

1 We are grateful to all the reviewers for their feedback. Below we provide responses to the main comments.

2 **Reviewer 1:**

3 "I don't feel NeurIPS is the appropriate venue to publish it. My main concern about the paper is with respect to its
4 appeal for the ML community at large. My impression is that its scope is rather limited. The presented examples
5 are of reduced dimension" Although we respect reviewer's opinion we disagree and we believe that our work is very
6 suitable for NeurIPS which is an interdisciplinary conference and MCMC is one of its subject areas. Many MCMC
7 papers and related Monte Carlo methods have been previously published in NeurIPS. Also not all inference problems in
8 statistics and ML are large scale or high-dimensional and certainly our method does not exclude applicability to high
9 dimensions.

10 **Reviewer 2:**

11 "It was not clear how the Adaptive MCMC (AM) described in the results works i.e. what is the objective function of
12 the adaptation. A reference is made to the supplement, but I didn't quite find it there (lines 226-227)." The underlying
13 objective of AM is the minimisation of the KL divergence between the target distribution and the proposal distribution,
14 as described earlier in section 4.1 in the tutorial paper of Andrieu and Thoms. Of course this optimisation is challenging
15 because we only observe correlated samples which at the early adaptation stages are really far from the target.

16 "Line 266 and Figure 2 top panel. This should really be the auto-correlation plot." Thank you, we will follow your
17 suggestion.

18 "The choice of ρ_t on lines 136, 137 is not well motivated. In similar situations a Robbins-Monro sequence is typically
19 used." We agree that the Robbins-Monro sequence is the one that ensures convergence in the limit. The motivation
20 behind the RMSprop sequence we used (that is very popular in Deep Learning together with similar adaptive learning
21 rate schemes such as Adam) is that in practise when you run for a fixed (relatively small) budget of stochastic optimisation/
22 adaptation iterations it tends to provide more effective optimisation. Note that in our experiments we adapt only
23 during burn-in, while at the collection of samples phase we keep the proposal fixed.

24 "The authors could give more theoretical justification for their MALA approximation to avoid computing the Hessian
25 on lines 172-174." So far we only empirically observe that the fast MALA scheme tends to provide stochastic gradients
26 with smaller variance leading to faster optimisation. We will try to analyse this theoretically by finding expressions of
27 the variance (at least for simple targets) for the exact Hessian and the fast scheme.

28 "I would like to see the results for Stan, which has a good implementation of NUTS, on the presented data sets.
29 That would help highlight the advantage of the proposed method over the current standard practice." Currently all
30 experiments are based on a MATLAB implementation and the NUTS version is precisely from the published 2014
31 article and follows Hoffman's implementation. We are going to provide code that reproduces all our results.

32 **Reviewer 3:**

33 "The specific algorithm (and simplifications) for MALA are clearly the "highlight", performance-wise, but may be
34 slightly lower impact simply due to implementation headaches (it can be done easily enough by people familiar with
35 deep learning or AD software, but that leaves out many practicing statisticians)." Thanks for the comment. For the
36 final version we plan to release a non automatic differentiation (AD) based MATLAB implementation for all proposed
37 algorithms. We also plan to provide pseudo-code showing how this fast MALA scheme can be implemented with the
38 minimum number of vector operations (few details about this are already part of the supplement) and with AD possibly
39 used only to compute the gradient of the log target.

40 "the one thing I find unsatisfying is the need to fall back on existing results for "optimal" acceptance rates when tuning
41 β . For RW and MALA proposals, maybe this is even reasonable (there's some theory about optimal preconditioners,
42 proposal covariances, mass matrices anyway), but it would be nice to address this more generally, or at least discuss
43 the appropriateness of falling back on these recommendations." We agree that this is a limitation. The ideal will be to
44 automatically "learn" what is the optimal average acceptance rate for a specific target. However, the standard heuristics
45 for adapting β worked well in all our experiments.

46 "For MALA and NUTS, what sort of preconditioning is done as a baseline? For NUTS as implemented in STAN
47 and PyMC3, I know a standard practice is to run a short pre-run which is used to estimate the mass matrix; this
48 is then used fixed with only the step size adapted online. Is something like this done here? What about for the
49 non-adaptive MALA?" For NUTS we use the algorithm as defined in the initial 2014 paper and it is based on the
50 corresponding implementation of Hoffman, which does not use preconditioners. Also the non-adaptive MALA is not
51 using a preconditioner. Notice our method regarding MALA essentially allows to obtain a full covariance preconditioner
52 using gradient-based adaptation.

53 "line 140, L is described as a "positive definite lower triangular matrix" ... presumably a typo (it is correct later for
54 MALA), but LL' is positive definite, while L is lower triangular." Thank you, we will clarify this.