = b_P + b_Q$ . Then,

$$KV(x) \leq \lambda_K V(x) + b_K \text{ and, for all } x \in V_r, K(x, A) \geq \epsilon_r \gamma_{Q,r}(A),$$

where  $\gamma_{Q,r}(A) = \int \gamma_r(dy)Q(y, A)$ . Moreover, let  $r \geq r_0$  be such that  $\lambda_K + 2b_K/(1+r) < 1$ . Then,  
for any  $x \in \mathbb{X}$  and  $k \in \mathbb{N}$ ,

$$\|K^k(x, \cdot) - \pi\|_V \leq c_K \{V(x) + \pi(V)\} \rho_K^k,$$

652 where

$$\rho_K = \frac{\log(1 - \epsilon_r) \log \bar{\lambda}_K}{\log(1 - \epsilon_r) + \log \bar{\lambda}_K - \log \bar{b}_K}, \quad c_K = (\lambda_K + b_K)(1 + \bar{b}_K / [(1 - \epsilon_r)(1 - \bar{\lambda}_K)]),$$

$$\bar{\lambda}_K = \lambda_K + 2b_K/(1+r), \quad \bar{b}_K = \lambda_K r + b_K.$$

653 *Proof.* By Lemma 10, it holds that for any  $(x, A) \in V_r \times \mathcal{X}$ ,  $K(x, A) \geq \epsilon_r \gamma_{Q,r}(A)$ . Moreover, for  
654 any  $x \in \mathbb{X}$ ,  $KV(x) = PQV(x) \leq \lambda_Q PV(x) + b_Q \leq \lambda_Q V(x) + b_Q + b_P$ . The proof is completed  
655 with [21, Theorem 19.4.1].  $\square$

656 *Proof of Theorem 2.* The proof consists of the 3 main steps:

657 1. Lemma 8 implies that for all  $r \geq r_R$ , the level sets  $V_r$  for the Markov kernel  $P_N$  are  $(1, \epsilon_{r,N}\gamma_r)$ -  
658 small for the Markov kernel  $P_N$ , where

$$\epsilon_{r,N} = (N-1)\pi(V_r)/[2w_{\infty,r} + N-2],$$

659 and  $\gamma_r(A) = \int \pi_{V_r}(dy)R(y, A)$  with  $\pi_{V_r}(B) = \pi(B \cap V_r)/\pi(V_r)$ , for any  $B \in \mathcal{X}$ .

660 2. Lemma 9 implies that for all  $x \in \mathbb{X}$ ,  $P_N V(x) \leq V(x) + b_{P_N}$ , where  $b_{P_N}$  is given in (12).

661 3. We finally show (see Lemma 11) that the Markov kernel  $K_N$  also satisfies a Foster-Lyapunov  
662 condition with the same drift function  $V$  as  $R$ , that is,  $K_N V \leq \lambda_R V + b_{K_N}$  with  $b_{K_N} = b_R + b_{P_N}$ .

663 We conclude by using Lemma 11. We choose  $r_N = r_R \vee \{4b_{K_N}/(1-\lambda_R) - 1\}$ . Then  $\lambda_R + 2b_{K_N}/(1 + r_N) \leq (1 + \lambda_R)/2 < 1$ , and Lemma 11 implies (3) with

$$\log \tilde{\kappa}_{K_N} = \frac{\log(1 - \epsilon_{r,N}) \log \bar{\lambda}_{K_N}}{\log(1 - \epsilon_{r,N}) + \log \bar{\lambda}_{K_N} - \log \bar{b}_{K_N}},$$

$$c_{K_N} = (\lambda_R + \bar{b}_{K_N})(1 + \bar{b}_{K_N} / [2(1 - \epsilon_{r,N})(1 - \bar{\lambda}_{K_N})]),$$

$$\bar{\lambda}_{K_N} = (1 + \lambda_R)/2, \quad \bar{b}_{K_N} = \lambda_R r_N + b_{K_N}.$$

665 Set  $b_{K_\infty} = \lim_{N \rightarrow \infty} b_{K_N} = b_R + b_{P_\infty}$ , where  $b_{P_\infty}$  is defined in (11),  $r_\infty = r_R \vee [4b_{K_\infty}/(1-\lambda_R) - 1]$   
666 and  $\epsilon_\infty = \pi(V_{r_\infty})$ . With these notations, we have

$$\begin{aligned}\log \tilde{\kappa}_{K_\infty} &= \frac{\log(1 - \epsilon_\infty) \log \bar{\lambda}_{K_\infty}}{\log(1 - \epsilon_\infty) + \log \bar{\lambda}_{K_\infty} - \log \bar{b}_{K_\infty}}, \\ c_{K_\infty} &= (\lambda_R + \bar{b}_{K_\infty})(1 + \bar{b}_{K_\infty}) / [(1 - \epsilon_\infty)(1 - \bar{\lambda}_{K_\infty})] \\ \bar{\lambda}_{K_\infty} &= (1 + \lambda_R)/2, \quad \bar{b}_{K_\infty} = \lambda_R r_\infty + b_{K_\infty}.\end{aligned}\tag{13}$$

667  $\square$

## 668 C Metropolis-Adjusted Langevin rejuvenation kernel

669 This section addresses the convergence of the Metropolis Adjusted Langevin algorithm (MALA) for  
670 sampling from a positive target probability density  $\pi$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ , where  $\mathcal{B}(\mathbb{R}^d)$  is the Borel  $\sigma$   
671 field of  $\mathbb{R}^d$  endowed with the Euclidean topology. For simplicity, let  $U = -\log \pi$  be the associated  
672 potential function. MALA is a Markov chain Monte Carlo (MCMC) method based on Langevin  
673 diffusion associated with  $\pi$ :

$$d\mathbf{X}_t = -\nabla U(\mathbf{X}_t)dt + \sqrt{2}dB_t, \tag{14}$$

674 where  $(B_t)_{t \geq 0}$  is a  $d$ -dimensional Brownian motion. It is known that under mild conditions this  
675 diffusion admits a strong solution  $(\mathbf{X}_t^{(x)})_{t \geq 0}$  for any starting point  $x \in \mathbb{R}^d$  and defines a Markov  
676 semigroup  $(P_t)_{t \geq 0}$  for any  $t \geq 0$ ,  $x \in \mathbb{R}^d$  and  $A \in \mathcal{B}(\mathbb{R}^d)$  by  $P_t(x, A) = P(P(\mathbf{X}_t^{(x)} \in A))$ .  
677 Moreover, this Markov semigroup admits  $\pi$  as its unique stationary measure, is ergodic and even  
678  $V$ -uniformly geometrically ergodic with additional assumptions on  $U$  (see [63, 46]). However,  
679 sampling a path solution of (14) is a real challenge in most cases, and discretizations are used instead  
680 to obtain a Markov chain with similar long-term behaviour. Here we consider the Euler-Maruyama  
681 discretization, which is given by (14), defined for all  $k \geq 0$  by

$$Y_{k+1} = Y_k - \gamma \nabla U(Y_k) + \sqrt{2\gamma} Z_{k+1}, \tag{15}$$

682 where  $\gamma$  is the step size of the discretization and.  $\{Z_k, k \in \mathbb{N}^*\}$  is a i.i.d. sequence of  $d$ -dimensional  
683 standard Gaussian random variables. This algorithm was proposed by [24, 55] and later studied  
684 by [28, 29, 50, 63]. According to [63], this algorithm is called the Unadjusted Langevin algorithm  
685 (ULA). A drawback of this method is that even if the Markov chain  $\{Y_k, k \in \mathbb{N}\}$  has a unique  
686 stationary distribution  $\pi_\gamma$  and is ergodic (which is guaranteed under mild assumptions about  $U$ ),  $\pi_\gamma$  is  
687 different from  $\pi$  most of the time. To solve this problem, in [65, 63] it is proposed to use the Markov  
688 kernel associated with the recursion defined by the Euler-Maruyama discretization (15) as a proposal  
689 kernel in a Metropolis-Hastings algorithm that defines a new Markov chain  $\{X_k, k \in \mathbb{N}\}$  by:

$$X_{k+1} = X_k + 1_{\mathbb{R}_+}(U_{k+1} - \alpha_\gamma(X_k, \tilde{Y}_{k+1}))\{\tilde{Y}_{k+1} - X_k\}, \tag{16}$$

690 where  $\tilde{Y}_{k+1} = X_k - \gamma \nabla U(X_k) + \sqrt{2\gamma} Z_{k+1}$ ,  $\{U_k, k \in \mathbb{N}^*\}$  is a sequence of i.i.d. uniform random  
691 variables on  $[0, 1]$  and  $\alpha_\gamma : \mathbb{R}^{2d} \rightarrow [0, 1]$  is the usual Metropolis acceptance ratio. This algorithm  
692 is called Metropolis Adjusted Langevin Algorithm (MALA) and has since been used in many  
693 applications.

694 Denote by  $r_\gamma$  the proposal transition density associated to the Euler-Maruyama discretization (15)  
695 with stepsize  $\gamma > 0$ , i.e., for any  $x, y \in \mathbb{R}^d$ ,

$$r_\gamma(x, y) = (4\pi\gamma)^{-d/2} \exp(-(4\gamma)^{-1}\|y - x + \gamma \nabla U(x)\|^2).$$

696 Then, the Markov kernel  $R_\gamma$  of the MALA algorithm (16) is given for  $\gamma > 0$ ,  $x \in \mathbb{R}^d$ , and  $A \in \mathcal{B}(\mathbb{R}^d)$   
697 by

$$\begin{aligned}R_\gamma(x, A) &= \int_{\mathbb{R}^d} 1_A(y) \alpha_\gamma(x, y) r_\gamma(x, y) dy + \delta_x(A) \int_{\mathbb{R}^d} \{1 - \alpha_\gamma(x, y)\} r_\gamma(x, y) dy, \\ \alpha_\gamma(x, y) &= 1 \wedge \frac{\pi(y) r_\gamma(y, x)}{\pi(x) r_\gamma(x, y)}.\end{aligned}\tag{17}$$

698 It is well-known, see e.g. [63], that for any  $\gamma > 0$ ,  $R_\gamma$  is reversible with respect to  $\pi$  and  $\pi$ -irreducible.

699 **H1.** The function  $U : \mathbb{R}^d \rightarrow \mathbb{R}$  is three times continuously differentiable. In addition,  $\nabla U(0) = 0$  and  
700 there exists  $L \geq 0$  and  $M \geq 0$  such that  $\sup_{x \in \mathbb{R}^d} \|D^2 U(x)\| \leq L$  such that  $\sup_{x \in \mathbb{R}^d} \|D^3 U(x)\| \leq M$ .

701 The condition  $\nabla U(0) = 0$  is satisfied (up to a translation) as soon as  $U$  has a local minimum, which  
702 is the case when  $\lim_{\|x\| \rightarrow +\infty} U(x) = +\infty$ , since  $U$  is continuous.

703 **H2.** There exist  $m > 0$  and  $K \geq 0$  such that for any  $x, y \in \mathbb{R}^d$ ,  $\|x\| \geq K$  and  $\|y\| = 1$ ,

$$D^2 U(x)\{y\}^{\otimes 2} \geq m.$$

704 Note that under **H1** and **H2**, for any  $x, y \in \mathbb{R}^d$ ,  $\|y\| = 1$ , it holds that

$$D^2 U(x)\{y\}^{\otimes 2} \geq m - (m + L)1_{B(0,K)}(x).$$

705 In the case  $K = 0$ , **H2** amounts to  $U$  being strongly convex and the convexity constant being equal to  
706  $m$ . However, if  $K > 0$ , **H2** is a slight strengthening of the condition of strong convexity at infinity  
707 considered in [18, 23]: there is  $m' > 0$  and  $K' \geq 0$  such that for each  $x, y \in \mathbb{R}^d$ ,  $\|x - y\| \geq K'$

$$\langle \nabla U(x) - \nabla U(y), x - y \rangle \geq m'\|x - y\|^2. \quad (18)$$

708 Indeed, if (18) holds for any  $x, y \in \mathbb{R}^d$  that  $\|x\| \vee \|y\| \geq K'$  instead of  $\|x - y\| \geq K'$ , then a simple  
709 calculation implies that **H2** holds with  $m \leftarrow m'$  and  $K \leftarrow K' + 1$ . Finally, while the condition (18)  
710 holds for  $x, y \in \mathbb{R}^d$ ,  $\|x - y\| \geq K'$ , is weaker than **H2**, it may be more convenient in many situations  
711 to check whether the latter holds.

