# Analysis of high-dimensional Continuous Time Markov Chains using the Local Bouncy Particle Sampler

Tingting Zhao

TINGTING.ZHAO@STAT.UBC.CA

Department of Statistics University of British Columbia Vancouver, BC, V6T 1N4, Canada

Alexandre Bouchard-Côté

BOUCHARD@STAT.UBC.CA

Department of Statistics University of British Columbia Vancouver, BC, V6T 1N4, Canada

**Editor:** XXX

#### Abstract

Sampling the parameters of high-dimensional Continuous Time Markov Chains (CTMC) is a challenging problem with important applications in many fields of applied statistics. In this work a recently proposed type of non-reversible rejection-free Markov Chain Monte Carlo (MCMC) sampler, the Bouncy Particle Sampler (BPS), is brought to bear to this problem. BPS has demonstrated its favourable computational efficiency compared with state-of-the-art MCMC algorithms, however to date applications to real-data scenario were scarce. An important aspect of practical implementation of BPS is the simulation of event times. Default implementations use conservative thinning bounds. Such bounds can slow down the algorithm and limit the computational performance. Our paper develops an algorithm with exact analytical solution to the random event times in the context of CTMCs. Our local version of BPS algorithm takes advantage of the sparse structure in the target factor graph and we also provide a framework for assessing the computational complexity of local BPS algorithms.

**Keywords:** CTMCs, Hamiltonian Monte Carlo (HMC), Piecewise-deterministic Markov Process (PDMPs), BPS, Local Bouncy Particle Sampler (LBPS), Generalized Linear Models (GLM).

#### 1. Introduction

CTMCs have widespread applications ranging from multi-state disease progression to phylogenetics (Yin and Zhang, 2012). However, the estimation of parameters in CTMCs is a challenging problem when incomplete data observations are only available at a finite number of time points. This is the case in a wide range of applications, for analyzing censored survival data (Kay, 1977), for describing panel data under Markov assumptions (Kalbfleisch and Lawless, 1985), for characterizing multi-state disease progression (Jackson et al., 2003), and for inferring evolutionary processes (Jukes and Cantor, 1969; Zhao et al., 2016) using biological sequences.

An (homogeneous) CTMC is a continuous-time stochastic process taking values on a finite or countable set. The parameters involved in a CTMC are used to characterize the

transitions between states and the distributions of the intervals between two consecutive transitions. The parameters are organized into a rate matrix. If the sample paths have been completely observed continuously over a finite time interval, statistical inference is straightforward. However, a more typical situation is that only partial observations of the states on a finite number of time points are available. For high dimensional rate matrices, efficient posterior inference is challenging. In particular, despite several algorithmic advances (Moler and Van Loan, 2003), matrix exponentiation, which is required to compute the marginal distributions of CTMCs, is still computational expensive.

In the following, we will make use of a flexible framework to parameterize rate matrices (Zhao et al., 2016), which subsumes much of the earlier parameterizations (Kimura, 1980; Hasegawa et al., 1985). This previous work used off-the-shelf Adaptive Hamiltonian Monte Carlo (AHMC) methods and did not exploit the sparsity often found in the parameterization of high-dimensional rate matrices.

In order to exploit sparsity, we make use of recently developed Monte Carlo schemes based on non-reversible PDMPs. In particular, we build our samplers based on the non-reversible rejection-free dynamics proposed in Peters et al. (2012) in the physics literature and later developed for statistical applications in Bouchard-Côté et al. (2018). It has been shown by Neal (2004); Sun et al. (2010); Chen and Hwang (2013); Bierkens (2016) that non-reversible MCMC algorithm can outperform reversible MCMC in terms of mixing rate and asymptotic performance. Related sampling schemes have been developed including continuous-time Monte Carlo algorithms and continuous-time Sequential Monte Carlo algorithms (Fearnhead et al., 2016; Pakman et al., 2017), but we focus on proof-of-concept applications of the BPS algorithm in this work.

In the BPS algorithm, the posterior samples of a variable of interest are continuously indexed by the position of a particle that moves along piecewise linear trajectories. When encountering a high energy barrier (low posterior density value), the particle is never rejected but instead the direction of its path is changed after a Newtonian collision. A key algorithmic requirement is to efficiently determine the bouncing time of the particle. Most existing work (Fearnhead et al., 2016; Vanetti et al., 2017; Pakman et al., 2017) use conservative bounds from thinning algorithm of an inhomogeneous Poisson Process (PP) to sample the collision time. Conservative bounds can lead to computational inefficiency of the algorithm. Bouchard-Côté et al. (2018) has obtained analytical solutions to the collision times under certain simple scenarios such as Gaussian distribution. Our paper provides an exact analytical solution for the important special case of CTMC posterior inference.

Moreover, if the target posterior distribution is composed of a product of factors, where each factor only depends on a small subset of the variables of sampling interest, BPS can be further sped up using a "local," sparse version of the algorithm. Whereas Bouchard-Côté et al. (2018) only provided real-data benchmarks for global BPS, we provide real data benchmarks for the local BPS algorithm.

To summarize, we make the following contributions:

- A novel sampler combining HMC and LBPS to efficiently sample from CTMCs while maintaining the techniques of auxiliary variables used in Zhao et al. (2016) under the Bayesian GLM parameterization of rate matrices for CTMCs.
- A characterization of a class of sparse factor graphs when LBPS is efficient.

- A running time analysis for one iteration of LBPS for both general factor graphs and factor graphs with sparse structure, and a running time analysis for the special case of our proposed algorithm LBPS-HMC in one iteration under a Bayesian GLM model with sparsity structure in the factor graph.
- A proof-of-concept application on protein evolution to demonstrate on real data the computational efficiency of LBPS compared with classical state-of-the-art HMC algorithms.
- Analytical solutions to the bouncing time for each factor of the factorized posterior density of CTMCs to boost the computational efficiency. In comparison, previous work use thinning algorithms to either obtain analytical conservative bounds (Bouchard-Côté et al., 2018) or construct approximate bounds via predictive models using local regression (Pakman et al., 2017).

The remaining of the paper is organized as follows. In Section 2, we introduce the notation and basics of CTMCs, review the Bayesian GLM rate matrix parameterization, provide the background of BPS and LBPS, and the target posterior density of the problem. In Section 3, based on Bayesian GLM rate matrix parameterization, we propose the Bayesian GLM chain General Time Reversible model (GTR), which characterizes the exchangeable rate between a pair of states based on the current pair and also its nearest neighbour pair defined according to physiochemical property similarities. This is also the model we use throughout our paper. In Section 4, we propose the notion of strong sparsity for a family of factor graphs, where LBPS can be efficient by taking advantage of the sparsity structure we defined. We also provide a running time analysis for one iteration of LBPS for general factor graphs and factor graphs with sparse structure according to our definition. In Section 5, the sampling scheme combining HMC and LBPS is proposed for CTMCs under our Bayesian GLM chain GTR model. For notation simplicity, we denote the combined sampling scheme as LBPS-HMC throughout this paper. We have also provided the running time analysis for one iteration of LBPS-HMC under a Bayesian GLM chain GTR model. Later, we compare our novel sampler to state-of-the-art methods using only HMC kernel via numerical experiments in Section 6. We use synthetic datasets to demonstrate the better scalability of our novel algorithm as the dimension of the parameters increases. We then provide a proof-of-concept real data application on protein evolution in Section 7. Finally, we discuss the advantages and potential extensions of our work in Section 8.

### 2. Problem setup and notations

#### 2.1 CTMCs notation

We first introduce some notation for CTMCs. More background on CTMCs can be found in Norris (1998); Guttorp and Minin (2018). We use the same notation as Zhao et al. (2016). An (homogeneous) CTMC is a continuous-time stochastic process  $\{X(t): t \geq 0\}$  taking values on a finite or countable set  $\mathcal{X}$ . Throughout the paper, we assume  $\mathcal{X}$  is finite and  $\mathcal{X} = \{1, 2, \ldots, |\mathcal{X}|\}$ . Denote  $\{X_n, n \geq 0\}$  as the sequence of states visited in the continuous time path  $\{X(t)\}$ , and let  $A_n$  be the corresponding times when the state changes. A rate matrix Q indexed by  $\mathcal{X}$  is used to describe the instantaneous transition rate for each pair

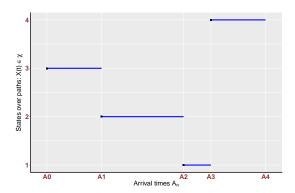


Figure 1: A fully observed realization of a single path of CTMCs at arrival time points  $A_1, A_2, A_3, A_4$  over a state space  $\{1, 2, 3, 4\}$ .

of distinctive states in  $\mathcal{X}$ . For example,  $q_{x,x'}$  represents the instantaneous rate between x and x', where  $x \in \mathcal{X}, x' \in \mathcal{X}, x \neq x'$ . The diagonal elements of Q are negative and enforce the constraint that each row sums to zero. The absolute values of the diagonal elements represent the rate parameter of an exponential distribution, used to characterize the waiting time spent on each state. We denote the initial state distribution as  $p_{\text{ini}}(\cdot)$ . Given the rate matrix Q, the transition probability matrix is  $P_Q(\Delta) = \exp(\Delta Q) = \sum_{j=0}^{\infty} (\Delta Q^j)/j!$ , for  $\forall \Delta \geq 0$ . We use  $\pi = (\pi_1, \pi_2, \dots, \pi_{|\mathcal{X}|})$  to represent the stationary distribution. For a stationary chain,  $\pi_i$  represents the proportion of time that the chain stays at state i in equilibrium. We denote  $\mathbf{z}(\mathbf{x}) \coloneqq (\mathbf{n}(\mathbf{x}), \mathbf{h}(\mathbf{x}), \mathbf{c}(\mathbf{x}))$  as the sufficient statistics of N fully observed CTMC paths (example shown in Figure 1), where  $n_x$  represents the number of paths started at state x,  $c_{x,x'}$  denotes the number of jumps from state x to x' and  $h_x$  indicates the total time spent in state x. The density of N fully observed paths given Q over time interval length  $\Delta$  is

$$f_{\mathbf{x}|\mathbf{Q},\Delta}(\boldsymbol{x}|\mathbf{Q},\Delta) := \left(\prod_{x\in\mathcal{X}} p_{\mathrm{ini}}(x)^{n_x}\right) \left(\prod_{(x,x')\in\mathcal{X}^{\mathrm{distinct}}} q_{x,x'}^{c_{x,x'}}\right) \times \left(\prod_{x\in\mathcal{X}} \exp\left(h_x q_{x,x}\right)\right).$$
(1)

Usually, these paths are only partially observed at finite number of points  $\tau_0, \tau_1, \ldots, \tau_n$  with states  $y_0, y_1, \ldots$  and  $y_n$ . We use  $\mathcal{Y}$  to represent the observed path. Assuming there is no error in the observation of the states, the density over this single path is

$$g_{y|Q}(\mathcal{Y}|Q) = p_{\text{ini}}(y_{\tau_0}) \prod_{k=1}^{K} (\exp(Q\Delta_k))_{y_{\tau_{k-1}}, y_{\tau_k}},$$
 (2)

 $\Delta_k = \tau_k - \tau_{k-1}$  is the length of the time interval between two consecutive observations.

The question of interest is often to estimate the rate matrix given partially observed paths. Pointwise evaluation of Equation 2 can be done in  $\mathcal{O}(|\mathcal{X}|^3)$  using diagonalization method to evaluate the matrix exponential. We reduce the computational cost by introducing the substitution mapping auxiliary variables (instantiated via uniformization algorithms) as described in Zhao et al. (2016). With the augmented sufficient statistics, the density of the paths is given by Equation 1.

#### 2.2 Bayesian GLM parameterization

We describe here a sparse rate matrix parameterization built in the framework of Zhao et al. (2016). This is useful since we propose Bayesian GLM chain GTR model in Section 3 on the basis of Bayesian GLM reversible rate matrix parameterization (Zhao et al., 2016). We first review briefly the Bayesian GLM unnormalized reversible rate matrix parameterization:

$$\theta_{\{x,x'\}}(\boldsymbol{w}^b) := \exp\{\langle \boldsymbol{w}^b, \boldsymbol{\phi}(\{x,x'\})\rangle\},$$
 (3)

$$\pi_x(\boldsymbol{w}^u) := \exp\left\{ \langle \boldsymbol{w}^u, \boldsymbol{\psi}(x) \rangle - A(\boldsymbol{w}^u) \right\},\tag{4}$$

$$A(\mathbf{w}^u) := \log \sum_{x \in \mathcal{X}} \exp \{\langle \mathbf{w}^u, \mathbf{\psi}(x) \rangle\},$$
 (5)

$$q_{x,x'}^{(\text{rev})}(\boldsymbol{w}) := \theta_{\{x,x'\}} \left(\boldsymbol{w}^b\right) \pi_{x'} \left(\boldsymbol{w}^u\right),$$
 (6)

The parameters of interest are  $\boldsymbol{w}$ , where  $\boldsymbol{w} = \begin{pmatrix} \boldsymbol{w}^u \\ \boldsymbol{w}^b \end{pmatrix}$ . We introduce the set of unordered distinct pairs of states as  $\mathcal{X}^{\text{unordered,dist.}} := \{\{x,x'\} \in \mathcal{X}^2 : x \neq x'\}$ . In Equation 3, we introduce bivariate feature functions  $\boldsymbol{\phi} : \mathcal{X}^{\text{unordered,dist.}} \to \mathbb{R}^{p_2}$ . In Equation 4, we introduce univariate feature functions  $\boldsymbol{\psi} : \mathcal{X} \to \mathbb{R}^{p_1}$ . We have  $p_1 + p_2 = p$  and  $\boldsymbol{w} \in \mathbb{R}^p$ . The weights corresponding to the *univariate features*  $\boldsymbol{\psi}$  are denoted as univariate weights  $\boldsymbol{w}^u$  and the weights related to the *bivariate features*  $\boldsymbol{\phi}$  are denoted as bivariate weights  $\boldsymbol{w}^b$ .