712 **Lemma 12.** Assume **H1** and **H2** hold. The function  $U$  satisfies for any  $x \in \mathbb{R}^d$ ,

$$\langle \nabla U(x), x \rangle \geq (m/2)\|x\|^2 - \tilde{C}1_{B(0,\tilde{K})}(x),$$

713 with  $\tilde{K} = 2K(1 + L/m)$  and  $\tilde{C} = L\tilde{K}^2$ .

714 Note that under **H1** and **H2**,  $m \leq L$ . Define for any  $\eta > 0$ ,  $V_\eta : \mathbb{R}^d \rightarrow [1, +\infty)$  for any  $x \in \mathbb{R}^d$  by

$$V_\eta(x) = \exp(\eta\|x\|^2). \quad (19)$$

715 The analysis of MALA is naturally related to the study of the ULA algorithm. More precisely, since  
716 for any  $x \in \mathbb{R}^d$  and  $A \in \mathcal{B}(\mathbb{R}^d)$ , the Markov kernel corresponding to ULA (15) is given by

$$Q_\gamma(x, A) = \int_{\mathbb{R}^d} 1_A(x - \gamma\nabla U(x) + \sqrt{2\gamma}z) g(z) dz.$$

717 To show that MALA satisfies a Lyapunov condition, we first state a drift condition for the ULA  
718 algorithm.

719 **Proposition 13.** Assume **H1** and **H2** and let  $\bar{\gamma} \in (0, m/(4L^2)]$ . Then, for any  $\gamma \in (0, \bar{\gamma}]$ ,  $x \in \mathbb{R}^d$ ,

$$Q_\gamma V_{\bar{\eta}}(x) \leq \exp(-\bar{\eta}m\gamma\|x\|^2/4)V_{\bar{\eta}}(x) + b_{\bar{\eta}, \bar{\gamma}}^U \gamma 1_{B(0, K^U)}(x),$$

720 where  $V_{\bar{\eta}}$  is defined in (19),  $\bar{\eta} = m/16$ ,  $K^U = \max(\tilde{K}, 4\sqrt{d/m})$ ,  $\tilde{K}$  is defined in Lemma 12 and

$$\begin{aligned} b_{\bar{\eta}, \bar{\gamma}}^U &= [\bar{\eta}\{m/4 + (1 + 16\bar{\eta}\bar{\gamma})(4\bar{\eta} + 2L + \bar{\gamma}L^2)\}(K^U)^2 + 4\bar{\eta}d] \\ &\times \exp(\bar{\gamma}\bar{\eta}\{m/4 + (1 + 16\bar{\eta}\bar{\gamma})(4\bar{\eta} + 2L + \bar{\gamma}L^2)\}(K^U)^2 + 4\bar{\eta}\bar{\gamma}d). \end{aligned}$$

721 *Proof.* The proof follows from [22, Proposition 6]. □

722 We now introduce for  $\bar{\gamma} > 0$  the auxiliary constant

$$C_{1, \bar{\gamma}} = 2(2^{1/2}M \vee \bar{\gamma}^{1/2}ML \vee 2L^2[1 \vee \bar{\gamma}^{1/2} \vee \bar{\gamma}L \vee (\bar{\gamma}L^{4/3})^{3/2}]). \quad (20)$$

723 For  $\bar{\gamma} \in (0, m^3/(4L^4)]$ , we also define  $C_{2, \bar{\gamma}}$  as

$$C_{2, \bar{\gamma}} = 2L + (\bar{\gamma}/2)L^2 + 2^{-3/2}\bar{\gamma}^{3/2}L^3 + \{2^{1/2}L^2 + (2^{1/2}L^2 + 2^{-3/2}\bar{\gamma}^{1/2})L^3\}^2(2^4/m^3).$$

724 Using Proposition 13, we state a drift condition for the MALA kernel  $R_\gamma$ .

725 **Proposition 14.** Assume **H 1** and **H 2**. Then, there exist  $\Gamma > 0$  (given in (21)) such that for any  
 726  $\bar{\gamma} \in (0, \Gamma]$ ,  $\gamma \in (0, \bar{\gamma}]$  and  $x \in \mathbb{R}^d$ ,

$$R_\gamma V_{\bar{\eta}}(x) \leq (1 - \varpi\gamma)V_{\bar{\eta}}(x) + b_{\bar{\eta}, \bar{\gamma}}^M \gamma 1_{B(0, K^M)}(x),$$

727 where  $V_{\bar{\eta}}$  is defined by (19),  $R_\gamma$  is the Markov kernel of MALA defined by (17),  $\bar{\eta} = \mathfrak{m}/16$ ,  $\varpi =$   
 728  $\bar{\eta}\mathfrak{m}(K^M)^2/16$ , and

$$\begin{aligned} \Gamma_{1/2} &= \min(1, \mathfrak{m}^3/(4L^4), d^{-1}) , \quad \Gamma = \min(\Gamma_{1/2}, 4/\{\mathfrak{m}\bar{\eta}(K^M)^2\}) , \\ K^M &= \max(2^4, 2K, K^U, \tilde{K}, 4b_{1/2}^{1/2}/(\mathfrak{m}\bar{\eta})^{1/2}) , \quad b_{1/2} = C_{2, \Gamma_{1/2}} d + \sup_{u \geq 1} \{ue^{-u/2^7}\} , \\ b_{\bar{\eta}, \bar{\gamma}}^M &= b_{\bar{\eta}, \bar{\gamma}}^U + \bar{\eta}\mathfrak{m}(K^M)^2 e^{\bar{\eta}(K^M)^2}/16 + C_{1, \bar{\gamma}} \bar{\gamma}^{1/2} \left\{ d + \sqrt{3}d^2 + (K^M)^2 \right\} , \end{aligned} \quad (21)$$

729 where  $K^U, b_{\bar{\eta}, \bar{\gamma}}^U$  are defined in Proposition 13, and  $\tilde{K}$  is defined in Lemma 12.

730 *Proof.* The proof follows from [22, Proposition 7].  $\square$

731 Quantitative bound on the mixing rate of the MALA sampler requires also the *minorization condition*  
 732 for the MALA kernel. The result below is due to [22, Proposition 12].

733 **Proposition 15.** Assume **H 1** and **H 2**. Then for any  $K \geq 0$  there exists  $\tilde{\Gamma}_K > 0$  (given in (22) below),  
 734 such that for any  $x, y \in \mathbb{R}^d$ ,  $\|x\| \vee \|y\| \leq K$ , and  $\gamma \in (0, \tilde{\Gamma}_K]$  we have

$$\|\delta_x R^{[1/\gamma]} - \delta_y R^{[1/\gamma]}\|_{\text{TV}} \leq 2(1 - \varepsilon(K)/2) ,$$

735 where

$$\begin{aligned} \varepsilon(K) &= 2\Phi\left(-(1 + 1/L)^{1/2}(3L)^{1/2}K\right) , \quad \tilde{\Gamma}_{1/2} = \mathfrak{m}/(4L^2) , \\ \tilde{\Gamma}_K &= \tilde{\Gamma}_{1/2} \wedge \left[ \frac{\varepsilon(K)}{2C_{1, \tilde{\Gamma}_{1/2}}(d + \sqrt{3}d^2 + K^2 + 2\tilde{b}_{\tilde{\Gamma}_{1/2}}^U/\mathfrak{m})} \right]^2 , \\ \tilde{b}_{\tilde{\Gamma}_{1/2}}^U &= 2d + [\max(\tilde{K}, 2\sqrt{(2d)/\mathfrak{m}})]^2 \left( \tilde{\Gamma}_{1/2} L^2 + 2L + \mathfrak{m}/2 \right) , \end{aligned} \quad (22)$$

736 where  $C_{1, \tilde{\Gamma}_{1/2}}$  is defined in (20),  $\tilde{K}$  is defined in Lemma 12, and  $\Phi(\cdot)$  is the cumulative distribution  
 737 function of the Gaussian distribution with zero mean and unit variance on  $\mathbb{R}$ .

738 Combining Proposition 14 and Proposition 15 yields the following ergodicity result in  $V_{\bar{\eta}}$ -norm.

739 **Theorem 16.** Assume **H 1** and **H 2**. Then, there exist  $\bar{\Gamma} > 0$  (defined in (23) below), such that for any  
 740  $\gamma \in (0, \bar{\Gamma}]$ , there exist  $C_{\bar{\Gamma}} \geq 0$  and  $\rho_{\bar{\Gamma}} \in [0, 1)$  (given in (23)) satisfying for any  $x \in \mathbb{R}^d$ ,

$$\|\delta_x R_\gamma^k - \pi\|_{V_{\bar{\eta}}} \leq C_{\bar{\Gamma}} \rho_{\bar{\Gamma}}^{\gamma k} \{V_{\bar{\eta}}(x) + \pi(V_{\bar{\eta}})\} ,$$

741 where  $\bar{\eta} = \mathfrak{m}/16$ ,

$$\begin{aligned} \log \rho_{\bar{\Gamma}} &= \frac{\log(1 - 2^{-1}\varepsilon(K_{\bar{\Gamma}})) \log \bar{\lambda}}{\log(1 - 2^{-1}\varepsilon(K_{\bar{\Gamma}})) + \log \bar{\lambda} - \log \bar{b}_{\bar{\eta}, \bar{\Gamma}}^M} , \\ \bar{\lambda} &= (1 + \lambda)/2 , \quad \lambda = e^{-\varpi} , \quad \bar{b}_{\bar{\eta}, \bar{\Gamma}}^M = \lambda b_{\bar{\eta}, \bar{\Gamma}}^M + M_{\bar{\Gamma}} , \quad \bar{\Gamma} = \Gamma \wedge \tilde{\Gamma}_{K_\Gamma} , \\ M_{\bar{\gamma}} &= \left( \frac{4b_{\bar{\eta}, \bar{\gamma}}^M(1 + \bar{\gamma})}{1 - \lambda} \right) \vee 1 , \quad K_{\bar{\gamma}} = (\log(M_{\bar{\gamma}})/\bar{\eta})^{1/2} , \quad \bar{\gamma} \in \{\bar{\Gamma}, \Gamma\} , \\ C_{\bar{\Gamma}} &= \rho_{\bar{\Gamma}}^{-1} \{ \lambda + 1 \} \{ 1 + \bar{b}_{\bar{\eta}, \bar{\Gamma}}^M / [1 - 2^{-1}\varepsilon(K_{\bar{\Gamma}})(1 - \bar{\lambda})] \} , \end{aligned} \quad (23)$$

742 and  $\varpi$  is given in Proposition 14.

743 *Proof.* The proof follows from [22, Theorem 2]. For completeness we repeat here the main steps of  
 744 the proof. Proposition 14 shows that there exist  $\Gamma > 0$  (given in (21)) such that for any  $\bar{\gamma} \in (0, \Gamma]$ ,  
 745  $\gamma \in (0, \bar{\gamma}]$  and  $x \in \mathbb{R}^d$ ,

$$R_\gamma V_{\bar{\eta}}(x) \leq (1 - \varpi\gamma)V_{\bar{\eta}}(x) + b_{\bar{\eta}, \bar{\gamma}}^M \gamma ,$$

746 where the constants  $\varpi$  and  $b_{\bar{\eta}, \bar{\gamma}}^M$  are given in Proposition 14. Hence, setting  $\lambda = e^{-\varpi} < 1$ , we obtain  
 747 by induction that

$$R_\gamma^{[1/\gamma]} V_{\bar{\eta}}(x) \leq \lambda V_{\bar{\eta}}(x) + b_{\bar{\eta}, \bar{\gamma}}^M (1 + \bar{\gamma}) .$$

748 Now we set  $M_{\bar{\gamma}}$  and  $K_{\bar{\gamma}}$  as in (23). Then Proposition 15 implies that for any  $\bar{\gamma} \in (0, \tilde{\Gamma}_{K_\Gamma}]$ , any  
 749  $x, y \in \{V_{\bar{\eta}}(\cdot) \leq M_{\bar{\gamma}}\}$ , and  $\gamma \in (0, \bar{\gamma}]$ ,

$$\|\delta_x R_\gamma^{[1/\gamma]} - \delta_y R_\gamma^{[1/\gamma]}\|_{\text{TV}} \leq 2(1 - \varepsilon(K_{\bar{\gamma}})) .$$

750 Now it remains to combine both statements with  $\bar{\gamma} = \Gamma \wedge \tilde{\Gamma}_{K_\Gamma}$  and apply [21, Theorem 19.4.1] to the  
 751 Markov kernel  $R_\gamma^{[1/\gamma]}$ .  $\square$

752 **Comparison with Ex<sup>2</sup>MCMC kernel.** Based on the results above, we first state the quantitative  
 753 mixing rate bounds for Ex<sup>2</sup>MCMC algorithm with MALA kernel  $R_\gamma(x, \cdot)$  applied as rejuvenation  
 754 kernel. The corresponding Markov kernel writes for  $x \in \mathbb{R}^d$  and  $A \in \mathcal{B}(\mathbb{R}^d)$  as

$$K_{N,\gamma}(x, A) = P_N R_\gamma(x, A) = \int P_N(x, dy) R_\gamma(y, A) ,$$

755 where  $R_\gamma(x, A)$  is defined in (17). Note also that, for  $r \geq 1$ , and  $V_{\bar{\eta}}$  defined in (19), the level sets

$$V_{\bar{\eta},r} = \{x : V_{\bar{\eta}}(x) \leq r\} = \{x : \|x\| \leq \sqrt{\log r / \bar{\eta}}\} .$$

756 The result above allows to state the following ergodicity result for  $K_{N,\gamma}$  kernel.

757 **Theorem 17.** Assume **H 1**, **H 2**, and **A1,A2** with  $V_{\bar{\eta}}$  defined in (19). Then there exist  $\bar{\Gamma}$  (defined in  
 758 (23)), such that for any  $\gamma \in (0, \bar{\Gamma}]$ ,  $x \in \mathbb{R}^d$ , and  $k \in \mathbb{N}$ ,

$$\|K_{N,\gamma}^k(x, \cdot) - \pi\|_V \leq c_{N,\gamma} \{\pi(V_{\bar{\eta}}) + V_{\bar{\eta}}(x)\} \tilde{\kappa}_{N,\gamma}^k ,$$

759 where  $V_{\bar{\eta}}$  is defined in (19), and the constants  $c_{N,\gamma}, \tilde{\kappa}_{N,\gamma} \in [0, 1]$  are given by

$$\log \tilde{\kappa}_{N,\gamma} = \frac{\log(1 - \epsilon_{r_N, N}) \log \bar{\lambda}}{\log(1 - \epsilon_{r_N, N}) + \log \bar{\lambda} - \log \bar{b}_{N,\gamma}} , \quad r_N = 1 \vee \{4b_{N,\gamma}/(1 - \lambda) - 1\} , \quad (24)$$

$$\epsilon_{r_N, N} = (N - 1)\pi(V_{\bar{\eta},r_N})/[2w_{\infty, r_N} + N - 2], \quad b_{N,\gamma} = b_{P_N} + \gamma \bar{b}_{\bar{\eta}, \bar{\Gamma}}^M ,$$

$$c_{N,\gamma} = (\lambda + \bar{b}_{N,\gamma})(1 + \bar{b}_{N,\gamma}/[2(1 - \epsilon_{r_N, N})(1 - \bar{\lambda})])$$

$$\bar{\lambda} = (1 + \lambda)/2 , \quad \bar{b}_{N,\gamma} = \lambda r_N + b_{N,\gamma} ,$$

760 where  $\lambda$  is defined in (23).

761 *Proof.* The proof follows from the combination of Theorem 2 and Proposition 14.  $\square$

762 Note that the result of Theorem 17 does not require to identify the small sets of the MALA rejuvenation  
 763 kernel  $R_\gamma(x, A)$ . Theorem 16 implies that the rate of convergence of MALA is  $\gamma \log \rho_{\bar{\Gamma}}$ . The  
 764 following statement allows to quantify the improvement in the convergence rate of  $K_{N,\gamma}(x, A)$   
 765 compared to  $R_\gamma(x, A)$ :

766 **Theorem 18.** Assume **H 1**, **H 2**, and assume additionally **A1,A2** with  $V_{\bar{\eta}}$  defined in (19). Then there  
 767 exist  $\bar{\Gamma}$  (defined in (23)), such that for any  $\gamma \in (0, \bar{\Gamma}]$ , it holds that

$$\lim_{N \rightarrow \infty} \frac{\gamma \log \rho_{\bar{\Gamma}}}{\log \tilde{\kappa}_{N,\gamma}} = C_{\infty,\gamma} \gamma \log(1 - 2^{-1} \varepsilon(K_{\bar{\Gamma}})) ,$$

768 where

$$C_{\infty,\gamma} = \frac{\log(1 - \epsilon_\infty) + \log \bar{\lambda} - \log \bar{b}_{\infty,\gamma}}{\{\log(1 - \epsilon_\infty)\}(\log(1 - 2^{-1} \varepsilon(K_{\bar{\Gamma}})) + \log \bar{\lambda} - \log \bar{b}_{\bar{\eta}, \bar{\Gamma}}^M)} , \quad \epsilon_\infty = \pi(V_{\bar{\eta},r_\infty}) ,$$

$$r_\infty = 1 \vee \{4b_{\infty,\gamma}/(1 - \lambda) - 1\} , \quad b_{\infty,\gamma} = b_{K_\infty} + \gamma \bar{b}_{\bar{\eta}, \bar{\Gamma}}^M , \quad \bar{b}_{\infty,\gamma} = \lambda r_\infty + b_{\infty,\gamma} ,$$

769 where  $\lambda, \bar{\lambda}$ , and  $\bar{b}_{\bar{\eta}, \bar{\Gamma}}^M$  are defined in (23).

770 *Proof.* The proof follows by combining the expressions (23) and (24).  $\square$

771 **Remark 19.** Note that the expressions for  $\bar{b}_{\bar{\eta}, \bar{\Gamma}}^M > 1$  and  $\bar{b}_{\infty, \gamma} > 1$  imply that  $|C_{\infty, \gamma}|$  remains  
 772 bounded for  $\gamma \in (0, \bar{\Gamma}]$  when  $\bar{\Gamma} \rightarrow 0$ . This regime is typical for high-dimensional scenarios, since  
 773  $\bar{\Gamma} \leq \Gamma \leq d^{-1}$ .

## 774 D Proof of Theorem 3

775 The proof relies on results of stochastic approximation with Markovian dynamics; see [7, 8]. For  
 776 reader's convenience, before going into the details, we give an outline of the proof. The motivation of  
 777 such algorithms is to find the roots of the function  $h : \Theta \rightarrow \mathbb{R}^q$ ,  $\Theta \subset \mathbb{R}^q$

$$h(\theta) = \int_{\mathbb{U} \times \mathbb{E}} H(\theta, u, e) \mu(de) \rho_\theta(de),$$

778 for families of functions  $\{H(\theta, u, e) : \Theta \times \mathbb{U} \times \mathbb{E} \rightarrow \Theta\}$ , a family of probability distributions  
 779  $\{\rho_\theta, \theta \in \Theta\}$  of  $(\mathbb{E}, \mathcal{E})$  and a probability distribution  $\mu$  on a space  $(\mathbb{U}, \mathcal{U})$ . These roots are not available  
 780 analytically and a way of finding them numerically consists of considering the controlled Markov  
 781 chain on  $\{(\Theta \times \mathbb{U})^\mathbb{N}, (\mathcal{B}(\Theta) \otimes \mathcal{U})^{\otimes \mathbb{N}}\}$  initialized at some  $(\theta_0, U_0) = (\vartheta, u) \in \Theta \times \mathbb{U}$  and defined  
 782 recursively for a sequence of stepsize  $\{\gamma_i, i \in \mathbb{N}\}$  by

$$\begin{aligned} U_{i+1} &\sim P_{\theta_i}(U_i, \cdot), \quad E_{i+1} \sim \rho_{\theta_i} \\ \theta_{i+1} &= \theta_i + \gamma_{i+1} H(\theta_i, U_{i+1}, E_{i+1}), \end{aligned}$$

783 where  $\{P_\theta, \theta \in \Theta\}$  is a family of Markov kernels such that for each  $\theta \in \Theta$ ,  $\mu P_\theta = \mu$ . The rationale  
 784 for this recursion goes as follows. Let us first rewrite the Robbins-Monro recursion

$$\theta_{i+1} = \theta_i + \gamma_{i+1} \{h(\theta_i) + \xi_{i+1}\},$$

785 where  $\xi_{i+1} = H(\theta_i, U_{i+1}, E_{i+1})$  is referred to as the "noise". Therefore,  $\{\theta_i\}$  is a noisy version  
 786 of the sequence  $\{\bar{\theta}_i\}$  defined as  $\bar{\theta}_{i+1} = \bar{\theta}_i + \gamma_{i+1} h(\bar{\theta}_i)$ . The convergence of such sequences has  
 787 been studied by many authors, starting with [47] under various conditions. A crucial step of such  
 788 convergence analysis consists of assuming that the sequence  $\{\theta_i\}$  remains bounded with probability 1  
 789 in a compact set of  $\Theta$ . This problem has traditionally can be circumvented by means of modifications  
 790 of the recursion. Indeed, one of the major difficulties specific to the Markovian dynamic scenario is  
 791 that  $\{\theta_i\}$  governs the ergodicity of the controlled Markov chain  $\{U_i\}$  and that stability properties of  
 792  $\{\theta_i\}$  require "good" ergodicity properties which might vanish whenever  $\{\theta_i\}$  approaches  $\partial\Theta$  often  
 793 away from the roots of  $h(\theta)$ , resulting in instability. Most existing results rely on modifications of the  
 794 updates designed to ensure a form of ergodicity of  $\{\xi_i\}$  which in turn ensures that  $\{\theta_i\}$  inherits the  
 795 stability properties of  $\{\bar{\theta}_i\}$ ; see e.g. [7, 10] and the discussion in [8, Section 3]. We follow here [10].  
 796 Let  $\{\mathcal{R}_i\}$  be a sequence of compact subsets of  $\Theta$  and consider the recursion:

$$\begin{aligned} U_{i+1} &\sim P_{\theta_i}(U_i, \cdot) \quad E_{i+1} \sim \rho_{\theta_i} \\ \theta_{i+1}^* &= \theta_i + \gamma_{i+1} H(\theta_i, U_{i+1}, E_{i+1}) \\ \theta_{i+1} &= \theta_{i+1}^* 1_{\mathcal{R}_{i+1}}(\theta_{i+1}^*) + \theta_{i+1}^{\text{proj}} 1_{\mathcal{R}_{i+1}^c}(\theta_{i+1}^*) \end{aligned}$$