The connection between Bayesian GLM model and GTR model is discussed in details by Zhao et al. (2016). Zhao et al. (2016) have proved that the Bayesian GLM model can be used to represent any rate matrices under GTR parameterization equivalently. Here, we review this construction since the notation will be needed to introduce our new sparse model in Section 3.

Recall that in a GTR model, Q is parameterized by the stationary distribution  $\pi = (\pi_1, \pi_2, \dots, \pi_{|\mathcal{X}|})$  and exchangeable parameters  $\theta_{x,x'}$ , where  $q_{x,x'} = \theta_{\{x,x'\}}\pi_{x'}$ . Under the reversibility assumption, we have  $\theta_{\{x,x'\}} = \theta_{\{x',x\}}$ . In total, there are  $p_2 = |\mathcal{X}|(|\mathcal{X}|-1)/2$  exchangeable parameters  $\theta_{\{x,x'\}}$  under reversibility. We define a map  $\eta$ :  $\mathcal{X}^{\text{unordered,dist.}} \to \{1,2,\dots,|\mathcal{X}|(|\mathcal{X}|-1)/2\}$  such that the ith element of  $bivariate\ features\ \phi^{\text{gtr}}:\mathcal{X}^{\text{unordered,dist.}} \to \mathbb{R}^{p_2}$  is defined as:

$$\boldsymbol{\phi}_i^{\text{gtr}}(\{x, x'\}) := \mathbf{1}(\eta(\{x, x'\}) = i), \tag{7}$$

where the *i*th feature  $\phi_i^{\text{gtr}}(\{x, x'\})$  is equal to one if and only if the pair of unordered states  $\{x, x'\}$  is mapped to *i*th exchangeable parameter via  $\eta(\{x, x'\})$ . To ensure the reversibility

of the rate matrix, we require that  $\eta(\{x, x'\}) = \eta(\{x', x\})$ . Under the definition of  $\phi^{\text{gtr}}$ , any rate matrices can be represented under Bayesian GLM rate matrix parameterization with feature function  $\phi^{\text{gtr}}$ .

#### 2.3 Bayesian inference

Under the Bayesian GLM rate matrix parameterization defined in Equation 3-6, the parameters of interest are  $\boldsymbol{w}$ . We place a prior on  $\boldsymbol{w}$  with density denoted by  $g_{\mathbf{w}}(\boldsymbol{w})$ . Following Zhao et al. (2016), we use a Normal distribution with mean zero and precision  $\kappa$  for  $g_{\mathbf{w}}(\boldsymbol{w})$ .

We assume the observations are partially observed states  $y_1, y_2, \ldots, y_k$  at a finite number of time points  $\tau_0, \tau_1, \ldots, \tau_k$  on each sample path. The density over one single path is given in Equation 2. Under the Bayesian GLM parameterization, we denote the density by  $g_{\mathbf{x}|\mathbf{w}}(\mathcal{X}|Q(\mathbf{w}))$ . Following Equation 2, we obtain:

$$g_{\mathbf{y}|\mathbf{w}}(\mathcal{Y}|Q(\boldsymbol{w})) = p_{\mathrm{ini}}(y_{\tau_0}) \prod_{k=1}^{K} \left(\exp(Q(\boldsymbol{w})\Delta_k)\right)_{y_{\tau_{k-1}}, y_{\tau_k}}, \tag{8}$$

Thus, the target posterior density is given by

$$g_{\text{wlv}}(\boldsymbol{w}|\boldsymbol{x}) \propto g_{\text{w}}(\boldsymbol{w})g_{\text{vlw}}(\mathcal{Y}|Q(\boldsymbol{w})) = \exp(-U(\boldsymbol{w})),$$
 (9)

where  $U(\boldsymbol{w})$  represents the negative log of the unnormalized posterior density function.

In Equation 9, the second term involves the computation of matrix exponential. Zhao et al. (2016) have shown that classical MCMC algorithms are inefficient to sample from Bayesian GLM models for CTMCs. However, state-of-the-art methods such as HMC require the gradient of the target posterior density. A single gradient evaluation of Equation 8 will take a running time of  $\Theta(|\chi|^5)$  (Zhao et al., 2016). To circumvent the computational difficulties, we follow the substitution mapping techniques introduced in Zhao et al. (2016) to augment the sufficient statistics of the partially observed sample path of CTMCs given the two end-point states at two consecutive observation time points and ensure that the sampling algorithm still converges to the desired equilibrium distribution. Given the partially observed states of the starting and ending states of a CTMC path over time interval  $\Delta$ , we can simulate a full path conditional on the end-points and time interval  $\Delta$  according to Equation 1 using cached uniformization techniques described in Zhao et al. (2016). Following Equation 1, the negative unnormalized posterior log-density for N full paths of CTMCs given w with augmented sufficient statistics is

$$U_{z}(\boldsymbol{w}) = \frac{1}{2}\kappa \|\boldsymbol{w}\|_{2}^{2} - \sum_{x \in \mathcal{X}} h_{x}q_{x,x}(\boldsymbol{w})$$

$$- \sum_{(x,x')\in\mathcal{X}^{\text{distinct}}} c_{x,x'} \log \left(q_{x,x'}(\boldsymbol{w})\right) - \sum_{x \in \mathcal{X}} n_{x} \log(\pi_{x}(\boldsymbol{w})),$$

$$= \frac{1}{2}\kappa \|\boldsymbol{w}\|_{2}^{2} + \sum_{(x,x')\in\mathcal{X}^{\text{distinct}}} h_{x}q_{x,x'}(\boldsymbol{w}),$$

$$- \sum_{(x,x')\in\mathcal{X}^{\text{distinct}}} c_{x,x'} \log \left(q_{x,x'}(\boldsymbol{w})\right) - \sum_{x \in \mathcal{X}} n_{x} \log(\pi_{x}(\boldsymbol{w})).$$

$$(10)$$

# 2.4 Background on BPS and LBPS

#### 2.4.1 Basics of BPS

The BPS belongs to the type of emerging continuous-time non-reversible MCMC sampling algorithms constructed from PDMPs (Davis, 1984). It was first proposed by Peters et al. (2012), formalized, developed, and generalized by Bouchard-Côté et al. (2018); Vanetti et al. (2017). We use three key components to describe PDMPs:

The deterministic dynamics: a system of differential equations that characterize the process' deterministic behaviour between jumps.

The event rate: a function that determines the intensity of jumps at each state.

The transition distribution: a probability measure that determines the next state of the process according to certain transition kernel.

BPS is a special case of PDMPs with certain choices of the three aforementioned components. We denote the state of PDMPs by  $Y = (\boldsymbol{W}, \boldsymbol{V})$ , which encodes the position and velocity of the particle. Let  $\zeta(\boldsymbol{w}, \boldsymbol{v}) = \zeta(\boldsymbol{v})\zeta(\boldsymbol{w})$ , where  $\zeta(\boldsymbol{v})$  represents the standard multivariate Gaussian distribution and  $\zeta(\boldsymbol{w})$  represents target posterior distribution of interest. The BPS keeps  $\zeta(\boldsymbol{w}, \boldsymbol{v})$  invariant. We denote  $U(\boldsymbol{w})$  as the associated energy function, which is the negative log of the unnormalized target posterior density function and is assumed continuously differentiable. In BPS, the specific choices for the three components are:

The deterministic dynamics:

$$\frac{d\mathbf{w}}{dt} = \mathbf{v}, \frac{d\mathbf{v}}{dt} = 0.$$

The event rate: the intensity  $\lambda(\boldsymbol{w}, \boldsymbol{v}) = \max\{0, \langle \nabla U(\boldsymbol{w}), \boldsymbol{v} \rangle\}$  of an inhomogeneous PP according to which the jumps take place.

The transition distribution: a Dirac centered at  $(\boldsymbol{w}, T_{w}(\boldsymbol{v}))$ , where

$$T_{\mathrm{w}}(oldsymbol{v}) = oldsymbol{v} - 2 rac{\langle 
abla U(oldsymbol{w}), oldsymbol{v} 
angle}{||
abla U(oldsymbol{w})||^2} 
abla U(oldsymbol{w}).$$

To ensure the ergodicity of the Markov chains, refreshment events also happen at random times while the velocity are refreshed according to certain distributions such as a standard Normal distribution. In summary, we present the BPS algorithm in Algorithm 1.

#### 2.4.2 Basics of LBPS

If the target posterior density  $\zeta(\boldsymbol{w})$  can be represented as a product of positive factors in Equation 11,

$$\zeta(\boldsymbol{w}) \propto \prod_{f \in \mathcal{F}} \gamma_f(N_f),$$
(11)

# Algorithm 1 BPS algorithm

#### 1: Initialization:

Initialize the particle position and velocity  $(\boldsymbol{w}^{(0)}, \boldsymbol{v}^{(0)})$  arbitrarily on  $\mathbb{R}^d \times \mathbb{R}^d$ . Set T to certain fixed trajectory length.

- 2: **for** i = 1, 2, ..., do
- 3: Sample the first arrival times  $\tau_{\text{ref}}$  and  $\tau_{\text{bounce}}$  of PP with intensity  $\lambda^{\text{ref}}$  and  $\lambda^{\text{bounce}}$  respectively, where  $\lambda^{\text{bounce}} = \max \left( \langle \boldsymbol{v}^{(i-1)}, \nabla U(\boldsymbol{w}^{(i-1)} + \boldsymbol{v}^{(i-1)}t) \rangle, 0 \right)$  and the value of  $\lambda^{\text{ref}}$  is pre-fixed.
- 4: Set  $\tau_i \leftarrow \min(\tau_{\text{bounce}}, \tau_{\text{ref}})$ .
- 5: Update the position of the particle via  $\boldsymbol{w}^{(i)} \leftarrow \boldsymbol{w}^{(i-1)} + \boldsymbol{v}^{(i-1)} \tau_i$ .
- 6: If  $\tau_i = \tau_{\text{ref}}$ , sample the next velocity  $\mathbf{v}^{(i)} \sim \mathcal{N}(0_d, I_d)$ .
- 7: If  $\tau_i = \tau_{\text{bounce}}$ , obtain the next velocity  $v^{(i)}$  by applying the transition function  $T_{\mathbf{w}^{(i)}}(v^{(i-1)})$ .
- 8: If  $t_i = \sum_{j=1}^i \tau_j \geqslant T$ , exit For Loop (line 2).
- 9: end for

where  $\mathcal{F}$  is the set of index for all factors in the target density and  $N_f$  represents the subset of variables connected to factor f, then a "local" version of BPS referred to as LBPS can be computationally efficient by taking advantage of the structural properties of the target density, especially if the target density has strong sparsity property described in definition 2. When LBPS is used, computationally cheaper refreshment scheme such as "local refreshment" (Bouchard-Côté et al., 2018) is often used by exploiting the structure for the factor graph. In the local refreshment scheme, one factor  $\gamma_f$  is chosen uniformly at random first and only the components of v with indices corresponding to variables  $N_f$  are resampled. The candidate collision time is recomputed only for the extended neighbour factors  $\gamma_{f'}$  of  $\gamma_f$ , where  $N_{f'} \cap N_f \neq \emptyset$ .