797 where, denoting  $\mathcal{F}_i = \sigma(U_0, \theta_j, j \leq i)$ ,  $\theta_{i+1}^{\text{proj}}$  is a random variable measurable w.r.t  $\mathcal{F}_i \vee \sigma(\theta_{i+1}^*)$ .  
 798 Most common practical projection mechanisms include  $\theta_{i+1}^{\text{proj}} = \theta_i$ , 'rejecting' an update outside the  
 799 current feasible set, and  $\theta_{i+1}^{\text{proj}} = \text{Proj}_{\mathcal{R}_{i+1}}(\theta_{i+1}^*)$ , where Proj is a measurable mapping  $\Theta \setminus \mathcal{R}_{i+1} \rightarrow$   
 800  $\mathcal{R}_{i+1}$ . In words, the expanding projections approach only ensures that  $\theta_i$  is in a feasible set  $\mathcal{R}_i$  but  
 801 does not involve potentially harmful 'restarts' as is the case with the adaptive reprojection strategy  
 802 of [7]. We use the results in [10] to show that the SA  $\{\theta_i\}$  'stays away' from  $\partial\Theta$  with probability  
 803 one for any initialization  $(\theta_0, u) \in \mathcal{R}_0 \times \mathbb{U}$  under appropriate conditions on  $\{H(\theta, u, e), (\theta, u, e) \in$   
 804  $\Theta \times \mathbb{U} \times \mathbb{E}\}$ ,  $\{P_\theta, \theta \in \Theta\}$  and  $\{\mathcal{R}_i\}$ . We denote throughout the probability distribution associated  
 805 to the process  $(\theta_i, U_i)_{i \geq 0}$  defined in Algorithm 1.1 and starting at  $(\theta_0, U_0) \equiv (\theta, u) \in \Theta \times \mathbb{U}$  as  
 806  $\mathbb{P}_{\theta, u}(\cdot)$  and the associated expectation as  $\mathbb{E}_{\theta, u}[\cdot]$ . The approach developed in [10] relies on the  
 807 existence of a Lyapunov function  $w : \Theta \rightarrow [0, \infty)$  for the recursion on  $\theta$  and the subsequent proof  
 808 that  $\{w(\theta_i)\}$  is  $\mathbb{P}_{\theta, u}$ -a.s. under some adequate level. For any  $M > 0$ , we define the level sets  
 809  $\mathcal{W}_M := \{\theta \in \Theta : w(\theta) \leq M\}$ . Consider the following assumptions:

810 **SA1.** There exists a continuously differentiable function  $w : \Theta \rightarrow [0, \infty)$  such that

811 (i) For all  $\theta, \theta' \in \Theta$ ,

$$\|\nabla w(\theta) - \nabla w(\theta')\| \leq C_w \|\theta - \theta'\|.$$

(ii) the projection sets are increasing subsets of  $\Theta$ , that is,  $\mathcal{R}_i \subset \mathcal{R}_{i+1}$  for all  $i \geq 0$ , and

$$\hat{\Theta} := \bigcup_{i=0}^{\infty} \mathcal{R}_i \subset \Theta,$$

(iii) there exists a constant  $M_0 > 0$  such that for any  $\theta \in \mathcal{W}_{M_0}^c \cap \hat{\Theta}$

$$\langle \nabla w(\theta), h(\theta) \rangle \leq 0$$

(iv) the family of random variables  $\{\theta_i^{\text{proj}}\}_{i \geq 1}$  satisfies for all  $i \geq 1$  whenever  $\theta_i^* \notin \mathcal{R}_i$

$$\theta_i^{\text{proj}} \in \mathcal{R}_i \quad \text{and} \quad w(\theta_i^{\text{proj}}) \leq w(\theta_i^*) \quad \mathbb{P}_{\theta, u} - \text{a.s.}.$$

812 (v) there exists constants  $c \in [0, \infty)$  and a non-decreasing sequence of constants  $\zeta_i \in [1, \infty)$   
813 satisfying  $\sup_{\theta \in \mathcal{R}_i} |\nabla w(\theta)| \leq c \zeta_i$  for all  $i \geq 0$ .

814 Hereafter, we denote the 'centred' version of  $\bar{H}(\theta, u, e) := H(\theta, u, e) - h(\theta)$ . For the stability  
815 results, we shall introduce the following general condition on the noise sequence. In general terms, it  
816 is related to the rate at which  $\{\theta_i\}$  may approach  $\partial\hat{\Theta}$  in relation to the growth of  $\|H(\theta, u, e)\|$  and  
817 the loss of ergodicity of  $P_\theta$ .

818 **SA2.** For any  $(\theta, u) \in \mathcal{R}_0 \times \mathbb{U}$  it holds that

819 (i)  $\mathbb{P}_{\theta, u}(\lim_{i \rightarrow \infty} \gamma_{i+1} \|\nabla w(\theta_i)\| \cdot \|H(\theta_i, U_{i+1}, E_{i+1})\| = 0) = 1$ ,

820 (ii)  $\mathbb{E}_{\theta, u} \left[ \sum_{i=0}^{\infty} \gamma_{i+1}^2 \|H(\theta_i, U_{i+1}, E_{i+1})\|^2 \right] < \infty$ ,

821 (iii)  $\mathbb{E}_{\theta, u} \left[ \sup_{k \geq 0} \left| \sum_{i=0}^k \gamma_{i+1} \langle \nabla w(\theta_i), \bar{H}(\theta_i, U_{i+1}, E_{i+1}) \rangle \right| \right] < \infty$ .

822 (iv)  $\lim_{\theta \rightarrow \partial\hat{\Theta}} w(\theta) = \infty$

**Theorem 20.** Assume SA1-SA2. Then, for any  $(\theta, u) \in \mathcal{R}_0 \times \mathbb{U}$

$$\mathbb{P}_{\theta, u} \left( \limsup_{i \rightarrow \infty} w(\theta_i) < \infty \right) = 1.$$

823 *Proof.* The proof is a simple adaptation of [10, Theorem 2.5].  $\square$

824 The condition  $\lim_{\theta \rightarrow \partial\hat{\Theta}} w(\theta) = \infty$  is weakened in [10, Section 2.2]. Verifiable conditions implying  
825 **SA2** are given in [10, Section 3, Condition 3.1]. They are summarized in the next assumption. In the  
826 assumptions below, it is implicitly assumed that **SA1** holds with constants  $(\zeta_i)_{i \geq 0}$ .

827 We denote  $\tilde{H}(\theta, u) = \int \bar{H}(\theta, u, e) \rho(de)$  and we consider the following assumptions:

828 **SA3.** For all  $\theta \in \hat{\Theta}$ , the solution  $g_\theta : \mathbb{U} \rightarrow \Theta$  to the Poisson equation  $g_\theta(u) - P_\theta g_\theta(u) \equiv \tilde{H}(\theta, u)$   
829 exists and for all  $i \geq 0$  the step size  $\Gamma_{i+1}$  is independent of  $\mathcal{F}_i$  and  $U_{i+1}$ . Moreover, there exist a  
830 measurable function  $V : \mathbb{U} \rightarrow [1, \infty)$  and constants  $c < \infty, \beta_H, \beta_g \in [0, 1/2]$  and  $\alpha_g, \alpha_H, \alpha_V \in$   
831  $[0, \infty)$  such that for all  $(\theta, u) \in \mathcal{R}_0 \times \mathbb{U}$

832 (i)  $\sup_{\theta \in \mathcal{R}_i} |\tilde{H}(\theta, u)| \leq c \zeta_i^{\alpha_H} V^{\beta_H}(u)$ ,

833 (ii)  $\mathbb{E}_{\theta, u} [V(U_i)] \leq c \zeta_i^{\alpha_V} V(u)$ ,

834 (iii)  $\sup_{\theta \in \mathcal{R}_i} [|g_\theta(u)| + |P_\theta g_\theta(u)|] \leq c \zeta_i^{\alpha_g} V^{\beta_g}(u)$ ,

835 (iv)  $\sum_{i=1}^{\infty} \gamma_{i+1} \zeta_i \mathbb{E}_{\theta, u} [|P_{\theta_i} g_{\theta_i}(U_i) - P_{\theta_{i-1}} g_{\theta_{i-1}}(U_i)|] < \infty$ ,

836 (v)  $\sum_{i=1}^{\infty} \gamma_i^2 \zeta_i^{2+2((\alpha_H+\beta_H\alpha_V)\vee(\alpha_g+\beta_g\alpha_V))} < \infty$ ,

837 (vi)  $\sum_{i=1}^{\infty} \gamma_{i+1} \gamma_i \zeta_i^{\alpha_H+\alpha_g+(\beta_H+\beta_g)\alpha_V} < \infty$ ,

838 (vii)  $\sum_{i=1}^{\infty} |\gamma_{i+1} - \gamma_i| \zeta_i^{1+\alpha_g+\beta_g\alpha_V} < \infty$ .

839 For geometrically ergodic Markov chain, these conditions may be shown to boil down to "uniform-in-  
840  $\theta$ " geometric ergodicity conditions and "smoothness" of the mapping  $\theta \mapsto P_\theta$ .

**MC1.** For any  $r \in (0, 1]$  and any  $\theta \in \hat{\Theta}$ , there exist constants  $M_{\theta,r} \in [0, \infty)$  and  $\rho_{\theta,r} \in (0, 1)$ , such that for any function  $\|f\|_{V^r} < \infty$

$$|P_\theta^k f(u) - \mu_\theta(f)| \leq V^r(u) \|f\|_{V^r} M_{\theta,r} \rho_{\theta,r}^k$$

841 for all  $k \geq 0$  and all  $u \in \mathbb{U}$ . Moreover, it holds that  $\sup_{\theta \in \mathcal{R}_i} M_{\theta,r} \leq c_r \zeta_i^{\alpha_M}$  and  
842  $\sup_{\theta \in \mathcal{R}_i} (1 - \rho_{\theta,r})^{-1} \leq c_r \zeta_i^{\alpha_\rho}$ .

**MC2.** For any  $\theta, \theta' \in \hat{\Theta}$ , there exist a constant  $D_{\theta,\theta',r} \in [0, \infty)$  and a constant  $\beta_D \in (0, \infty)$  independent of  $\theta, \theta'$  and  $r$  such that for any function  $\|f\|_{V^r} < \infty$

$$\|P_\theta f - P_{\theta'} f\|_{V^r} \leq \|f\|_{V^r} D_{\theta,\theta',r} |\theta - \theta'|^{\beta_D}.$$

843 Moreover,  $\sup_{(\theta, \theta') \in \mathcal{R}_i^2} D_{\theta,\theta',r} \leq c_r^D \zeta_i^{\alpha_D}$  for some constant  $c_r^D \in [0, \infty)$  depending only on  $r \in$   
844  $(0, 1]$ .

**MC3.** SA3-(i) and (ii) hold with constants  $\alpha_H, \beta_H$  and  $\alpha_V$ , and there exist constants  $c < \infty, \alpha_\Delta \in [0, \infty)$  and  $\beta_\Delta > 0$  such that

$$\sup_{(\theta, \theta') \in \mathcal{R}_i^2} \left\| \tilde{H}(\theta, \cdot) - \tilde{H}(\theta', \cdot) \right\|_{V^{\beta_H}} \leq c \zeta_i^{\alpha_\Delta} |\theta - \theta'|^{\beta_\Delta}.$$

845 Up to this point, we have only considered the stability of the stochastic approximation process with  
846 expanding projections. Indeed, after showing the stability we know that the projections can occur  
847 only finitely often (almost surely), and the noise sequence can typically be controlled. Given this,  
848 the stochastic approximation literature provides several alternatives to show the convergence; see  
849 [40, 15]. We formulate below a convergence result following from [7].

850 **SA4.** The set  $\Theta \subset \mathbb{R}^d$  is open, the mean field  $h : \Theta \rightarrow \mathbb{R}^d$  is continuous, and there exists a  
851 continuously differentiable function  $\hat{w} : \Theta \rightarrow [0, \infty)$  such that

(i) there exists a constant  $M_0 > 0$  such that

$$\mathcal{L} := \{\theta \in \Theta : \langle \nabla \hat{w}(\theta), h(\theta) \rangle = 0\} \subset \{\theta \in \Theta : \hat{w}(\theta) < M_0\}$$

852 (ii) there exists  $M_1 \in (M_0, \infty]$  such that  $\{\theta \in \Theta : \hat{w}(\theta) \leq M_1\}$  is compact.  
853 (iii) for all  $\theta \in \Theta \setminus \mathcal{L}$ , the inner product  $\langle \nabla \hat{w}(\theta), h(\theta) \rangle < 0$  and the closure of  $\hat{w}(\mathcal{L})$  has an empty  
854 interior.

855 **Theorem 21.** Assume SA4 holds, and let  $\mathcal{K} \subset \Theta$  be a compact set intersecting  $\mathcal{L}$ , that is,  $\mathcal{K} \cap \mathcal{L} \neq$   
856  $\emptyset$ . Suppose that  $(\gamma_i)_{i \geq 1}$  is a sequence of non-negative real numbers satisfying  $\lim_{i \rightarrow \infty} \gamma_i = 0$   
857 and  $\sum_{i=1}^{\infty} \gamma_i = \infty$ . Consider the sequence  $(\theta_i)_{i \geq 0}$  taking values in  $\Theta$  and defined through the  
858 recursion  $\theta_i = \theta_{i-1} - 1 + \gamma_i h(\theta_{i-1}) + \gamma_i \varepsilon_i$  for all  $i \geq 1$ , where  $(\varepsilon_i)_{i \geq 1}$  take values in  $\mathbb{R}^d$ . If  
859 there exists an integer  $i_0$  such that  $\{\theta_i\}_{i \geq i_0} \subset \mathcal{K}$  and  $\lim_{m \rightarrow \infty} \sup_{n \geq m} |\sum_{i=m}^n \gamma_i \varepsilon_i| = 0$ , then  
860  $\lim_{n \rightarrow \infty} \inf_{x \in \mathcal{L} \cap \mathcal{K}} \|\theta_n - x\| = 0$ .

861 We have now all the necessary elements to prove Theorem 3. For simplicity, we set  $\alpha_k = \alpha_\infty$  for  
862 any  $k \in \mathbb{N}$  and  $\gamma_k = 1/(1+k)^\iota$  where  $\iota \in (1/2, 1]$ . In this case, the state space is  $\mathbb{U} = \mathbb{X}^M$  and  
863  $\mathbb{E} = \mathbb{Z}^{(N-1) \cdot M}$ ,  $U_k = (Y_k[j])_{j=1}^M$ ,  $E_k = (Z_k^{2:N}[j])_{j=2}^N$ . With  $u = (y[j])_{j=1}^M$  and  $e = (z^{2:N}[j])_{j=1}^M$ ,  
864  $H(\theta, u, e)$  is given by

$$H(\theta, u) = M^{-1} \sum_{mj=1}^N \{\alpha_\infty H^f(\theta, y[j], z^{2:N}[j]) + (1 - \alpha_\infty) H^b(\theta, z^{2:N}[j])\}.$$

865 where  $H^f$  and  $H^b$  are defined respectively in (4) and (5). In this case, the Markov kernel  $P_\theta$  is given  
866 for any nonnegative function  $f$ ,

$$P_{\theta,N} f(y[1], \dots, y[M]) = \int \prod_{j=1}^N K_{\theta,N}(y[j], dy[j]) f(\tilde{y}[1], \dots, \tilde{y}[M]),$$

867 and  $K_{\theta,N}$  is defined in (2.2) with  $\lambda \leftarrow \lambda_\theta$  and  $w \leftarrow w_\theta$ . By construction, for any  $\theta \in \Theta$ ,  $P_\theta$  has a  
868 unique stationary distribution which is given by  $\mu = \pi^{\otimes M}$ . Using Theorem 6, and, for all  $\theta \in \Theta$ ,

$$H^f(\theta, x^{1:N}) = \Pi_{\theta,N} [\nabla_\theta \log \lambda_\theta](x^{1:N})$$

869 we get that

$$h(\theta) = -\alpha_\infty \nabla_\theta \text{KL}(\pi || \lambda_\theta) - (1 - \alpha_\infty) \nabla_\theta \text{KL}(\lambda_\theta || \pi) .$$

870 Recall that  $\Theta = \mathbb{R}^q$ . To check **SA1**, we set

$$w(\theta) = \alpha_\infty \text{KL}(\pi || \lambda_\theta) - (1 - \alpha_\infty) \text{KL}(\lambda_\theta || \pi) , \text{ for } \theta \in \Theta.$$

871 and for  $i \in \mathbb{N}$ ,  $\zeta_i = \log(i + 1)$ . The subset  $\mathcal{R}_i$  is a ball centered at 0 and of radius  $r_i$  where  $r_i$  is  
872 chosen so that  $\sup_{\|\theta\| \leq r_i} \nabla w(\theta) \leq c\zeta_i$  (such  $r_i$  exists using **A3**). It is easily checked that **SA1** is  
873 satisfied thanks to **A3** (note in particular that  $\nabla w$  is globally Lipschitz under the stated conditions).  
874 Conditions **SA3-(v)-(vi)-(vii)** are automatically satisfied.

875 We consider the drift function for the Markov kernel  $P_{\theta,N}$

$$V(y[1], \dots, y[M]) = \sum_{i=1}^M V(y[i])$$

876 where  $V$  is the drift function in **A1**. **MC1** follows from Theorem 2 under **A4**. It is important to note  
877 that it is essential to have explicit controls on the drift and reduction conditions here. Conditions **MC**  
878 2 and **MC2** follow from **A3**. The precise tuning of constants is done along the same lines as [10,  
879 Section 5.3].

## 880 E Numerical experiments

### 881 E.1 Metrics

882 **ESTV** To compute Empirical sliced total variation distance (ESTV), we perform 25 random one-  
883 dimensional projections and then perform Kernel Density Estimation there for reference and produced  
884 samples. We then take the TV-distance between two distributions over 1D grids of 1000 points.  
885 We consider the value averaged over the projections to show the divergence between the MCMC  
886 distribution and the reference distribution.

887 **EMD** We compute the EMD as the transport cost between sample and reference points in  $L_2$  using  
888 the algorithm proposed in [14]. Then we report the EMD rescaled by the target dimension  $d$ .

889 **ESS** ESS (effective sample size) measures how many independent samples from target yield  
890 (approximately) the same variance for estimating the mean of some function. The closer ESS is  
891 to 1, the better is the sampler. Following [26], we compute ESS component-wise for multivariate  
892 distributions. Namely, given a sample  $\{Y_t\}_{t=1}^M$ ,  $Y_t \in \mathbb{R}^d$  of size  $M$ , for  $i = 1, \dots, d$ , we compute

$$\text{ESS}_i = \frac{1}{1 + \sum_{k=1}^M \rho_k^{(i)}} .$$

893 Here  $\rho_k^{(i)} = \frac{\text{Cov}(Y_{t,i}, Y_{t+k,i})}{\text{Var}(Y_{t,i})}$  is the autocorrelation at lag  $k$  for  $i$ -th component. We replace  $\rho_k^{(i)}$  by its  
894 sample counterpart  $\hat{\rho}_k^{(i)}$ , and report  $\text{ESS} = d^{-1} \sum_{i=1}^d \widehat{\text{ESS}}_i$ , where

$$\widehat{\text{ESS}}_i = \frac{1}{1 + \sum_{k=1}^M \hat{\rho}_k^{(i)}} .$$

### 895 E.2 Mixtures of Gaussians

896 **Equally weighted Gaussians in two dimension** The target density is

$$p_\beta(x) \propto \sum_{i=1}^3 \beta_i \exp\left\{-\|x - \mu_i\|^2/(2\sigma^2)\right\} . \quad (25)$$

897 Here we choose  $\sigma = 1$ ,  $\beta_i = 1/3$ , and  $\mu_i$ ,  $i \in \{1, 2, 3\}$  as vertices of an equilateral triangle with  
898 side length  $4\sqrt{3}$  and center  $(0, 0)$ . The contour representation of (25) can be found in Figure 2. We  
899 compare 3 sampling strategies:

- 900 • i-SIR algorithm with  $N = 3$  particles and  $\mathcal{N}(0, 4I)$  proposal distribution;  
 901 • MALA with step size  $\gamma = 0.5$ , tuned to obtain acceptance rate 0.67;  
 902 • Ex<sup>2</sup>MCMC algorithm with the same parameters as i-SIR and 3 consecutive MALA steps with  
 903  $\gamma = 0.5$  as rejuvenations.

904 We generate 100 observations within each sampler and represent them in Figure 2. For the MALA  
 905 sampler, we generate 300 samples and select every 3th to maintain compatibility with the Ex<sup>2</sup>MCMC  
 906 setup. Note that in this example, the variance of the global proposals in i-SIR should be relatively  
 907 large to cover well all modes of the (25) mixture. However, since the modes are narrow, the step  
 908 size of MALA cannot be very large to obtain a sensible acceptance rate. Therefore, Figure 2 shows  
 909 the drawbacks of the two approaches: i-SIR covers all modes of the target, but the chain often gets  
 910 stuck at a certain point, which affects the variability of the samples. MALA allows a better local  
 911 exploration of each mode, but does not cover the whole support of the target. The Ex<sup>2</sup>MCMC  
 912 algorithm combines the advantages of both methods by combining i-SIR-based global exploration  
 913 with MALA -based local exploration.

914 Now, the mixture model of (25) is modified with the weights parameters  $\beta = (\beta_1, \beta_2, \beta_3) =$   
 915  $(2/3, 1/6, 1/6)$  and same values of  $\mu_i$  and  $\sigma$ . To compare the quality of the methods, we perform the  
 916 following procedure

- 917 • starting with the initial distribution  $\mathcal{N}(0, 4I)$ , we generate the trajectory  $(X_1, \dots, X_n)$  for different  
 918 values of  $n \in [25, 800]$  for each of the compared methods (i-SIR MALA, Ex<sup>2</sup>MCMC ). Sampler  
 919 hyperparameters are the same as above, and the burn-in period equals 50;
- 920 • We perform the kernel density estimate (KDE)  $\hat{p}_n$  based on the observations  $(X_1, \dots, X_n)$ ,  
 921 and compute the total variation distance between  $\hat{p}_n$  and the target density  $p_\beta$ , and the forward  
 922  $\text{KL}(\hat{p}_n || p_\beta)$ . Then we average the results over 100 independent runs of each sampler.

923 Now we use the same values for the means and covariances but set the mixing weights to  $\beta =$   
 924  $(\beta_1, \beta_2, \beta_3) = (2/3, 1/6, 1/6)$ . To compare the different sampling methods, we perform the following  
 925 procedure.

- 926 • starting from the initial distribution  $\mathcal{N}(0, 4I)$ , we generate the trajectory  $(X_1, \dots, X_n)$  for different  
 927 values of  $n \in [25, 800]$  for each of the compared methods (i-SIR MALA, Ex<sup>2</sup>MCMC ). The  
 928 hyperparameters of the sampler are the same as above, and the burn-in period is 50;
- 929 • We perform kernel density estimation (KDE)  $\hat{p}_n$  based on the observations  $(X_1, \dots, X_n)$  and  
 930 calculate the total variation distance between  $\hat{p}_n$  and the target density  $p_\beta$ , as well as the forward  
 931 value  $\text{KL}(\hat{p}_n || p_\beta)$ . We then average the results over 100 independent runs of each sampler.

932 The results for each sampler are given in Figure 3b. We also provide a simple illustration to  
 933 the statements of (2) and Theorem 2. Starting from the initial distribution  $\xi \sim \mathcal{N}(0, 4I)$ , we  
 934 draw 500 independent chains of length 50 for each of the compared methods. Using these 500  
 935 observations, we create a KDE  $\hat{p}_n$  for the density corresponding to the distribution of  $\xi Q^n$  for  
 936 different  $n \in \{5, \dots, 50\}$  and Q corresponding to i-SIR MALA or Ex<sup>2</sup>MCMC. Then we calculate  
 937 the total variation distance between  $\hat{p}_n$  and the target density  $p_\beta$ . Corresponding plots can be  
 938 found in Figure 3a. Note that Ex<sup>2</sup>MCMC significantly outperforms the results of both MALA and  
 939 i-SIR. Indeed, the inhomogeneous mixture model is a complicated target for the Langevin-based  
 940 methods. The trajectories generated by MALA tend to remain in a single mode of mixture (25),  
 941 which reduces the reliability of the estimates and requires the generation of long trajectories even for  
 942  $d = 2$ . At the same time, it is difficult for i-SIR type methods without local exploration trajectories  
 943 to quickly cover all the modes.

### 944 E.3 Normalizing flow RealNVP

945 We use the RealNVP architecture ([20]) for our experiments with adaptive MCMC. The key element  
 946 of RealNVP is a coupling layer, defined as a transformation  $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$ :

$$y_{1:d} = x_{1:d}$$

$$y_{d+1:D} = x_{d+1:D} \odot \exp(s(x_{1:d}) + t(x_{1:d}))$$

947 where  $s$  and  $t$  are some functions from  $\mathbb{R}^D$  to  $\mathbb{R}^D$ . Thus, it is clear that the Jacobian of such a  
 948 transformation is a triangular matrix with nonzero diagonal terms. We use fully connected neural  
 949 networks to parameterize the functions  $s$  and  $t$ .