Now we present a brief description of LBPS algorithm and we will see how sparsity plays an important role in improving the computational efficiency. Detailed description about an efficient implementation of LBPS via the priority queue can be found in Bouchard-Côté et al. (2018). If the target unnormalized posterior density can be factorized according to Equation 11, its associated energy function is

$$U(\boldsymbol{w}) = \sum_{f \in \mathcal{F}} U_f(N_f),\tag{12}$$

where  $U_f(N_f) = -\log(\gamma_f(N_f))$ . We define the local intensity  $\lambda_f(\boldsymbol{w}, \boldsymbol{v})$  and local transition function  $T_{\rm w}^f(\boldsymbol{v})$  for factor  $\gamma_f$  as:

$$\lambda_f(\boldsymbol{w}, \boldsymbol{v}) = \max\{0, \langle \nabla U_f(\boldsymbol{w}), \boldsymbol{v} \rangle\}$$
(13)

$$T_{\mathbf{w}}^{f}(\boldsymbol{v}) = \boldsymbol{v} - 2 \frac{\langle \nabla U_f(\boldsymbol{w}), \boldsymbol{v} \rangle}{||\nabla U_f(\boldsymbol{w})||^2} \nabla U_f(\boldsymbol{w}). \tag{14}$$

It is worth noting that for variables that are not the neighbour variables for factor  $\gamma_f$ ,  $T_{\mathbf{w}}^f(\boldsymbol{v})_k = \boldsymbol{v}_k$ . In LBPS, the next collision time is the first arrival time of a PP with intensity  $\lambda(\boldsymbol{w},\boldsymbol{v}) = \sum_{f\in\mathcal{F}} \lambda_f(\boldsymbol{w},\boldsymbol{v})$ . We can sample the first arrival time via the superposition algorithm for a PP. We sample  $\tau_f$  for each factor  $\gamma_f$  from a PP with intensity  $\lambda_f(\boldsymbol{w},\boldsymbol{v}) = \max\{0, \langle \nabla U_f(\boldsymbol{w}), \boldsymbol{v} \rangle\}$ . The first arrival time for PP with intensity  $\lambda(\boldsymbol{w},\boldsymbol{v})$  is

 $\tau = \min_{f \in T} \tau_f$ . Once a bounce event takes place, LBPS only needs to manipulate a subset of the variables and factors. To be specific, if we denote  $f^*$  as the factor index such that  $\tau_{f^*} = \min_{f \in \mathcal{F}} \tau_f$ , we use  $\gamma_{f^*}$  to denote the collision factor. Following the priority queue implementation described in Bouchard-Côté et al. (2018), LBPS will update the position of neighbour variables for the collision factor  $\gamma_{f^*}$  and then apply the corresponding local transition function  $T_{\rm w}^{f^*}(\boldsymbol{v})$  to update the velocity. Next, the algorithm will sample the candidate bounce time of the next event for all the extended neighbour factors  $\gamma_{f'}$  for factor  $\gamma_{f^*}$  such that  $N_{f'} \cap N_{f^*} \neq \emptyset$ . If strong sparsity is satisfied for a family of factor graphs, then the number of neighbour variables for the collision factor grows much slower than the dimension of the parameter space and is negligible compared with the dimension of the parameter space. The number of operations needed to update the position of the neighbour variables is negligible compared to the dimension of the parameter space. When computing the candidate collision time for the extended neighbour factors, the number of factors involved is also negligible compared to the total number of factors in the factor graph, which is desirable in LBPS. We need to point out that although LBPS manipulates a subset of variables and factors, each local bounce will lead to changes in all the variables not just the neighbour variables connected with the current collision factor.

# 3. Proposed Bayesian GLM chain GTR model

In Section 2.2, we introduced a function  $\eta(\{x,x'\})$  which embeds state transitions into the integer. We will call this  $\eta$  an "order function" into the  $\eta$ -provided order. Here we use such order function to model the exchangeable rates between each pair of states through the current pair and its nearest neighbour pair. Thus, we would like to find an order of each pair of distinctive states that is biologically meaningful so that they can share statistical strength. Take protein evolution as an example. Most frequent amino acids exchanges take place between residues with similar physicochemical properties. Under the Bayesian GLM chain GTR model we proposed below, we assume that mutations between two pairs of amino acids sharing one common state are expected to have similar exchangeable rates if the unshared distinctive states between the two pairs have similar physicochemical properties. We would like to find an ordering  $\eta(\{x,x'\})$  such that amino acid pairs that are nearest neighbours share the most similar physicochemical properties given one common state. The Nearest Neighbour Pairwise Amino Acid Ordering (NNPAAO) Algorithm (described in Appendix F) proposed in Section 7 provides an ordering of amino acid pairs that satisfies such property. We denote it as  $\eta_{\text{dist}}(\{x, x'\})$  since the order is determined based on the Euclidean distance between amino acids defined in Equation 18. Given this ordering, neighbour pairs share more biological similarities than non-neighbour pairs.

In this paper, based on the intuition of  $\eta_{\text{dist}}(\{x, x'\})$ , we define a novel Bayesian GLM chain GTR model with feature function

$$\phi_i^{\text{chain}}(\{x, x'\}) := \mathbf{1}(\eta(\{x, x'\}) \in \{i, i+1\}),$$
 (15)

where  $i = 1, 2, ..., |\mathcal{X}|(|\mathcal{X}| - 1)/2$ . If we set  $\eta((x, x')) = \eta_{\text{dist}}(\{x, x'\})$ , the intuition of the Bayesian GLM chain GTR model is that the value of each exchangeable parameter  $\theta_{\{x, x'\}}$ 

depend on  $\{x, x'\}$  and its neighbour pair denoted as  $\{x'', x'''\}$ . Its neighbour pair satisfies that  $|\eta(\{x'', x'''\}) - \eta(\{x, x'\})| = 1$ . Equivalently, we can define the Bayesian **chain** GTR model through the exchangeable parameters  $\theta_{\{x, x'\}}(\boldsymbol{w}^b)$  by plugging in the definition of  $\phi_i^{\text{chain}}(\{x, x'\})$  in Equation 15:

$$\theta_{\{x,x'\}}(\boldsymbol{w}^b) = \exp\left\{\langle \boldsymbol{w}^b, \boldsymbol{\phi}^{\text{chain}}(\{x,x'\})\rangle\right\}$$

$$= \begin{cases} \exp\left(\boldsymbol{w}^b_{\eta(\{x,x'\})} + \boldsymbol{w}^b_{\eta(\{x,x'\})-1}\right), & \text{if } \eta(\{x,x'\}) = 2, 3, \dots, (|\mathcal{X}|(|\mathcal{X}|-1)/2), \\ \exp\left(\boldsymbol{w}^b_{\eta(\{x,x'\})}\right), & \text{if } \eta(\{x,x'\}) = 1. \end{cases}$$

For simplicity, we focus on the simple example of chains, but it does not need to be a chain. Other sparse graphs would also be suitable. We provide a general characterization of the running time analysis for both arbitrary factor graphs and factor graphs with sparse structure in Section 4.2. The chain GTR model is general since it is able to represent all rate matrices. We provide this proof in Appendix A.

### 4. Characterization of factor graphs where LBPS can be efficient

#### 4.1 Strong sparsity of factor graphs

A factor graph is bipartite graph used to represent the factorization of a function. We proposed our notion of sparsity for factor graphs. The sparsity definition is useful in analyzing the computational performance of LBPS. The standard notion of sparsity of a complete graph G = (V, E) refers to  $|E| = o(|V|^2)$ , where |E| represents the number of edges and |V| represents the number of vertices in the graph. However, the definition of the sparsity defined in our paper is different and more stringent than the conventional definition of sparsity. The motivation for proposing the sparsity definition for factor graphs is we expect LBPS to be efficient when strong sparsity (shown in definition 2) is satisfied.

**Definition 1.** Given a factorization of a function

$$f(\boldsymbol{w}) = \prod_{f=1}^{m} \gamma_f(N_f),$$

the factor graph  $G = (w, \Gamma, E)$  consist of variables  $w = \{w_1, w_2, \ldots, w_p\}$ , set of factors  $\Gamma = \{\gamma_1, \gamma_2, \ldots, \gamma_m\}$  and edges E. An undirected edge exists between factor  $\gamma_f$  and variable  $w_k$  iff  $w_k \in N_f$ . The neighbour variables for factor  $\gamma_f$  is denoted as  $N_f$ ,  $N_f \subset \{w_1, w_2, \ldots, w_p\}$ . The neighbour factors for variable  $w_k$  is  $S_k := \{\gamma_f : w_k \in N_f\}$ . The extended neighbour variables for factor  $\gamma_f$  is  $\overline{N}_f := \{w_k : w_k \in N_{f'} \text{ such that } N_{f'} \cap N_f \neq \emptyset\}$ . The extended neighbour factors for factor  $\gamma_f$  is  $\overline{S}_f := \{\gamma_{f'} : N_{f'} \cap N_f \neq \emptyset\}$ .

In Figure 2, we use the following factor graph to illustrate definition 1:

$$f(\mathbf{w}) = \gamma_1(w_1)\gamma_2(w_1, w_2)\gamma_3(w_2)\gamma_4(w_2, w_3)\gamma_5(w_3, w_4).$$

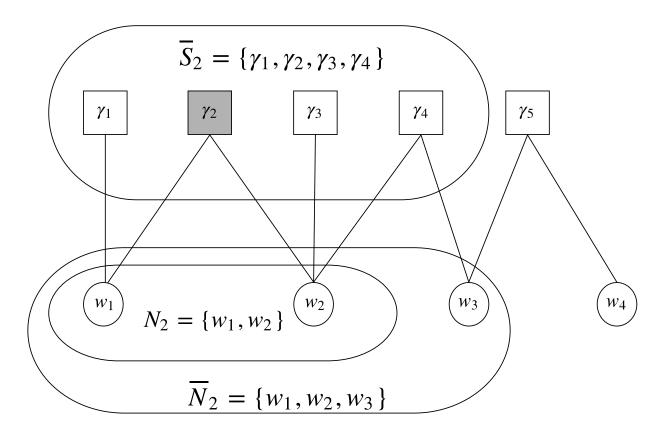


Figure 2: For factor  $\gamma_2$ , its neighbour variables  $N_2 = \{w_1, w_2\}$ . Its extended neighbour variables  $\overline{N}_2 = \{w_1, w_2, w_3\}$ . Its extended neighbour factors  $\overline{S}_2 = \{\gamma_1, \gamma_2, \gamma_3, \gamma_4\}$ .

Given definition 1, we define the **strong sparsity** for factor graphs below. We need a family of factor graphs, where the number of factors and the dimension of the variables can go to infinity. For simplicity, we omit m from notation  $\overline{N}_f^{(m)}$  and  $\overline{S}_f^{(m)}$  but use  $\overline{N}_f$  and  $\overline{S}_f$  instead.

**Definition 2.** A factor graph  $G = (\boldsymbol{w}, \Gamma, E)$  consist of  $\boldsymbol{w} = \{w_1, w_2, \dots, w_p\}$ ,  $\Gamma = \{\gamma_1, \gamma_2, \dots, \gamma_m\}$  and edges E, the **strong sparsity** requires that  $\max_f |\overline{S}_f| = o(m)$ , for  $1 \le k \le p$ , and  $\max_f |\overline{N}_f| = o(p)$ , for  $1 \le f \le m$ .

The strong sparsity reflects that for any factor in the factor graph, the number of its extended neighbour factors and extended neighbour variables grows much slower than the dimension of the parameter space and is ultimately negligible compared with the dimension of the parameter space and the total number of factors in the factor graph. We will show

later that the factor graph under our proposed sampling scheme satisfies the strong sparsity property.

### 4.2 Running time analysis of LBPS for factor graphs

Recently, there have been research (Deligiannidis et al., 2018; Andrieu et al., 2018) investigating the scaling limits of samplers constructed from PDMPs including BPS and Zig-zag processes. However, up to date there is no running time analysis for LBPS. Our running time analysis provides a framework for assessing the computational cost for one iteration of LBPS under both a general factor graph  $G = (\boldsymbol{w}, \Gamma, E)$  and also factor graphs satisfying strong sparsity property in definition 2.

In  $G = (\boldsymbol{w}, \Gamma, E)$ , we have variables  $\boldsymbol{w} = \{w_1, w_2, \dots, w_p\}$  and factors  $\Gamma = \{\gamma_1, \gamma_2, \dots, \gamma_m\}$ . We introduce some notations which are needed in the analysis. Denote  $|\overline{N}_*| = \max_f |\overline{N}_f|, |N_*|$  =  $\max_f |N_f|$  and  $|\overline{S}_*| = \max_f |\overline{S}_f|$ . Since  $N_f \subset \overline{N}_f$ , we have  $|N_f| \leqslant |\overline{N}_f|$  and  $|N_*| \leqslant |\overline{N}_*|$ . Denote the cost for computing the collision time for a factor  $\gamma_f$  as  $c_f$  and  $c_* = \max_f c_f$ . Let  $c_{U_f}$  denote the running time for computing  $\nabla U_f(\boldsymbol{w})$  and  $c_{U_*} = \max_f c_{U_f}$ . Now, we analyze the running time of LBPS via a priority queue implementation (Bouchard-Côté et al., 2018) for factor graph  $G = (\boldsymbol{w}, \Gamma, E)$ .

- 1. Compute the collision time for all factors and build a priority queue:  $\mathcal{O}(mc_* + m \log m)$ .
- 2. (a) If a collision takes place,
  - i. Update the position of extended neighbour variables:  $\mathcal{O}(\left|\overline{N}_*\right|)$ .
  - ii. Update the velocity of neighbour variables according to Equation 14:  $\mathcal{O}(|N_*| + c_{U_*}) \leq \mathcal{O}(|\overline{N}_*| + c_{U_*})$ .
  - iii. Add new samples to the trajectory list  $L_k$ :  $\mathcal{O}(|N_*|)$ . We use  $L_k$  to denote a list of triplets  $\left(w_k^{(i)}, v_k^{(i)}, t_k^{(i)}\right)$ , with  $w_k^{(0)}$  and  $v_k^{(0)}$  representing the initial position, and velocity and  $t_k^{(0)} = 0$ . For i > 0,  $v_k^{(i)}$  represents the velocity after the ith collision or refreshment event and  $t_k^{(i)}$  represents the time for the ith event.
  - iv. Compute the collision time for extended neighbour factors:  $\mathcal{O}(|\overline{S}_*|c_*)$ .
  - v. Insert the extended neighbour factors and the corresponding collision time into the priority queue:  $\mathcal{O}(\log m)$ . The elements of the priority queue are stored according to an increasing order of the collision time.
  - (b) If a refreshment event takes place and a local refreshment scheme (described in Section 2.4.2) is used:
    - i. Pick a factor uniformly at random and refresh the velocity:  $\mathcal{O}(|N_*|)$ .
    - ii. Update the position of the neighbour variables for the selected factor, compute the collision time for the extended neighbour factors and update the corresponding collision time in the priority queue:  $\mathcal{O}(|N_*| + |\overline{S}_*| c_* + |\overline{S}_*| \log m) \leq \mathcal{O}(|\overline{N}_*| + |\overline{S}_*| c_* + |\overline{S}_*| \log m)$ .
    - iii. Sample the next refreshment time  $\mathcal{O}(1)$ .

Thus, assuming there are  $L_1$  collision events and  $L_2$  refreshment events when the particle travels along a fixed trajectory length, the total running time for LBPS is:

$$\mathcal{O}\left(mc_* + L_1\left(\left|\overline{N}_*\right| + c_{U_*} + \left|\overline{S}_*\right|c_* + \log m\right) + L_2\left(\left|\overline{N}_*\right| + \left|\overline{S}_*\right|c_* + \left|\overline{S}_*\right|\log m\right)\right).$$

For a factor graph with no sparsity, the total running time for LBPS is

$$\mathcal{O}\bigg(L_1(p + c_{U_*} + mc_* + \log m) + L_2(p + mc_* + m\log m)\bigg).$$

For a factor graph with strong sparsity, the total running time for LBPS is

$$\mathcal{O}\left(mc_* + L_1\left(p^{\alpha_1} + c_{U_*} + m^{\alpha_2}c_* + \log m\right) + L_2(p^{\alpha_1} + m^{\alpha_2}c_* + m^{\alpha_2}\log m)\right),\tag{16}$$

where  $0 \le \alpha_1 < 1, 0 \le \alpha_2 < 1$ .

However, we need to mention that our running time analysis is part of a larger picture since we do not discuss the computational cost to obtain approximately one independent sample. We will leave the scaling limit of  $L_1$  and  $L_2$  with dimension to future work.

### 5. Methodology

#### 5.1 Achieving strong sparsity through LBPS-HMC alternative

In this section, we first provide an overview of our proposed sampling scheme summarized in Algorithm 2. We provide strong sparsity property of a family of factor graphs in definition 2, where LBPS can be efficient. Our sampling scheme combines both HMC and LBPS denoted as LBPS-HMC instead of using only LBPS since under the combined sampling scheme, the strong sparsity in definition 2 is satisfied for the target factor graph representation of the potential energy function augmented with the sufficient statistics z(x) := (n(x), h(x), c(x)) defined in Equation 10. If we only use LBPS kernel, the strong sparsity of the corresponding factor graph representation is not satisfied. Due to the parameterization of the stationary distribution of  $w^u$  in Equation 4, for any factor in the factor graph, its extended neighbour variables and extended neighbour factors include all the variables and factors in the factor graph. Thus, LBPS will not be efficient. We explain this in details in Section 5.2. We provide the analytical solution to the collision time for different categories of factors in Section 5.3.1 and Section 5.3.2. Finally, we provide the running time analysis for one iteration of LBPS-HMC under our proposed Bayesian GLM chain GTR rate matrix parameterization.

#### 5.2 Factor graph representation under LBPS-HMC

In this section, we show the motivation for proposing a combined sampling scheme using HMC to update the univariate weights  $\mathbf{w}^u$  and using LBPS to update the bivariate weights  $\mathbf{w}^b$  instead of using only LBPS to update  $\mathbf{w} = (\mathbf{w}^u, \mathbf{w}^b)$ . We refer to the second sampling scheme using only LBPS as a naive sampling scheme. We illustrate this through the factor graph representation for the potential energy function with augmented sufficient statistics defined in Equation 10 under the combined sampling scheme and naive sampling scheme

#### **Algorithm 2** Proposed sampling scheme LBPS-HMC for CTMCs

#### 1: Initialization:

Initialize weights  $\mathbf{w}^0 = (\mathbf{w}^u, \mathbf{w}^b)$  from N(0, 1), where  $\mathbf{w}^u$  represents the weights corresponding to univariate features to obtain the stationary distribution and  $\mathbf{w}^b$  represents the weights associated with bivariate features used to compute the exchangeable parameters.

- 2: **for** i = 1, 2, ..., N, **do**
- 3: Compute the rate matrix Q given  $\mathbf{w}^{(i-1)}$  under Bayesian GLM rate matrices.
- 4: Use an end-point sampler (Tataru and Hobolth, 2011; Rao and Teh, 2013) to simulate a path given time interval  $\Delta_e$  with two consecutive observations observed at time points e = (t, t+1) of the time series according to Equation 1 using the cached uniformization technique.
- 5: Compute the aggregate sum of the sufficient statistics of all time series obtained in Step 4.
- 6: Update univariate weights  $\mathbf{w}^u$  for the stationary distribution via HMC:

$$(\boldsymbol{w}^{u})^{(i+1)} \mid (\boldsymbol{w}^{u}, \boldsymbol{w}^{b})^{(i)}, Z^{(i)} \sim \text{HMC}(\cdot \mid (\boldsymbol{w}^{u}, \boldsymbol{w}^{b})^{(i)}, Z^{(i)}, L, \epsilon),$$

where L and  $\epsilon$  are tuning parameters representing the number of leapfrog jumps and step size in HMC. Function evaluation and gradient calculation required by HMC are described in Equation 20 and Equation 21 in Appendix B.

7: Update bivariate weights  $w^b$  used to compute the exchangeable parameters:

$$\left(\boldsymbol{w}^{b}\right)^{(i+1)} \mid \left(\boldsymbol{w}^{u}, \boldsymbol{w}^{b}\right)^{(i)}, Z^{(i)} \sim \text{LBPS}(\cdot \mid \left(\boldsymbol{w}^{u}, \boldsymbol{w}^{b}\right)^{(i)}, Z^{(i)}, T),$$

where T is the tuning parameter in LBPS representing the fixed length of the trajectory.

#### 8: end for

separately. Here, without loss of generality, we omit the Normal factors coming from the prior distribution of w.

For simplicity, under both sampling schemes, we illustrate the problem using DNA evolution since there are only four states in the state space, where  $|\mathcal{X}| = \{A, C, G, T\}$ . Under the naive sampling scheme, we are considering a Bayesian GLM rate matrix with GTR features instead of the chain GTR features since the parameterization for bivariate weights  $w^b$ is simpler. Each exchangeable parameter  $\theta_{x,x'} = \exp(w_{x,x'})$  only depends on one element of the bivariate weights. While in the chain GTR model, each exchangeable parameter depends on two bivariate weight elements except  $\theta_{\{x,x'\}}$  depends on one bivariate weight  $\mathbf{w}_1^b$  when  $\eta(\{x,x'\})=1$  according to the definition in Equation 15. The root cause that makes the naive sampling scheme undesirable is that the parameterization of the stationary distribution makes  $w^u$  common neighbour variables for all factors. This is true for either models with GTR features or chain GTR features. Thus, under the naive sampling scheme, we present the factor graph under a Bayesian GLM rate matrix with GTR features for DNA evolution in Figure 3. In this figure,  $\mathbf{w}^u = (w_A, w_C, w_G, w_T)$  is used to obtain the stationary distribution. For example,  $\pi_A(\mathbf{w}) = \exp(w_A)/(\exp(w_A) + \exp(w_C) + \exp(w_G) + \exp(w_T))$ . The bivariate weights  $\mathbf{w}^b = (w_{AC}, w_{AG}, w_{AT}, w_{CG}, w_{CT}, w_{GT})$  are used to obtain the exchangeable parameters, for example,  $\theta_{AC} = \exp(w_{AC})$ . As shown in Equation 10, all factors such as  $C_{x,x'}$ ,  $H_{x,x'}$  depend on rate matrix element  $q_{x,x'}(\boldsymbol{w}) = \theta_{x,x'}(\boldsymbol{w}^b)\pi_{x'}(\boldsymbol{w}^u)$ , thus all factors have neighbour variables  $w^u$ . This indicates that for any factor, its extended neighbour bour factors includes all other factors since they all share common neighbour variables  $w^u$ . In both models, the number of factors is  $\mathcal{O}(|\mathcal{X}|^2)$ . Once a collision takes place, the candi-

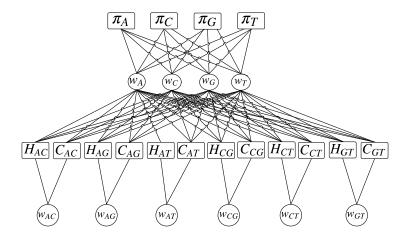
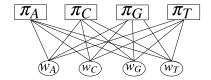


Figure 3: Factor graph under Bayesian GLM representation of GTR model for DNA evolution when using only LBPS to update  $\boldsymbol{w}$ , where  $H_{x,x'}$  and  $C_{x,x'}$  represent the sojourn time factor and transition count factor for states (x,x'), where  $x \neq x' \in \mathcal{X} = \{A,C,G,T\}$ . In the graph, we leave out the other half of the transition count factors and sojourn time factors for the symmetric pair of states (x',x) such as  $H_{CA},C_{CA},H_{GA},C_{GA},\ldots$  for representation simplicity.

date collision time for an order of  $\mathcal{O}(|\mathcal{X}|^2)$  extended neighbour factors need to be computed under a naive sampling scheme. This can be computationally expensive under a large state space  $\mathcal{X}$  and is not desirable.

To resolve the computational issue, we propose a sampling scheme using HMC to sample  $\boldsymbol{w}^u$  and LBPS to sample  $\boldsymbol{w}^b$  respectively. Under our proposed sampling scheme and Bayesian GLM chain GTR parameterization, the factor graph is shown in Figure 4. In this situation, we provide the factor graph for a model with chain GTR features instead of GTR features since later in the real data analysis of the protein evolution, we are considering a model with chain GTR features. We illustrate here why a model with chain GTR features has only a constant number of extended neighbour factors for any bivariate factor in DNA evolution when  $|\mathcal{X}| = 4$ . This conclusion generalizes naturally to protein evolution when  $|\mathcal{X}| = 20$ . Thus, the strong sparsity of the factor graph is satisfied under Bayesian GLM chain GTR using LBPS-HMC.

Regardless of the size of the state space  $\mathcal{X}$ , we prove that for any  $\{x, x'\} \in \mathcal{X}^{\text{unordered,dist}}$ , any sojourn time factor  $H_{x,x'}$  or transition count  $C_{x,x'}$  has at most twelve extended neighbour factors and three extended neighbour variables under Bayesian GLM chain GTR using the proposed sampling scheme. For the pairs of states  $\{x, x'\}$ , when  $\eta(\{x, x'\}) \neq |\mathcal{X}|(|\mathcal{X}|-1)/2$ , the twelve extended neighbour factors of factor  $H_{x,x'}$  or  $C_{x,x'}$  include the factors only connected with  $\boldsymbol{w}_i^b$  and  $\boldsymbol{w}_{i+1}^b$ , where  $i = \eta(\{x, x'\}) - 2, \eta(\{x, x'\}) - 1$ , and  $\eta(\{x, x'\})$ . Similarly, for the pair of states  $\{x, x'\}$  such that  $\eta(\{x, x'\}) = |\mathcal{X}|(|\mathcal{X}|-1)/2, H_{x,x'}$  or  $C_{x,x'}$  has eight extended neighbour factors, where each of the eight factors is only connected with



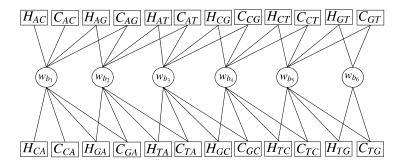


Figure 4: Factor graph when using Bayesian GLM representation of our proposed chain GTR model for DNA evolution when using HMC to update univariate weights  $\boldsymbol{w}^u$  and LBPS to update bivariate weights  $\boldsymbol{w}^b$ , where  $\boldsymbol{w} = (\boldsymbol{w}^u, \boldsymbol{w}^b)$ .

 $m{w}_i^b$  and  $m{w}_{i+1}^b$ , where  $i=\eta(\{x,x'\})-2$  and  $\eta(\{x,x'\})-1$ . All factors have three extended neighbour variables  $m{w}_{\eta(\{x,x'\})-1}^b, m{w}_{\eta(\{x,x'\})}^b$  and  $m{w}_{\eta(\{x,x'\})+1}^b$ , except that factor  $H_{x'',x'''}$  or  $C_{x'',x'''}$  has two extended neighbour variables  $m{w}_1^b$  and  $m{w}_2^b$  when  $\eta(\{x'',x'''\})=1$ .

Thus, when updating the position of extended neighbour variables or computing the candidate bounce time for extended neighbour factors of the collision factor, there is only a small constant number of extended neighbour variables and extended neighbour factors involved. We can conclude that the factor graph in Figure 4 satisfies the desirable strong sparsity in definition 2. In the next section, we will show that we can derive analytical solution to the candidate bounce time for each category of factors.

# 5.3 Analytical bounce time solution to related factors

In this section, we decompose the potential energy  $U(\mathbf{w})$  in Equation 10 with augmented sufficient statistics into a sum of factors and organized them into the following categories:

Normal factor:  $\frac{1}{2}\kappa \|\boldsymbol{w}\|_2^2$ ,

Sojourn time factor:  $H_{x,x'} \coloneqq h_x q_{x,x'}\left(\boldsymbol{w}\right)$ , for all  $(x,x') \in \mathcal{X}^{\text{distinct}}$ ,

**Transition count factor:**  $C_{x,x'} := -c_{x,x'} \log (q_{x,x'}(\boldsymbol{w}))$ , for all  $(x,x') \in \mathcal{X}^{\text{distinct}}$ ,

Initial count factor:  $\pi_x := -n_x \log(\pi_x(\boldsymbol{w}))$ , for all  $x \in \mathcal{X}$ , with  $\pi_x$  representing the initial count factor but  $\pi_x(\boldsymbol{w})$  denotes the stationary distribution for state x given  $\boldsymbol{w}$ .