950 In all experiments with normalizing flows, we use the optimizer Adam ([37]) with  $\beta_1 = 0.9$ ,  $\beta_2 =$   
 951 0.999 and weight decay 0.01 to avoid overfitting.

#### 952 E.4 Distributions with complex geometry

953 In this section, we study the sampling quality from high-dimensional distributions, whose density  
 954 levels have high curvature (Banana shaped and Funnel distributions, details below). With such  
 955 distributions, standard MCMC algorithms like MALA or i-SIR, fail to explore fully the density  
 956 support.

957 The corresponding densities are given for  $x \in \mathbb{R}^d$  by

$$\begin{aligned} p_f(x) &= Z^{-1} \exp \left( -x_1^2/2a^2 - (1/2)e^{-2bx_1} \sum_{i=2}^d \{x_i^2 + 2bx_1\} \right), \quad d \geq 2, \\ p_b(x) &= Z^{-1} \exp \left( - \sum_{i=1}^{d/2} \{x_{2i}^2/2a^2 - (x_{2i-1} - bx_{2i} + a^2b)^2/2\} \right), \quad d = 2k, k \in \mathbb{N}. \end{aligned} \quad (26)$$

958 where  $Z$  is a normalizing constant. We set  $a = 2$ ,  $b = 0.5$  for funnel and  $a = 5$ ,  $b = 0.02$  for banana-  
 959 shape distributions, respectively. For MALA we use an adaptive step size tuning strategy to maintain  
 960 acceptance rate approximately 0.5. For i-SIR and Ex<sup>2</sup>MCMC algorithms we use wide Gaussian  
 961 global proposal  $\mathcal{N}(0, \sigma_p^2 \mathbf{I})$  with  $\sigma_p^2 = 4$  for Funnel and  $\sigma_p^2 = 9$  for Banana-shape distribution.

962 For FlEx<sup>2</sup>MCMC use a simple RealNVP-based normalizing flow [20] with 4 hidden layers. Note  
 963 that for  $p_f(x)$  the energy landscape in the region with  $x_1 < 0$  is steep, so the distributions (26) are  
 964 hard to capture, especially when the dimension  $d$  is large. Moreover, due to the complex geometry  
 965 of the distribution support, we cannot hope that local samplers (MALA) or global samplers (i-SIR )  
 966 alone will give good results. In this example, we want to compare FlEx<sup>2</sup>MCMC with i-SIR MALA  
 967 and the HMC-based NUTS sampler [34]. We also add a vanilla version of the Ex<sup>2</sup>MCMC algorithm  
 968 to the comparison. To generate the ground-truth samples, we use the explicit reparametrisation of (26).  
 969 Indeed, given a random vector  $(Z_1, \dots, Z_d) \sim \mathcal{N}(0, \mathbf{I})$ , we consider its transformation  $(X_1, \dots, X_d)$   
 970 under the formulas

$$\begin{cases} X_1 = aZ_1 \\ X_i = e^{bX_1} Z_i, \quad i \in \{2, \dots, d\} \end{cases}.$$

971 It is easy to check that  $(X_1, \dots, X_d)$  follows the density  $p_f(x)$ ,  $x \in \mathbb{R}^d$ . Similarly, for  $d = 2k$   
 972 consider the transformation

$$\begin{cases} Y_{2i} = aZ_{2i} \\ Y_{2i-1} = Y_{2i} + bY_{2i}^2 - ba^2, \quad i \in \{1, \dots, k\} \end{cases}.$$

973 Then  $(Y_1, \dots, Y_d)$  follows the density  $p_b(x)$ ,  $x \in \mathbb{R}^d$ . We provide the average computation time  
 974 for NUTS, adaptive i-SIR and FlEx<sup>2</sup>MCMC algorithms in Table 1 and Table 2 for the Funnel  
 975 and Banana-shape distributions, respectively, averaged over 50 runs. Note that different runs of  
 976 NUTS algorithm yields high variance of the running time, especially for the Funnel distribution and  
 977 dimensions  $d \geq 50$ .

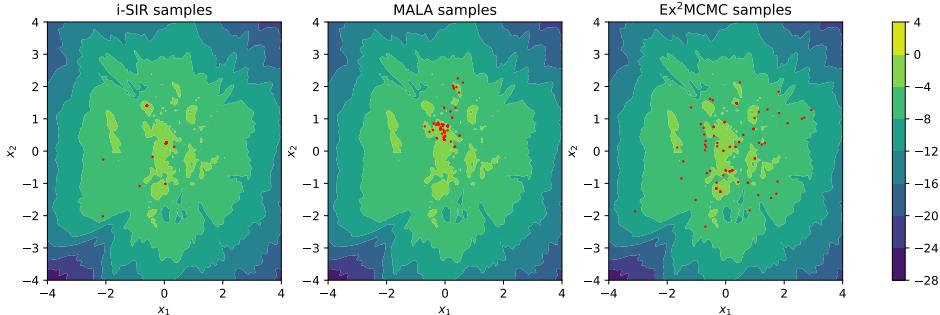
978 We give the computation time for the above algorithms and additional implementation details in  
 979 Appendix E.4. The implementation of FlEx<sup>2</sup>MCMC is based on the use of 5 MALA steps as  
 980 rejuvenation steps.

Method	$d = 10$	$d = 20$	$d = 50$	$d = 100$	$d = 200$
NUTS	$33.4 \pm 8.2$	$41.1 \pm 12.3$	$61.6 \pm 30.2$	$82.3 \pm 73.2$	$88.4 \pm 59.5$
Adaptive i-SIR	$38.1 \pm 3.2$	$39.4 \pm 2.8$	$45.3 \pm 2.5$	$59.8 \pm 0.7$	$80.4 \pm 0.4$
FlEx <sup>2</sup> MCMC	$46.8 \pm 3.2$	$48.2 \pm 2.8$	$54.2 \pm 2.5$	$68.8 \pm 0.8$	$89.5 \pm 0.5$

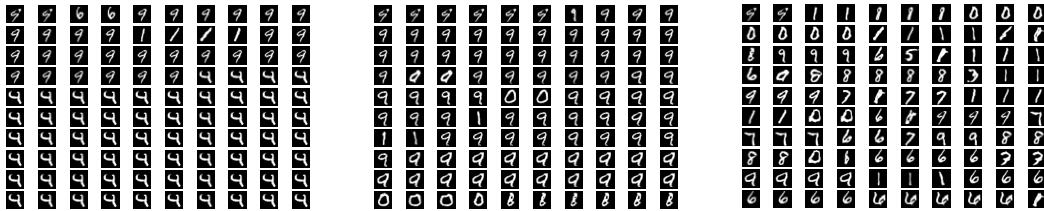
Table 1: Computational time for the Funnel distribution.

Method	$d = 20$	$d = 40$	$d = 60$	$d = 80$	$d = 100$
NUTS	$27.6 \pm 1.8$	$32.1 \pm 1$	$34.2 \pm 0.5$	$35.2 \pm 0.5$	$35.9 \pm 0.4$
Adaptive i-SIR	$24.5 \pm 0.2$	$26.8 \pm 0.3$	$28.5 \pm 0.2$	$30.1 \pm 0.2$	$32.8 \pm 0.2$
FlEx <sup>2</sup> MCMC	$39.3 \pm 0.5$	$41.8 \pm 0.3$	$43.5 \pm 0.3$	$45.1 \pm 0.3$	$47.8 \pm 0.4$

Table 2: Computational time for the Banana-shape distribution.



(a) JS-GAN: latent space visualizations



(b) i-SIR samples

(c) MALA samples

(d) Ex<sup>2</sup>MCMC samples

## 981 E.5 GANs as energy-based models

### 982 E.5.1 MNIST results

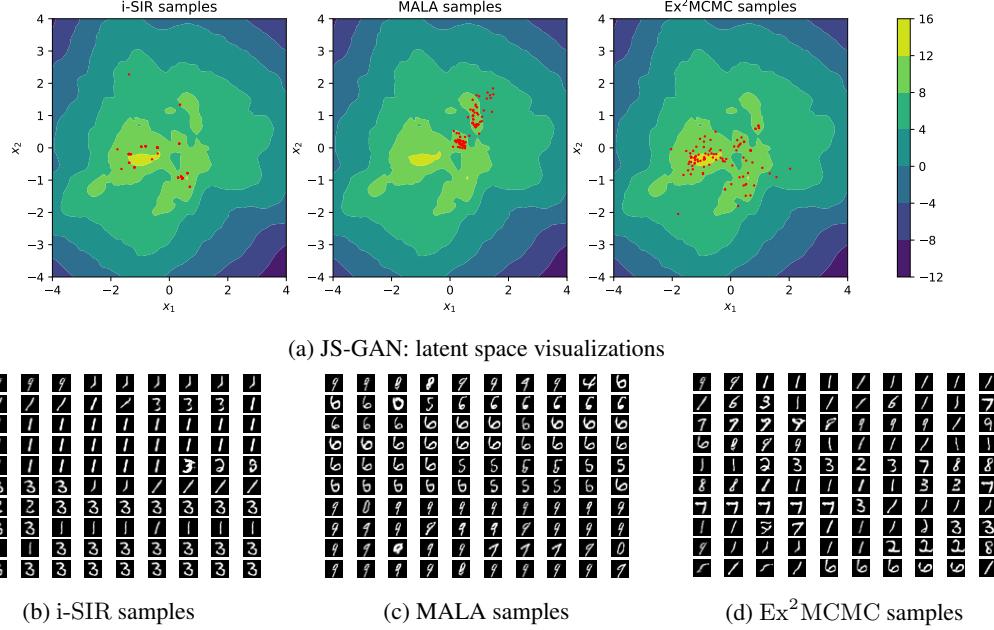
983 For this example, we consider both the Wasserstein GAN (WGAN) setup with energy function  
984  $E_W(z)$  and the classical Jensen-Shannon GAN with energy function  $E_{JS}(z)$ . In both cases, we use  
985 fully connected networks with 3 convolutional layers for discriminator and 3 linear + 3 convolutional  
986 layers for generator. For WGAN training, we use gradient penalty regularisation, following [31].  
987 We provide additional visualisations of the latent space and samples along a given trajectory for  
988 Jensen-Shannon GAN in Appendix E.5.1 and for Wasserstein GAN in Appendix E.5.1. Sampling  
989 hyperparameters are summarized in Table 3. For fair comparison, we take each 3-rd sample produced  
990 by the MALA, when running this algorithm separately. Both for WGAN-GP and vanilla GAN  
991 experiments we apply i-SIR and Ex<sup>2</sup>MCMC with wide Gaussian global proposal  $\mathcal{N}(0, \sigma_p^2)$ . The  
992 particular values of  $\sigma_p^2$  are specified in Table 3.

### 993 E.5.2 Cifar-10 results

994 We consider two popular GAN architectures, DC-GAN [58] and SN-GAN [48]. Below we provide  
995 the details on experimental setup and evaluation for both of the models.

Method	# iterations	MALA step size $\gamma$	# particles, $N$	$\sigma_p^2$	# MALA steps
JS-GAN	100	0.02	10	9	3
WGAN-GP	100	0.02	10	9	3

Table 3: MNIST hyperparameters.



## 996 E.6 Training and sampling details.

997 For DC-GAN and SN-GAN experiments, we took the implementation and training script of the  
 998 models from Mimicry repository <https://github.com/kwotsin/mimicry>. Both models were  
 999 trained on a single GPU GeForce GTX 1060 for approximately 20 hours.

1000 Both for DC-GAN and SN-GAN, the latent dimension is equal to  $d = 128$ . Following [17], for both  
 1001 models we consider sampling from the latent spatial distribution

$$p(z) = e^{-E_{JS}(z)} / Z, \quad z \in \mathbb{R}^d, \quad E_{JS}(z) = -\log p_0(z) - \text{logit}(D(G(z))),$$

1002 where  $\text{logit}(y) = \log(y/(1-y))$   $y \in (0, 1)$  is the inverse of the sigmoid function and  $p_0(z) =$   
 1003  $\mathcal{N}(0, I)$ .