We derive the analytical solution to the collision time for sojourn time factors and transition count factors respectively. The analytical collision time solution for a Normal distribution can be found in Bouchard-Côté et al. (2018), so we focus on the other factors.

#### 5.3.1 Local sojourn time factors

We first sample the energy gap  $-\log(E) > 0$  with  $E \sim \mathcal{U}(0,1)$ . It represents the difference between the energy denoted as  $E_0$  at the current position of a particle and a higher energy denoted as  $E_0 - \log(E)$  in the parameter space. Our task is to determine the collision time that the particle can travel along the energy ladder given its current velocity until it hits the energy barrier  $E_0 - \log(E)$  in the *i*th iteration.

Under a reversible model, the corresponding potential energy of the particle is:

$$U(\boldsymbol{w}_0) = h_x^{(i)} q_{x,x'}^{(\text{rev})}(\boldsymbol{w}_0)$$
  
=  $h_x^{(i)} \pi_{x'} \exp(\langle \boldsymbol{w}_0, \boldsymbol{\phi}(\{x, x'\}) \rangle).$ 

The potential energy after a time interval  $\Delta$ :

$$U(\boldsymbol{w}_0 + \boldsymbol{v}\Delta) = h_x^{(i)} q_{x,x'}^{(\text{rev})} (\boldsymbol{w}_0 + \boldsymbol{v}\Delta)$$
$$= h_x^{(i)} \pi_{x'} \exp(\langle \boldsymbol{w}_0 + \boldsymbol{v}\Delta, \boldsymbol{\phi}(\{x, x'\})\rangle).$$

We observe that the particle is travelling to a higher energy area if and only if  $\langle v, \phi(\{x, x'\}) \rangle > 0$ . Therefore, we set

$$-\log(E) = U(\boldsymbol{w}_0 + \boldsymbol{v}\Delta) - U(\boldsymbol{w}_0)$$
$$= h_x^{(i)} q_{x,x'}^{(\text{rev})}(\boldsymbol{w}_0) \left( \exp(\langle \boldsymbol{v}\Delta, \boldsymbol{\phi}(\{x, x'\})\rangle) - 1 \right).$$

We denote  $c = -\log(E) > 0$  and obtain:

$$\Delta = \begin{cases} \frac{1}{\langle \boldsymbol{v}, \boldsymbol{\phi}(\{x, x'\}) \rangle} \log \left( \frac{c}{h_x^{(i)} q_{x, x'}^{(\text{rev})}(\boldsymbol{w}_0)} + 1 \right) & \text{if } \langle \boldsymbol{v}, \boldsymbol{\phi}(\{x, x'\}) \rangle > 0, \\ \infty, & \text{otherwise.} \end{cases}$$

#### 5.3.2 Local factors for transition counts

Similarly, the transition count factor for pairs of states (x, x') in the *i*th iteration is  $-c_{x,x'}^{(i)} \log \left(q_{x,x'}^{(\text{rev})}(\boldsymbol{w}_0)\right)$ , and the gradient is  $-c_{x,x'}^{(i)} \phi(\{x,x'\})$ . Thus, the corresponding potential energy of the particle is:

$$U(\boldsymbol{w}_0) = -c_{x,x'}^{(i)} \log \left( q_{x,x'}^{(\text{rev})}(\boldsymbol{w}_0) \right)$$
$$= -c_{x,x'}^{(i)} \left( \log \pi_{x'} + \langle \boldsymbol{w}_0, \boldsymbol{\phi}(\{x, x'\}) \rangle \right)$$

The potential energy after a time interval  $\Delta$ :

$$U(\boldsymbol{w}_0 + \boldsymbol{v}\Delta) = -c_{x,x'}^{(i)} \left( \log(\pi_{x'}) + \langle \boldsymbol{w}_0 + \boldsymbol{v}\Delta, \phi(\{x,x'\}) \rangle \right)$$

We sample  $E \sim \mathcal{U}(0,1)$ , set  $c = -\log(E)$ ,  $U(\boldsymbol{w}_0 + \boldsymbol{v}\Delta) - U(\boldsymbol{w}_0) = c$ , and obtain:

$$\Delta = \begin{cases} -\frac{c}{c_{x,x'}^{(i)} \langle \boldsymbol{v}, \boldsymbol{\phi}(\{x,x'\}) \rangle} & \text{if } \langle \boldsymbol{v}, \boldsymbol{\phi}(\{x,x'\}) \rangle \leqslant 0, \\ \infty, & \text{otherwise.} \end{cases}$$

# 5.4 Running time analysis of LBPS-HMC

In Section 4.2, we have provided the running time analysis of LBPS under a general factor graph G. In this section, we provide the running time analysis of LBPS-HMC alternative under Bayesian GLM chain GTR model for CTMCs. Recall that the computational cost for a factor graph G with sparse property is given in Equation 16:

$$\mathcal{O}\bigg(mc_* + L_1(p^{\alpha_1} + c_{U_*} + m^{\alpha_2}c_* + \log m) + L_2(p^{\alpha_1} + m^{\alpha_2}c_* + m^{\alpha_2}\log m)\bigg),$$

assuming there are  $L_1$  collision events and  $L_2$  refreshment events when the particle travels along a fixed trajectory length.

In the special case of using LBPS-HMC under Bayesian GLM chain GTR for CTMCs, we have  $c_* \leq \mathcal{O}(1), c_{U_*} \leq \mathcal{O}(1), \left|\overline{N}_*\right| \leq \mathcal{O}(1)$  and  $\left|\overline{S}_*\right| \leq \mathcal{O}(1)$ . The computational cost we derive not only holds for Bayesian GLM chain GTR model but also holds for any Bayesian GLM model that satisfies  $c_* \leq \mathcal{O}(1), c_{U_*} \leq \mathcal{O}(1), \left|\overline{N}_*\right| \leq \mathcal{O}(1)$  and  $\left|\overline{S}_*\right| \leq \mathcal{O}(1)$ . Under the Bayesian GLM chain GTR rate matrix parameterization, we have  $m \leq \mathcal{O}\left(|\mathcal{X}|^2\right)$  and  $p \leq \mathcal{O}\left(|\mathcal{X}|^2\right)$ .

- 1. Sample the auxiliary sufficient statistics using end-point sampler:  $\mathcal{O}(|\mathcal{X}|^3)$ .
- 2. Update the univariate weights  $\boldsymbol{w}^u$  using HMC:  $\mathcal{O}(L_0|\mathcal{X}|)$ , where  $L_0$  represents the number of leapfrog steps in one HMC iteration. When updating  $\boldsymbol{w}^u$ , the number of non-zero entries for all possible features is  $|\mathcal{X}|$ .
- 3. Update the bivariate weights  $\boldsymbol{w}^b$  using LBPS:  $\mathcal{O}(L_1 \log |\mathcal{X}| + L_2 \log |\mathcal{X}|)$ . This is obtained by plugging  $\alpha_1 = \alpha_2 = 0$  into Equation 16, where  $\alpha_1 = \alpha_2 = 0$  is obtained since we have  $|\overline{N}_*| \leq \mathcal{O}(1)$  and  $|\overline{S}_*| \leq \mathcal{O}(1)$ . A local refreshment scheme is adopted.

Thus, the total cost of one iteration of our LBPS-HMC alternative takes

$$\mathcal{O}(L_0|\mathcal{X}| + L_1 \log |\mathcal{X}| + L_2 \log |\mathcal{X}| + |\mathcal{X}|^3). \tag{17}$$

In comparison, one HMC iteration for Bayesian GLM chain GTR takes  $\mathcal{O}\left(L|\mathcal{X}|^2 + |\mathcal{X}|^3\right)$  since the total number of non-zero entries for all possible features of Bayesian GLM chain GTR is  $|\mathcal{X}|(|\mathcal{X}|+1)/2$  using only HMC, where L is the number of leapfrog steps in one HMC iteration. In our experiments, we find that  $|\mathcal{X}|$ ,  $\log |\mathcal{X}|$  are both lower order terms compared with  $|\mathcal{X}|^2$  using HMC. This also explains why our LBPS-HMC is more computationally efficient than HMC. The running time analysis discussed here is not a whole story since we do not provide information on how  $L_1$  and  $L_2$  should be scaled in order to obtain a constant number of effective samples. To provide the readers an idea of how L in HMC scales, consider the case of normally distributed random vectors of dimension p with an identity covariance matrix. This is not the same as our problem or any real practical problem, but it may shed

light on how the number of leapfrog steps L scales with dimensionality for some problems. To reach a nearly independent point, the number of leapfrog updates under the independent and identically distributed normal case grows as  $\mathcal{O}\left(p^{\frac{1}{4}}\right)$  (Neal et al., 2011), where p is the dimension of the parameters. For LBPS, under the independent and identically distributed normal distribution with dimension p or weakly dependent case, we expect  $L_1$  to grow as O(p).

# 6. Experiments

All numerical experiments are run on Compute Canada resources "cedar" with Intel "Broadwell" CPUs at 2.1Ghz and an E5-2683 v4 model. A detailed description of "cedar" can be found at https://docs.computecanada.ca/wiki/Cedar. The algorithm is implemented in Python 3.6.4. The implementation of the code can be found at https://github.com/zhaottcrystal/rejfreePy\_main. Across all experiments, the Effective Sample Size (ESS) is evaluated via the R package mcmcse (Flegal et al., 2012). We set the first 30% of the posterior samples as the burnin period to be discarded.

All our synthetic datasets are simulated under a Bayesian GLM chain GTR rate matrix parameterization. We assign a standard Normal distribution as the prior over both the univariate weights and the bivariate weights, which are both generated from  $\mathcal{U}(0,1)$  across all simulation studies. The refreshment rate of LBPS is set to one and the number of leapfrog jumps and the stepsize of each jump are set as L=40 and  $\epsilon=0.001$  in HMC. The combination of the two tuning parameters are simply chosen among the ones with the best computational efficiency given multiple combination of the parameter values.

All synthetic experiments generate 500 sequences given a total length of the observation time interval as three and the states of all sequences are observed every 0.5 unit time under a Bayesian GLM chain GTR with different dimensions of rate matrices. Later in this paper, we use "LBPS" for short to represent the combined sampling scheme LBPS-HMC, which uses HMC to sample the univariate weights and LBPS to sample the bivariate weights. When HMC sampling scheme is mentioned later, it refers to the sampling scheme using HMC to sample all the weight parameters  $\mathbf{w} = (\mathbf{w}^u, \mathbf{w}^b)$ .

In both the experiments using synthetic datasets in Section 6 and the real pair of data sequences in Section 7, we conduct computational efficiency comparison via ESS per second and we also check the correctness of the two samplers. For the synthetic data experiments, we perform Exact Invariance Test (EIT) in the spirit of Geweke (2004), compare the density of the posterior samples between the two samplers and compute the Absolute Relative Difference (ARD) defined in Appendix I. Since the two algorithms share the same limiting distribution, as we increase the number of iterations of the Markov chains, after an initial burnin period, the distribution of the posterior samples obtained from the two algorithms will be more and more similar and both will become closer to the target posterior distribution. For higher dimensional rate matrices, since the total number of parameters is  $\mathcal{O}(|\mathcal{X}|^2)$ , it is hard to display the density plots for all parameters, we use ARD as the metric to describe the similarities between the posterior samples from the two algorithms. Smaller values of ARD indicate more similarities. In Section 7 with the real data sequences, we only use ARD to check the similarity between the distributions of the posterior samples from the two algorithms and show the results in Appendix I.

#### 6.1 Correctness check of LBPS and HMC using exact invariance test

To check the correctness of our software implementation and derivation of analytic collision times, we use Kolmogorov-Smirnov tests to formally check that marginal-conditional and successive conditional distributions coincide, in the spirit of Geweke (2004). We performed tests on both LBPS and HMC kernels on various dimensions of the rate matrices (see Appendix D).

To further validate the correctness of LBPS and HMC, we compare the distribution of the posterior samples collected from LBPS and HMC respectively. The shared synthetic dataset is generated under an 8-by-8 Bayesian GLM chain GTR. The prior distribution for  $\boldsymbol{w}$  is standard Normal distribution. We obtain 40,000 and 10,000 posterior samples from LBPS and HMC separately with first 30% of the samples discarded. As we have a longer chain for both algorithms, we expect their density plots should be closer. We use the density plot of all exchangeable parameters in Figure 5 to illustrate this. The boxplot is provided in Figure 9 in Appendix E. We also provide summary statistics of the posterior samples among different parameters.

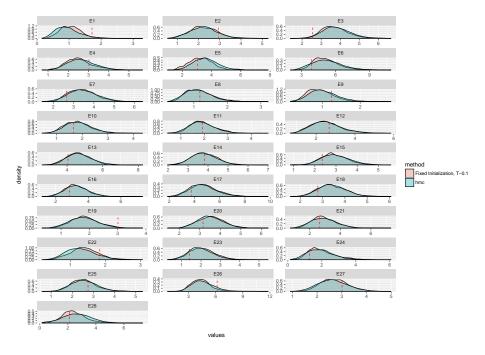


Figure 5: Density plot of posterior density comparison across exchangeable parameters of an 8-by-8 rate matrix between LBPS and HMC.

#### 6.2 Computational efficiency comparison between LBPS and HMC

In order to explore the scalability of our algorithm as the dimension of the parameters increases, we compared the ESS per second in the  $\log_{10}$  scale among all the parameters for different size of the rate matrices. The larger values of the ESS per second indicate better

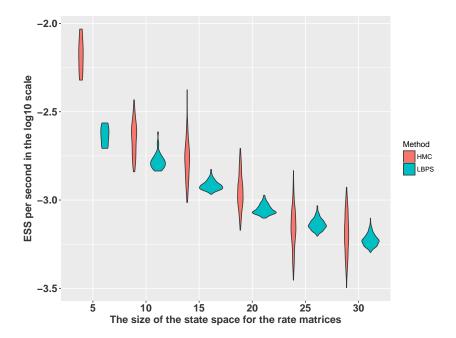


Figure 6: ESS per second in the  $\log_{10}$  with the dimension of the rate matrices ranging from 5-by-5, 10-by-10, ... to 30-by-30 with a step size of 5.

computational efficiency of the algorithm. The dimension of the rate matrices ranges from 5-by-5, 10-by-10, ... to 30-by-30. We obtain a total number of 60,000 posterior samples from LBPS and 10,000 samples from HMC. To speed up the actual running time of the experiments, for rate matrices of dimension lower than 15-by-15, the trajectory length is fixed as 0.1. For rate matrices with higher dimensions, the trajectory length is shortened to 0.05. According to the Bayesian GLM chain GTR model parameterization, the number of exchangeable parameters is  $|\mathcal{X}|(|\mathcal{X}|-1)/2$ . The result is displayed in Figure 6.

From Figure 6, we can see that when the dimension of the weight parameters is low, for a 5-by-5 rate matrix, HMC has better performance than LBPS. As the dimension increases, it shows that the minimum of the ESS per second for LBPS outperforms HMC. The larger variation in ESS per second of HMC in the violin plot in Figure 6 indicates that HMC is not efficient in exploring certain directions of the parameter space. This result suggests that LBPS has superior computational performance as the dimension increases. We suggest that LBPS should be adopted for high-dimensional CTMCs inference problems.

# 6.3 ESS per second using different trajectory lengths of LBPS comparisons

The computation efficiency of HMC has been found to depend highly on the choices of the tuning parameters, which are the number of leapfrog jumps L and the size of each jump  $\epsilon$ . Various strategies (Wang et al., 2013; Hoffman and Gelman, 2014; Zhao et al., 2016) have been developed to effectively tune the parameters. LBPS also involves a tuning

parameter, which is the fixed trajectory length. Thus, it is of interest to examine whether the computational efficiency of LBPS is sensitive to different choices of the trajectory length.

We simulate a synthetic dataset from a 20-by-20 rate matrix with 190 exchangeable parameters. A 20-by-20 rate matrix is chosen since it has the same dimension as protein evolution in the real data analysis. The synthetic dataset is the same when we perform sampling using LBPS with different trajectory lengths.

We use ESS per second as the metric to evaluate the computation efficiency of LBPS with different trajectory lengths. Our summary statistics include the minimum, first quantile, mean, median, third quantile and maximum of the ESS per second across 190 exchangeable parameters. 40,000 posterior samples of the parameters of interest are obtained using our algorithm. The actual walltime (in seconds) of our algorithm with fixed trajectory lengths at 0.025, 0.1, 0.15, 0.2 and 0.25 are shown in Table 1. Since the first 30% of the samples are discarded, the actual running time used to calculate the ESS per second are scaled by 70% of the total walltime in Table 1.

Trajectory Lengths	0.025	0.1	0.15	0.2	0.25
Walltime (in seconds)	240677	282233	338242	381565	403932

Table 1: Actual walltime of LBPS for 40,000 iterations with trajectory lengths at 0.025, 0.1, 0.15, 0.2, 0.25.

The results are shown in Figure 7. We have found that except the maximum of the ESS per second, the computational efficiency of LBPS is similar with different values of the trajectory length. Especially the minimum of ESS per second is very robust to different choices of trajectory lengths. The longer the trajectory lengths, the better performance of the maximum of the ESS per second. Similar conclusions are achieved in our real data analysis. The minimum of the ESS per second is more important than the maximum since the performance of a sampler depends on the direction that is the hardest to sample.

#### 7. Real Data Analysis

#### 7.1 Background

We use the real dataset from Zhao et al. (2016) with 641 amino acid sequences and each sequence has 415 sites from the protein kinase domain family. It is available at <a href="https://github.com/zhaottcrystal/rejfreePy\_main/tree/master/Dataset">https://github.com/zhaottcrystal/rejfreePy\_main/tree/master/Dataset</a>. In phylogenetics, the evolutionary process is often inferred using multiple homologous biological sequences under a evolutionary tree with the same rate matrix across the tree. For simplicity, we estimate the rate matrix from a pair of sequences. we pick randomly a pair of amino acid sequences to study the rate matrix from its posterior distribution.

In Section 2.2, we have provided intuition behind the chain GTR model since we assume that the exchangeable rates between a pair of states is affected by its neighbour pairs with similar biological properties. Then another question arises as how to pick a reasonable ordering of the pairs of amino acids to allow neighbour pairs share similar biological

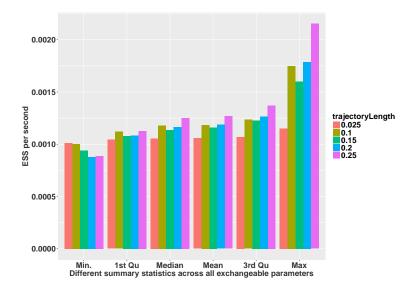


Figure 7: Effective sample size per second of the summary statistics across exchangeable parameters with respect to different trajectory lengths of LBPS.

properties. According to He et al. (2011), we determine the ordering of pairs of amino acids regarding their closeness based on their pairwise Euclidean distance defined by their physiochemical properties. The distance satisfies nonnegative, reflective, symmetric, and triangle properties and is given by Grantham (1974) as:

$$D_{ij} = \left(\alpha(c_i - c_j)^2 + \beta(p_i - p_j)^2 + \gamma(v_i - v_j)^2\right)^{\frac{1}{2}},\tag{18}$$

where c= composition, p= polarity, v= molecular volume, and  $\alpha,\beta,\gamma$  are defined in Grantham (1974). Larger pairwise distance indicates less similarity between pairs of amino acids, where mutation between them may be deteriorative. On the contrary, smaller distance indicates more similarities and tendencies to mutate between them. The pairwise Euclidian distance between amino acids is given Table 3 in He et al. (2011). We provide this table in Appendix F of our paper.

Given the distance in Appendix F, the ordering of pairs of amino acids regarding their closeness is determined according to our proposed NNPAAO Algorithm 3 in Appendix G. Some notations are introduced to help understand the algorithm. Denote AminoAcidDist as the pairwise Euclidean distance between amino acids shown in Appendix F,  $\mathcal{X}$  as the state space of 20 amino acids and  $\mathcal{X}^{\text{unordered,dist.}}$  as the set of unordered pairs of distinct amino acids, where  $|\chi|=20$  and  $|\mathcal{X}^{\text{unordered,dist.}}|=190$ . The algorithm outputs a dictionary AminoAcidsPairRank with keys from  $\mathcal{X}^{\text{unordered,dist.}}$  and the value associated with each key represents the rank of this amino acid pair. Neighbour pairs indicate closeness for similarities.

The algorithm first picks the amino acid pair with the smallest positive Euclidean distance from Appendix F, which is pair "IL", where D(I, L) = 5. The nearest neighbour to

"IL" is chosen among two candidate pairs which has either the smallest nonzero distance to "I" or "L", which are "IM" with D(I, M)=10 or "LM" with D(L, M)=15. Among them, we pick the pair with a smaller distance and this pair is defined as the nearest neighbour to "IL", which is "IM". Next our current pair is set as "IM" and similarly, we find the nearest neighbour to "IM" and we keep searching the nearest neighbour for the current pair until we have iterated over all pairs in  $\mathcal{X}^{\mathrm{unordered, dist.}}$ .

#### 7.2 Numerical results: computational efficiency comparison

We provide the correctness check via ARD between the two samplers in Appendix I. We compare the computational efficiency of LBPS and HMC via comparing the summary statistics of ESS per second across 190 exchangeable parameters of a 20-by-20 reversible rate matrix using the protein kinase domain family dataset in Zhao et al. (2016). Figure 8 demonstrates the computational advantages of using LBPS over HMC. The trajectory lengths of LBPS is chosen as 0.1, 0.15 and 0.2. We find that when the trajectory length is set at 0.2, LBPS outperforms HMC across all summary statistics of ESS per second. The minimum of ESS per second of HMC is markedly worse than LBPS, this implies that the inefficiency of HMC to explore certain direction of the parameter space. Under a fixed trajectory length of LBPS, there is not big difference between the first quantile and third quantile of the ESS per second across all parameters. It indicates that LBPS has similar computational efficiency across different directions of the parameter space. In Table 2, we provide the ESS runing HMC for 811380 seconds (9.4 days) compared with the ESS of LBPS only running 81.25% of the walltime of HMC, which further demonstrates the superior performance of LBPS compared with HMC. A traceplot scaled by running time of a selected exchangeable parameter showing better mixing using LBPS compared with HMC is provided in Figure 10 in Appendix H.

	Min.	1st Quantile	Median	Mean	3rd Quantile	Max.
LBPS (0.2)	773	1656	1864	1845	2108	2532
HMC	95	626	740	736	840	1233

Table 2: Summary of ESS across exchangeable parameters between HMC and LBPS with trajectory length 0.2.

#### 8. Discussion

In this paper, we have developed a computationally efficient sampling scheme combining LBPS and HMC to explore high dimensional CTMCs under a novel Bayesian GLM chain GTR rate matrix. In this model, we assume that the mutation rates of the amino acid evolutionary processes depends on its neighbour pairs with similar physiochemical properties.

We provide a framework for assessing the running time of our algorithm. We define the notion of sparsity through extended neighbourhood of each factor in the factor graph where LBPS can be efficient while taking advantage of the sparse structure. We analyze

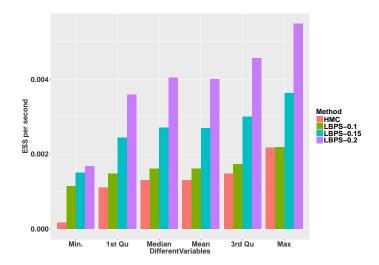


Figure 8: Summary statistics of ESS per second of across exchangeable parameters with different trajectory lengths of LBPS at 0.1, 0.15 and 0.2, compared with HMC.

the computational cost of one iteration of our algorithm for general factor graphs with no sparsity and factor graphs with strong sparsity.

In terms of the computational performance, we find that HMC may have difficulty in exploring with certain directions of the parameter space since there is a big discrepancy between ESS per second among different parameters and the minimum of ESS per second for HMC is much worse than LBPS. In the contrary, the computational performance of LBPS is similar across different directions of the parameter space. In the real application, LBPS outperforms HMC for all chosen summary statistics of ESS per second and especially the minimum of the ESS per second.

Moreover, via protein evolution, we provide the first proof-of-concept real data application for LBPS to explore high dimensional CTMCs compared with HMC. The experiment results suggest that LBPS may have great potential to explore high dimensional CTMCs. This shows promises for evolutionary processes such as codon evolution (Anisimova and Kosiol, 2008) with a large state space. Our sampling algorithm also has the potential to be applied to co-evolution of groups of interacting amino acids residues which is also modelled under large rate matrices.

For phylogenetic applications, the current limitation of our algorithm is that it only works for unnormalized, reversible rate matrices since only under this situation, the analytical solution to the collision time can be obtained. It is customary in phylogenetics to use normalized rate matrix, where the normalization ensures that the expected number of mutations is one given one unit time. To resolve the issue, for example, for pairwise sequences, we can first set the branch length to one unit length and obtain posterior distribution for an unnormalized reversible rate matrix. Later we can normalize the rate matrix by multiplying the normalization coefficient and the normalizing coefficient can be served as the estimate for the branch length.

It is also possible to combine our method with Bayesian variable selection techniques (Ishwaran et al., 2005) and a binary version of BPS (Pakman, 2017) to detect the features that have a large effect on the underlying CTMCs. To further facilitate the computational efficiency of our sampling algorithm, it is worth exploring adaptive version of the algorithm to choose the tuning parameter in LBPS wisely.

# Appendix A. Connection between Bayesian GLM GTR model and Bayesian GLM chain GTR model

We have defined  $\phi^{\text{gtr}}$  in Section 2.2 such that any GTR rate matrices can be represented via the Bayesian GLM rate matrix parameterization. Our objective is to prove that there exists  $\boldsymbol{w}_*$  such that for any rate matrices under the Bayesian GLM parameterization, for  $\forall \boldsymbol{w}, q_{x,x'}^{(\text{rev})}(\boldsymbol{w})$  can be represented under the Bayesian GLM chain GTR parameterization.