1004 **Evaluation protocol** We perform  $n = 100$  iterations of the algorithms MALA, i-SIR, Ex<sup>2</sup>MCMC  
 1005 and FIE<sup>2</sup>MCMC. For both the vanilla Ex<sup>2</sup>MCMC algorithm (Algorithm 2) and FIE<sup>2</sup>MCMC we  
 1006 use the Markov kernel (17), which corresponds to 3 MALA steps, as the rejuvenation kernel. The  
 1007 step size  $\gamma$  given for the algorithm Ex<sup>2</sup>MCMC corresponds to its rejuvenation kernel MALA. For  
 1008 more experimental details, see Table 4. For i-SIR and Ex<sup>2</sup>MCMC algorithms we use  $\mathcal{N}(0, \sigma_p^2 I)$   
 1009 with  $\sigma_p^2 = 1$  as a global proposal distribution.

1010 We run  $M = 500$  independent chains for each of the above MCMC algorithms. Then, for the  
 1011  $j$ -th iteration, we compute the average value of the energy function  $E(z)$  averaged over  $M$  chains.  
 1012 Hyperparameters are specified in Table 4. Energy profiles for different algorithms for DC-GAN and  
 1013 SN-GAN are provided in Figure 14 and Figure 11, respectively. Note that in both cases Ex<sup>2</sup>MCMC or  
 1014 FIE<sup>2</sup>MCMC algorithms yields lower energy samples. We visualize 10 randomly chosen trajectories  
 1015 obtained with each sampling methods in Figure 12-Figure 13 for SN-GAN and Figure 15-Figure 16  
 1016 for DC-GAN, respectively. For each trajectory we visualize every 10-th sample. Both architectures  
 1017 indicate the same findings: MALA typically is not available to escape the mode of the corresponding  
 1018 target density  $p(z)$  during one particular run. i-SIR travels well across the support of  $p(z)$ , yet the  
 1019 corresponding energy values are higher then the ones of Ex<sup>2</sup>MCMC or FIE<sup>2</sup>MCMC. Some i-SIR  
 1020 trajectories can get trapped in one particular image due to the absence of local exploration moves. At  
 1021 the same time, Ex<sup>2</sup>MCMC as illustrated in Figure 13-13a and Figure 16-16a, can both exploit the  
 1022 particular mode of the distribution and perform global moves over the support of  $p(z)$ . Of course,  
 1023 these global moves are more likely to occur during the first sampling iterations.

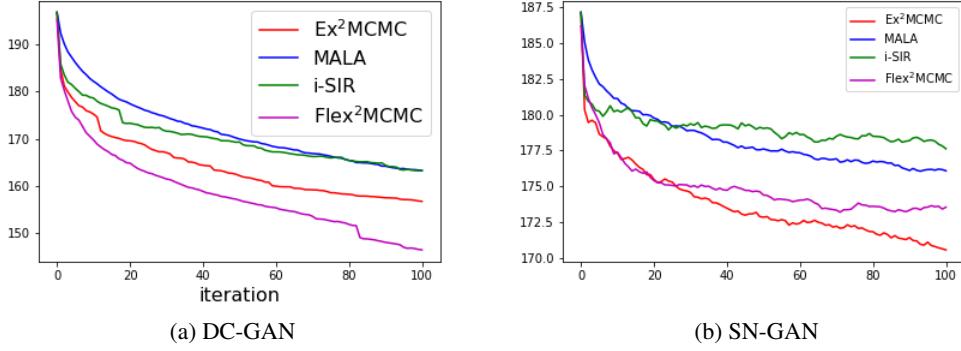


Figure 10: Energy profile for DC-GAN and SN-GAN architectures on CIFAR-10 dataset.

GAN type	# iterations	MALA step size $\gamma$	# particles, $N$	$\sigma_p^2$	# MALA steps
SNGAN	100	$5 \times 10^{-3}$	10	1	3
DCGAN	100	$10^{-3}$	10	1	3

Table 4: CIFAR-10 hyperparameters.

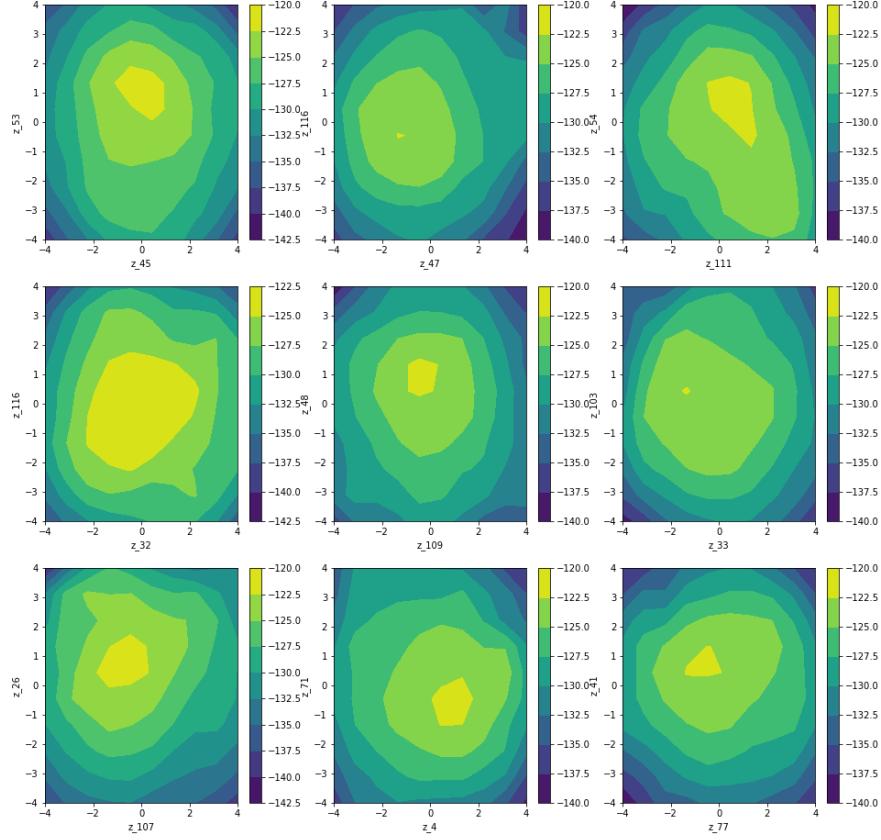
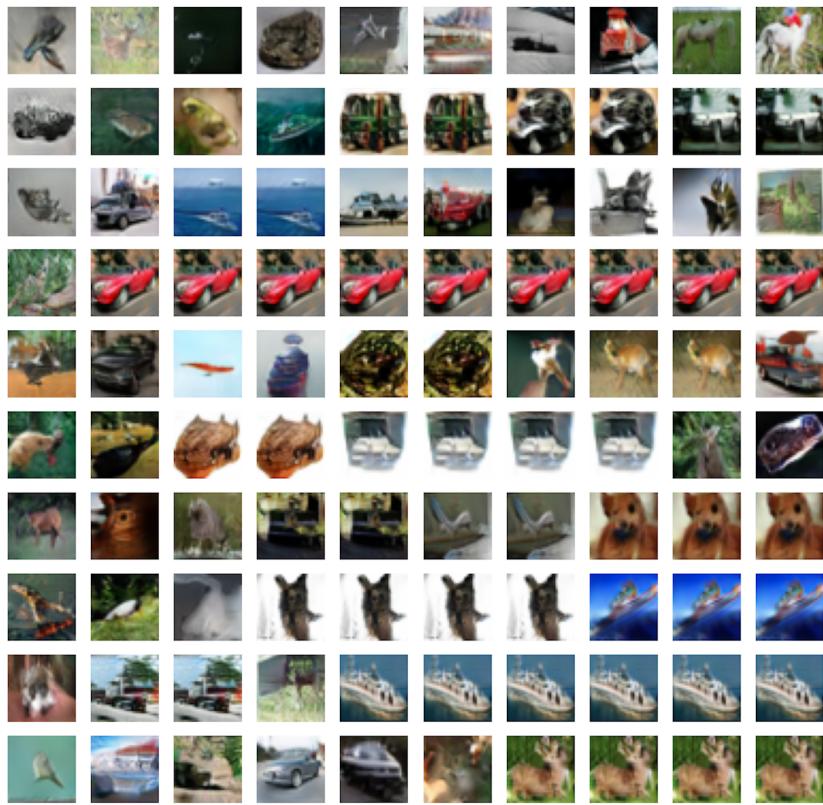
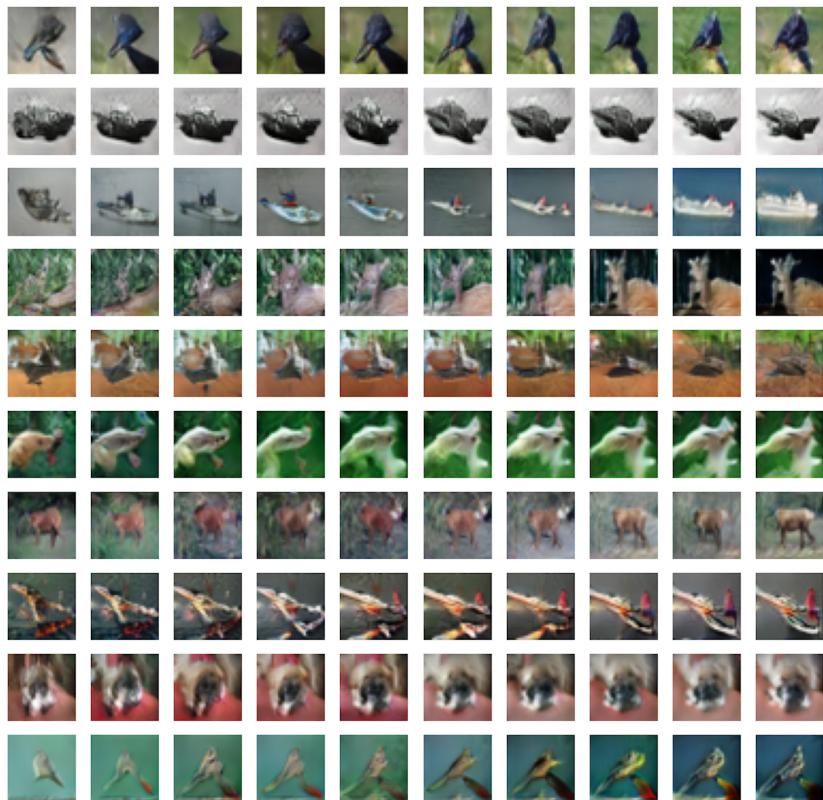