Zhao et al. (2016) have shown that for any rate matrix Q, we can find weights  $\boldsymbol{w}$  such that  $Q_{(x,x')} = q_{x,x'}(\boldsymbol{w}) = \exp\left\{\langle \boldsymbol{w}^b, \boldsymbol{\phi}^{\text{gtr}}(\{x,x'\})\rangle\right\} \pi_{x'}(\boldsymbol{w}^u)$  under the Bayesian GLM GTR model. In the Bayesian GLM chain GTR model, we show that there exists  $\boldsymbol{w}_*$  such that for  $\forall \boldsymbol{w} = \begin{pmatrix} \boldsymbol{w}^u \\ \boldsymbol{w}^b \end{pmatrix}$ ,

$$q_{x,x'}^{\text{gtr}}(\boldsymbol{w}) = \exp \left\{ \langle \boldsymbol{w}^b, \boldsymbol{\phi}^{\text{gtr}}(\{x,x'\}) \rangle \right\} \pi_{x'}(\boldsymbol{w}^u)$$

$$= q_{x,x'}^{\text{chain}}(\boldsymbol{w}_*)$$

$$= \exp \left\{ \langle \boldsymbol{w}_*^b, \boldsymbol{\phi}^{\text{chain}}(\{x,x'\}) \rangle \right\} \pi_{x'}(\boldsymbol{w}_*^u)$$

Thus, we can set  $\boldsymbol{w}_*^u = \alpha \boldsymbol{w}^u$  for any  $\alpha \in \mathbb{R}$  so that  $\pi_{x'}(\boldsymbol{w}^u) = \pi_{x'}(\boldsymbol{w}_*^u)$ , for  $\forall x' \in \mathcal{X}$ . By plugging the definition of  $\boldsymbol{\phi}^{\text{gtr}}$  and  $\boldsymbol{\phi}^{\text{chain}}$  in Equation 7 and Equation 15 respectively, we require that:

$$\begin{cases} \left(\boldsymbol{w}_{*}^{b}\right)_{\eta(\{x,x'\})} + \left(\boldsymbol{w}_{*}^{b}\right)_{\eta(\{x,x'\})-1} = \left(\boldsymbol{w}^{b}\right)_{\eta(\{x,x'\})}, & \text{for } \eta(\{x,x'\}) = 2,3,\dots,\mathcal{X}|(|\mathcal{X}|-1)/2), \\ \left(\boldsymbol{w}_{*}^{b}\right)_{1} = \left(\boldsymbol{w}^{b}\right)_{1}, & \text{for } \eta(\{x,x'\}) = 1. \end{cases}$$

Set  $p = \mathcal{X}|(|\mathcal{X}| - 1)/2)$ , it is equivalent to solve

$$\begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_1 = \begin{pmatrix} \boldsymbol{w}^b \end{pmatrix}_1$$

$$\begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_1 + \begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_2 = \begin{pmatrix} \boldsymbol{w}^b \end{pmatrix}_2$$

$$\begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_2 + \begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_3 = \begin{pmatrix} \boldsymbol{w}^b \end{pmatrix}_3$$

$$\vdots$$

$$\begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_{p-1} + \begin{pmatrix} \boldsymbol{w}_*^b \end{pmatrix}_p = \begin{pmatrix} \boldsymbol{w}^b \end{pmatrix}_p .$$

We obtain the solution  $\boldsymbol{w}_*^b = \boldsymbol{B} \boldsymbol{w}^b$ , where

$$\boldsymbol{B}_{ij} = \begin{cases} (-1)^{i+j-2}, & \text{if } j \leqslant i \text{ and } i, j = 1, 2, \dots, (|\mathcal{X}|(|\mathcal{X}|-1)/2), \\ 0, & \text{otherwise.} \end{cases}$$

Thus, the solution exists, where  $w_* = \begin{pmatrix} w_*^u \\ w_*^b \end{pmatrix} = \begin{pmatrix} \alpha w^u \\ B w^b \end{pmatrix}$ .

# Appendix B. Gradient Computation for Univariate Weights in the HMC step under combined sampling scheme

In our proposed sampling scheme combining HMC and LBPS described in Algorithm 2, we use HMC to update only the univariate weights  $\mathbf{w}^u$  and LBPS to update the bivariate weights  $\mathbf{w}^b$ , respectively. Since HMC requires the gradient information for  $\mathbf{w}^u$ , we provide the gradient computation for  $\mathbf{w}^u$  in the HMC step of our combined sampling scheme under a reversible, unnormalized rate matrix.

We review that the augmented joint density given the sufficient statistics z(y) := (n(y), h(y), c(y)) for a sample CTMC path y under a reversible, unnormalized rate matrix is:

$$\log f_{\mathbf{w}|\mathbf{z},\mathbf{y}}(\boldsymbol{w}|\boldsymbol{z},\boldsymbol{y}) := -\frac{1}{2}\kappa \|\boldsymbol{w}\|_{2}^{2} - \sum_{x \in \mathcal{X}} h_{x} \sum_{x' \in \mathcal{X}: x \neq x'} q_{x,x'}(\boldsymbol{w}),$$

$$+ \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} c_{x,x'} \log \left(q_{x,x'}(\boldsymbol{w})\right) + \sum_{x \in \mathcal{X}} n_{x} \log(\pi_{x}(\boldsymbol{w})). \quad (19)$$

Under the combined sampling scheme,  $\boldsymbol{w}^u$  and  $\boldsymbol{w}^b$  are updated respectively. In the *i*th iteration of the combined sampling scheme LBPS-HMC, while using HMC kernel, only  $\boldsymbol{w}^u$  is the parameter we need to update. The exchangeable parameters are fixed at the values obtained in the (i-1)th iteration denoted as  $\theta_{\{x,x'\}}\left(\left(\boldsymbol{w}^b\right)^{i-1}\right)$ . To simplify the notation, we denote  $\theta_{\{x,x'\}}\left(\left(\boldsymbol{w}^b\right)^{i-1}\right)$  as  $\theta_{\{x,x'\}}$ . Thus, Equation 19 can be rewritten as:

$$\log f_{\text{w}|\text{z},\text{y}}(\boldsymbol{w}|\boldsymbol{z},\boldsymbol{y}) = -\frac{1}{2}\kappa \|\boldsymbol{w}^u\|_2^2 - \sum_{x \in \mathcal{X}} h_x \sum_{x' \in \mathcal{X}: x \neq x'} \pi_{x'}(\boldsymbol{w}^u) \theta_{\{x,x'\}}$$

$$+ \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} c_{x,x'} \log \left( \pi_{x'}(\boldsymbol{w}^u) \theta_{\{x,x'\}} \right) + \sum_{x \in \mathcal{X}} n_x \log(\pi_x(\boldsymbol{w}^u))$$
(20)

Thus, the gradient of the augmented joint density with respect to  $w^u$  is:

$$\nabla \log f_{\mathbf{w}|\mathbf{z},\mathbf{y}}(\boldsymbol{w}|\boldsymbol{z},\boldsymbol{y}) = -\kappa \boldsymbol{w}^{u} - \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} h_{x} q_{x,x'}(\boldsymbol{w}^{u}) \left( \boldsymbol{\psi}(x') - \sum_{x \in \mathcal{X}} \pi_{x}(\boldsymbol{w}^{u}) \boldsymbol{\psi}(x) \right)$$

$$+ \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} c_{x,x'} \left( \boldsymbol{\psi}(x') - \sum_{x \in \mathcal{X}} \pi_{x}(\boldsymbol{w}^{u}) \boldsymbol{\psi}(x) \right)$$

$$+ \sum_{x \in \mathcal{X}} n_{x} \left( \boldsymbol{\psi}(x) - \sum_{x \in \mathcal{X}} \pi_{x}(\boldsymbol{w}^{u}) \boldsymbol{\psi}(x) \right),$$

$$(21)$$

where:

$$\nabla A(\boldsymbol{w}) = \sum_{x \in \mathcal{X}} \psi(x) \pi_x(\boldsymbol{w}).$$

# Appendix C. Gradient Computation for Univariate and Bivariate Weights using only HMC kernel as benchmark sampling algorithm

To investigate the computational efficiency of our combined sampling scheme using HMC and LBPS, we choose the benchmark sampling algorithm using only HMC kernel to update all the parameters, which is  $\mathbf{w} = (\mathbf{w}^u, \mathbf{w}^b)$ .

Since HMC requires the gradient of w, we derive it in Equation 19:

$$\nabla \log f_{\mathbf{w}|\mathbf{z},\mathbf{y}}(\boldsymbol{w}|\boldsymbol{z},\boldsymbol{y}) = -\kappa \boldsymbol{w} - \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} h_x q_{x,x'}(\boldsymbol{w}) \left( \boldsymbol{\psi}(x') - \sum_{x \in \mathcal{X}} \pi_x(\boldsymbol{w}) \boldsymbol{\psi}(x) \right)$$

$$- \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} h_x q_{x,x'}(\boldsymbol{w}) \boldsymbol{\phi}(\{x, x'\})$$

$$+ \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} c_{x,x'} \left( \boldsymbol{\psi}(x') - \sum_{x \in \mathcal{X}} \pi_x(\boldsymbol{w}) \boldsymbol{\psi}(x) \right)$$

$$+ \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}: x \neq x'} c_{x,x'} \boldsymbol{\phi}(\{x, x'\}) + \sum_{x \in \mathcal{X}} n_x \left( \boldsymbol{\psi}(x) - \sum_{x \in \mathcal{X}} \pi_x(\boldsymbol{w}) \boldsymbol{\psi}(x) \right).$$

# Appendix D. EIT Results

In Table 3 and Table 4, it is shown that our proposed new algorithm has passed the EIT since all the pvalues are bigger than a threshold of 0.05/n, where n represents the number of parameters in the test while taking multiple comparisons into account. For example, n=10 while testing the exchangeable parameters. Similarly, we display the test results for using HMC in Table 5 and Table 6.

Parameter index	1	2	3	4	5
Pvalue	0.640	0.200	0.329	0.114	0.167
Test statistics	0.060	0.087	0.077	0.097	0.090
Test	KS	KS	KS	KS	KS

Table 3: EIT for  $\mathbf{w}^u$  using LBPS-HMC, which uses HMC for the stationary distribution weights and LBPS for  $\mathbf{w}^b$  for exchangeable parameters.

Parameter index	1	2	3	4	5	6	7	8	9	10
Pvalue	0.200	0.504	0.441	0.061	0.571	0.838	0.382	0.709	0.238	0.382
Test statistics	0.087	0.067	0.070	0.107	0.063	0.050	0.073	0.057	0.083	0.073
Test	KS									

Table 4: EIT for exchangeable parameters using LBPS-HMC, which uses HMC for the stationary distribution and LBPS for the weights of the exchangeable parameters.

Parameter index	1	2	3	4	5
Pvalue	0.967	0.200	0.640	0.838	0.640
Test statistics	0.040	0.087	0.060	0.050	0.060
Test	KS	KS	KS	KS	KS

Table 5: EIT for stationary distribution elements using HMC for both the stationary distribution and the weights of the exchangeable parameters.

Parameter index	1	2	3	4	5	6	7	8	9	10
Pvalue	0.139	0.076	0.504	0.281	0.838	0.640	0.936	0.504	0.139	0.238
Test statistics	0.093	0.103	0.067	0.080	0.050	0.060	0.043	0.067	0.093	0.083
Test	KS									

Table 6: EIT for exchangeable parameters using HMC for both the stationary distribution and the weights of the exchangeable parameters.

# Appendix E. Posterior samples comparison between LBPS and HMC

In this section, we provide the summary statistics including the minimum, first quantile, median, mean, third quantile and maximum of the ARD which is defined in Equation 22, across all exchangeable parameters between samples obtained from one single run of LBPS-HMC and HMC in Table 7 under the 8-by-8 Bayesian GLM chain GTR model described in Section 3. We also provide the boxplot comparison of the posterior samples collected from LBPS and HMC under the same set-up in Figure 9.

Min.	1st Quantile	Median	Mean	3rd Quantile	Max.
1.4%	1.16%	2.70%	3.05%	3.63%	12.06%

Table 7: Summary of ARD across exchangeable parameters between HMC and LBPS.

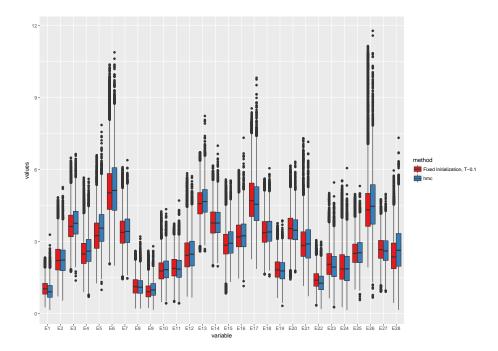


Figure 9: Boxplot of posterior samples comparison for exchangeable parameters of an 8-by-8 rate matrix between LBPS and HMC.