Figure 11: Energy profile for random axis pairs, SN-GAN

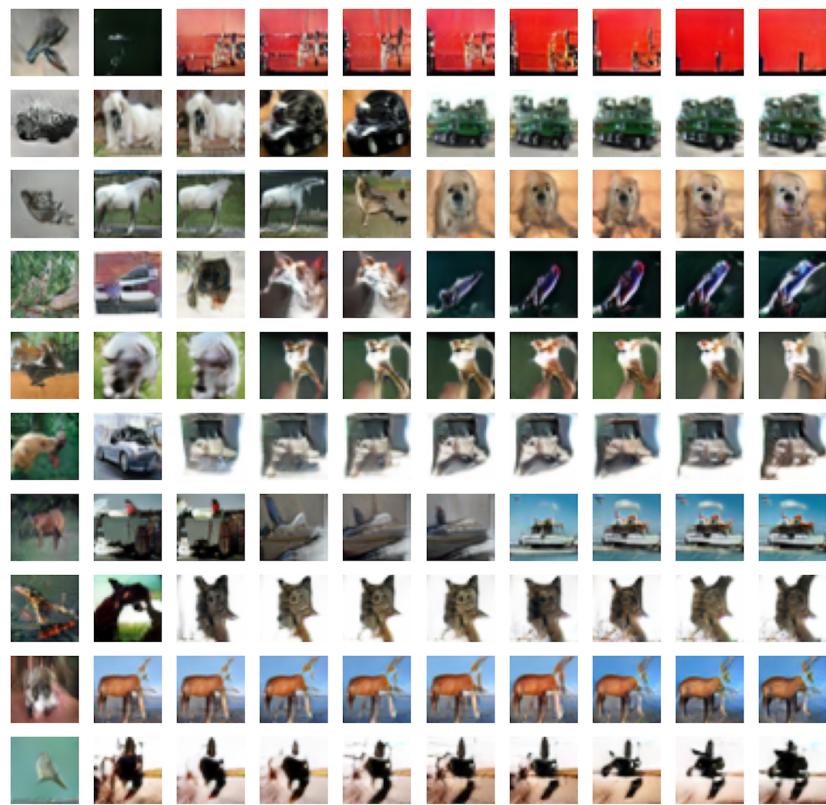


(a) i-SIR samples

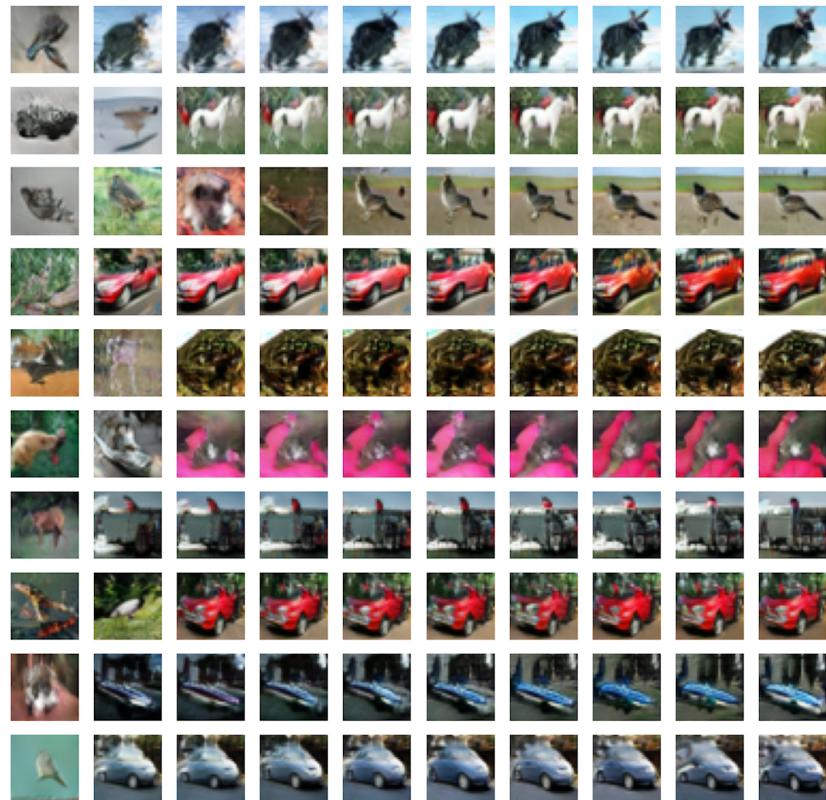


(b) MALA samples

Figure 12: i-SIR and MALA samples, SN-GAN.



(a) Ex<sup>2</sup>MCMC samples



(b) FlEx<sup>2</sup>MCMC samples

Figure 13: Ex<sup>2</sup>MCMC and FlEx<sup>2</sup>MCMC samples, SN-GAN.

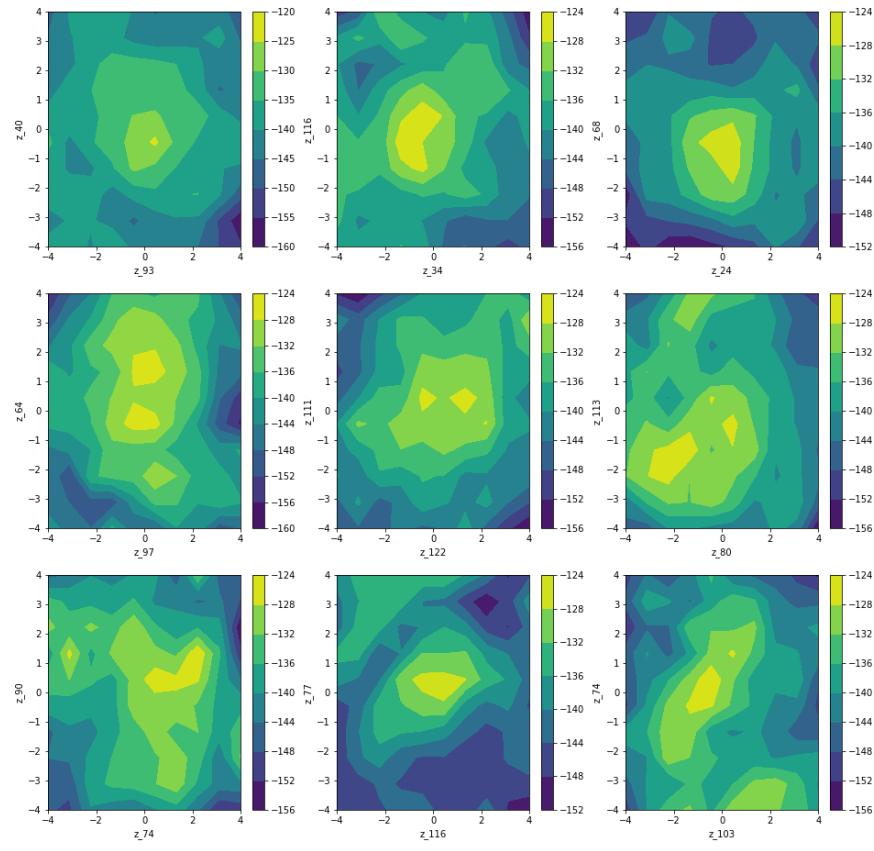


Figure 14: Energy profile for random axis pairs, DC-GAN

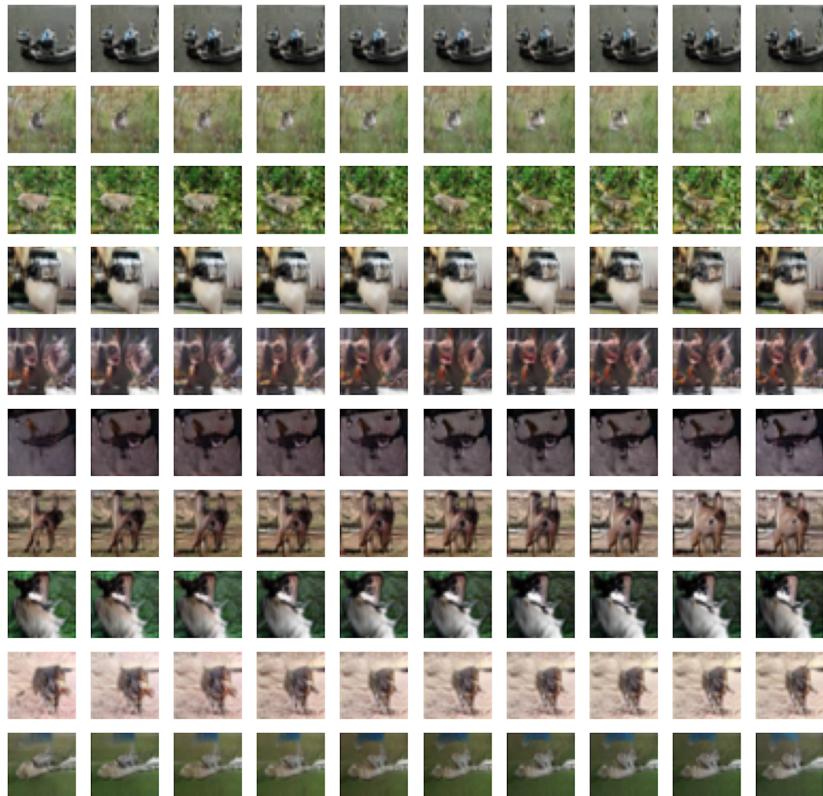
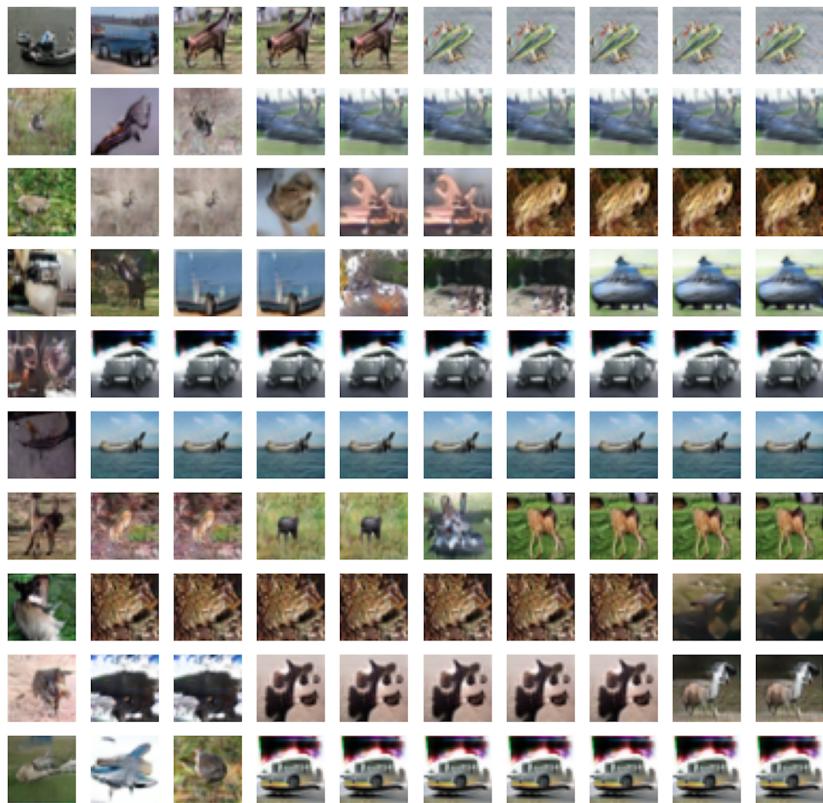
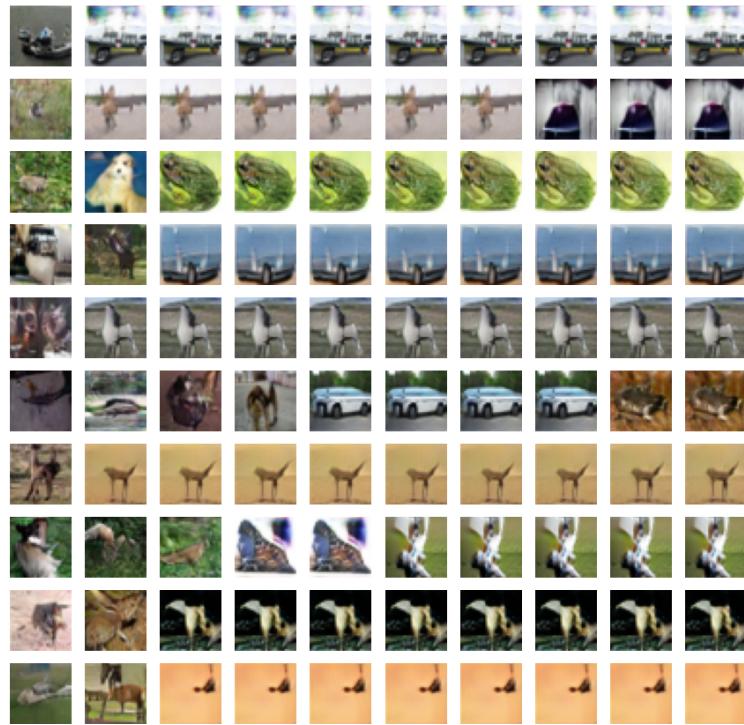


Figure 15: i-SIR and MALA samples, DC-GAN.



(a) Ex<sup>2</sup>MCMC samples



(b) FlEx<sup>2</sup>MCMC samples

Figure 16: Ex<sup>2</sup>MCMC and FlEx<sup>2</sup>MCMC samples, DC-GAN.