Appendix F. Euclidean Distance between Pairs of Amino Acids

1.	l.																			_ 1
		115	, ,	, ,	, ,	, ,	, ,	, ,	, ,	' '			, ,	, ,	, ,				215	0
C	194	174	154	180	149	139	202	154	170	159	205	198	195	112	169	198	196	192	0	215
>	55	84	96	96	69	133	97	152	121	109	20	32	64	124	89	29	21	0	192	88
M	36	87	101	91	81	142	95	160	126	127	28	15	84	135	87	10	0	21	196	29
I	33	94	109	97	88	149	102	168	134	135	21	ಬ	94	142	95	0	10	29	198	61
Ь	110	22	92	103	38	91	103	108	93	42	114	86	27	74	0	95	87	89	169	147
$\infty$	144	88	89	110	28	46	121	65	80	26	155	145	66	0	74	142	135	124	112	177
A	112	98	91	112	228	1111	106	126	107	09	113	96	0	66	27	94	84	64	195	148
П	36	66	113	102	92	153	107	172	138	138	22	0	96	145	86	ಬ	15	32	198	61
দ্র	22	100	116	26	103	158	102	177	140	153	0	22	113	155	114	21	28	20	205	40
IJ	147	86	87	125	59	80	127	94	86	0	153	138	09	26	42	135	127	109	159	184
田	122	40	29	54	65	42	26	45	0	86	140	138	107	80	93	134	126	121	170	152
D	160	81	61	96	85	23	101	0	45	94	177	172	126	65	108	168	160	152	154	181
K	85	32	53	26	78	94	0	101	26	127	102	107	106	121	103	102	92	26	202	110
Z	143	89	46	98	65	0	94	23	42	80	158	153	111	46	91	149	142	133	139	174
L	92	47	42	71	0	65	28	85	65	59	103	95	28	28	38	88	81	69	149	128
$\mathbb{R}$	22	53	43	0	71	98	26	96	54	125	26	102	112	110	103	26	91	96	180	102
O	66	24	0	43	42	46	53	61	29	87	116	113	91	89	92	109	101	96	154	130
Η	83	0	24	29	47	89	32	81	40	86	100	66	98	88	22	94	87	84	174	115
X	0	83	66	22	92	143	82	160	122	147	22	36	112	144	110	33	36	55	194	37
	Y	Н	o	$\mathbb{R}$	Н	Z	K	Ω	闰	IJ	ſΞı	Ц	A	$\infty$	Ъ	П	$\mathbb{Z}$	>	C	M

# Appendix G. Nearest Neighbour Pariwise Amino Acid Ordering

#### **Algorithm 3** Nearest neighbour pairwise amino acid ordering

```
1: Initialization:
     \mathcal{X}^{\text{dynamic, distinct}} = \mathcal{X}^{\text{distinct}}, counter=0.
     AminoAcidPairRank=DynamicSupportForAminoAcid =dict(), which is an empty dictionary.
     Set \alpha = \beta = i_0 = i_1 = \text{NULL}.
 2: for i \in \mathcal{X} do
        DynamicSupportForAminoAcid[i] = \mathcal{X} \setminus \{i\}.
 4: end for
 5: while \mathcal{X}^{\text{dynamic, distinct}} is not empty do
        if i_0 is not NULL and i_1 is not NULL then
 6:
           Find support of i_0: \alpha=DynamicSupportForAminoAcid[i_0].
 7:
           Find support of i_1: \beta=DynamicSupportForAminoAcid[i_1].
 8:
 9:
        end if
        if neither of \alpha or \beta is empty then
10:
11:
           rowDist = AminoAcidDist[i_0, \alpha]
12:
           colDist = AminoAcidDist[\beta, i_1]
           rowMinInd = \alpha[argmin(rowDist)]
13:
           colMinInd = \beta[argmin(colDist)]
14:
           if AminoAcidDist[i_0, rowMinInd] \leq AminoAcidDist[colMinInd, i_1] then
15:
16:
              i_1 = \text{rowMinInd}
17:
           else if AminoAcidDist[i_0, rowMinInd] > AminoAcidDist[colMinInd, i_1] then
18:
              i_0 = \text{colMinInd}
19:
           end if
        else if \alpha is NULL and \beta is not NULL then
20:
           colDist=AminoAcidDist[\beta, i_1]
21:
           \operatorname{colMinInd} = \beta [\operatorname{argmin}(\operatorname{colDist})]
22:
23:
           i_0 = \text{colMinInd}
        else if \alpha is not NULL and \beta is NULL then
24:
25:
           rowDist = AminoAcidDist[i_0, \alpha]
26:
           rowMinInd = \alpha [argmin(rowDist)]
27:
           i_1 = \text{rowMinInd}
28:
        else
           Find the amino acid pair SmallestPair \in \mathcal{X}^{dynamic,distinct} with nonzero minimum in
29:
           AminoAcidDist.
30:
           i_0= index of the first amino acid in SmallestPair in \mathcal{X}.
           i_1 = \text{index of the second amino acid in SmallestPair in } \mathcal{X}.
31:
32:
33:
        AminoAcidPairRank[i_0, i_1] = AminoAcidPairRank[i_1, i_0] = counter, counter ++.
34:
        AminoAcidDist[i_0, i_1]=AminoAcidDist[i_1, i_0]=\infty.
35:
        if i_1 \in DynamicSupportForAminoAcid[i_0] then
           Remove i_1 from DynamicSupportForAminoAcid[i_0]
36:
        end if
37:
        if i_0 \in DynamicSupportForAminoAcid[i_1] then
38:
39:
           Remove i_0 from DynamicSupportForAminoAcid[i_1]
40:
        end if
       s_0 = \mathcal{X}_{i_0} + \mathcal{X}_{i_1} (obtain amino acid pair s_0), s_1 = \mathcal{X}_{i_1} + \mathcal{X}_{i_0} (obtain amino acid pair s_1). if s_0 \in \mathcal{X}^{\text{dynamic, ordered, distinct}} then
41:
42:
          Remove s_0 from \mathcal{X}^{\text{dynamic, distinct}}
43:
44:
        if s_1 \in \mathcal{X}^{\text{dynamic, distinct}} then
45:
           Remove s_1 from \mathcal{X}^{\text{dynamic, distinct}}
                                                           34
46:
        end if
47:
48: end while
49: Return AminoAcidPairRank.
```

# Appendix H. Traceplot comparison for real data analysis

We provide the traceplot on a randomly selected exchangeable parameter using LBPS compared with HMC. There is a small difference between posterior means with ARD 0.43%.

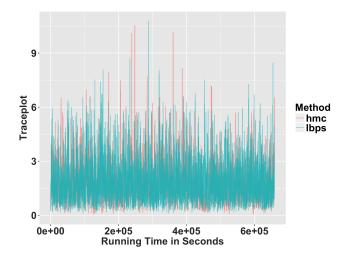


Figure 10: Traceplot for the 145th exchangeable parameter using LBPS-HMC compared with HMC with ARD 0.43% in terms of the posterior mean.

# Appendix I. Correctness check via ARD between two samplers

Theoretically, we have the same prior distribution and likelihood for the parameters (weights) of interest, LBPS-HMC and HMC share the same posterior distribution. Thus, if we have longer chains for both algorithms, their density of the posterior samples should be closer and closer. We evaluate this via ARD:

$$ARD(x,y) = \frac{|x-y|}{\max(x,y)}, \text{ for } \forall x,y>0.$$
 (22)

Figure 12 and Table 8 demonstrate the ARD across all exchangeable parameters between HMC and LBPS with trajectory length 0.2.

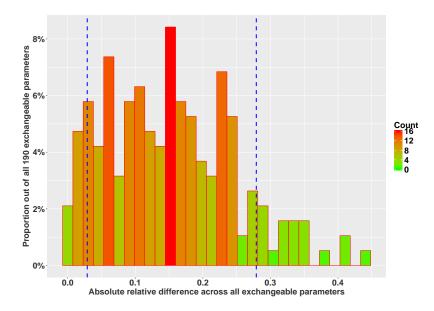


Figure 11: ARD across exchangeable parameters between LBPS compared with HMC with blue dotted lines representing the 10% and 90% quantile of ARD.

Min.	1st Quantile	Median	Mean	3rd Quantile	Max.
0.3%	7.5%	14.8%	15.1%	22.0%	44.4%

Table 8: Summary of ARD across exchangeable parameters.

Since the maximum value of ARD is 44.4%, so we further check whether the value will decrease as expected when we have a longer chain. We have found that the difference between the two algorithms in terms of ARD becomes smaller and smaller when we run both algorithms for a longer period. This is also the case when conducting our simulation studies. We validate this by providing a violin plot in Appendix G comparing the ARD from

the current run with the results from a longer run. We denoted the ARD from the current run as "shorter" and the ARD from a longer run as "longer", where HMC had a walltime of 12 days and LBPS had a walltime of 10 days. Each method had a walltime around two days longer than the current run. It has been validated that when we have longer chains, the ARD values will decrease, indicating smaller difference between the two methods. The maximum of ARD dropped from 44% to 37% when we have a longer walltime.

In order to figure out the situations when ARD values are larger for certain parameters, we demonstrate the quantile of ESS in Table 9 and the actual ESS in Table 10 of the exchangeable parameters that have ARD bigger than its 95% quantile between LBPS and HMC. We have found that a relatively bigger difference in the posterior mean between the two methods happens in situations when either they both have low ESS for this parameter compared to other parameters or when LBPS has high ESS such as 97.9% quantile and 94.2% quantile across exchangeable parameters but HMC does not obtain an equivalently high ESS for the same parameter within the same amount or even slightly longer walltime.

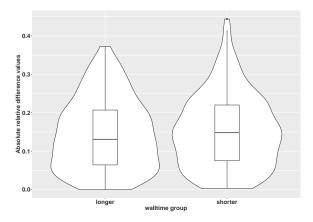


Figure 12: Violin plot of ARD across exchangeable parameters between a shorter run and a longer run, which has 2 days longer walltime.

LBPS	2.63%	2.105%	8.42%	5.79%	77.4%	94.2%	75.8%	97.9~%	0.526%	3.16%
HMC	16.84%	0.526%	15.79%	2.63%	53.7%	84.2%	73.7%	98.4%	1.053%	5.26%

Table 9: The quantiles of ESS of exchangeable parameters that have ARD bigger than its 95% quantile.

LBPS	1071	1051	1248	1147	2145	2330	2115	2451	773	1118
HMC	544	152	604	488	927	1032	1090	591	190	718

Table 10: The ESS of exchangeable parameters that have ARD bigger than 95% quantile of ARD.

#### References

- Christophe Andrieu, Alain Durmus, Nikolas Nüsken, and Julien Roussel. Hypercoercivity of piecewise deterministic markov process-monte carlo. arXiv preprint arXiv:1808.08592, 2018.
- Maria Anisimova and Carolin Kosiol. Investigating protein-coding sequence evolution with probabilistic codon substitution models. *Molecular biology and evolution*, 26(2):255–271, 2008.
- Joris Bierkens. Non-reversible metropolis-hastings. *Statistics and Computing*, 26(6):1213–1228, 2016.
- Alexandre Bouchard-Côté, Sebastian J Vollmer, and Arnaud Doucet. The bouncy particle sampler: A nonreversible rejection-free markov chain monte carlo method. *Journal of the American Statistical Association*, pages 1–13, 2018.
- Ting-Li Chen and Chii-Ruey Hwang. Accelerating reversible markov chains. Statistics & Probability Letters, 83(9):1956–1962, 2013.
- Mark HA Davis. Piecewise-deterministic markov processes: A general class of non-diffusion stochastic models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 353–388, 1984.
- George Deligiannidis, Daniel Paulin, and Arnaud Doucet. Randomized hamiltonian monte carlo as scaling limit of the bouncy particle sampler and dimension-free convergence rates. arXiv preprint arXiv:1808.04299, 2018.
- Paul Fearnhead, Joris Bierkens, Murray Pollock, and Gareth O Roberts. Piecewise deterministic markov processes for continuous-time monte carlo. arXiv preprint arXiv:1611.07873, 2016.
- James M Flegal, J Hughes, D Vats, et al. mcmcse: Monte carlo standard errors for mcmc. Riverside, CA and Minneapolis, MN. R package version, pages 1–0, 2012.
- John Geweke. Getting it right: Joint distribution tests of posterior simulators. *Journal of the American Statistical Association*, 99(467):799–804, 2004.
- Rao Grantham. Amino acid difference formula to help explain protein evolution. *Science*, 185(4154):862–864, 1974.

- Peter Guttorp and Vladimir N Minin. Stochastic modeling of scientific data. CRC Press, 2018.
- M. Hasegawa, H. Kishino, and T. Yano. Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. *Journal of Molecular Eevolution*, 22:160–174, 1985.
- Matthew X He, Miguel A Jiménez-Montaño, and Paolo E Ricci. Amino acids, euclidean distance and symmetric matrix. In *Proceedings of the International Conference on Bioinformatics & Computational Biology (BIOCOMP)*, page 1. The Steering Committee of The World Congress in Computer Science, Computer Engineering and Applied Computing (WorldComp), 2011.
- Matthew D Hoffman and Andrew Gelman. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *Journal of Machine Learning Research*, 15(1): 1593–1623, 2014.
- Hemant Ishwaran, J Sunil Rao, et al. Spike and slab variable selection: frequentist and bayesian strategies. *The Annals of Statistics*, 33(2):730–773, 2005.
- Christopher H Jackson, Linda D Sharples, Simon G Thompson, Stephen W Duffy, and Elisabeth Couto. Multistate markov models for disease progression with classification error. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 52(2):193–209, 2003.
- T.H. Jukes and C.R. Cantor. *Evolution of Protein Molecules*. New York: Academic Press., 1969.
- JD Kalbfleisch and Jerald F Lawless. The analysis of panel data under a markov assumption. Journal of the American Statistical Association, 80(392):863–871, 1985.
- B. R. Kay. Proportional hazard regression models and the analysis of censored survival data. *Applied Statistics*, 26(3):227–237, 1977.
- M. Kimura. A simple method for estimating evolutionary rate of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution*, 16:111–120, 1980.
- Cleve Moler and Charles Van Loan. Nineteen dubious ways to compute the exponential of a matrix, twenty-five years later. SIAM review, 45(1):3–49, 2003.
- Radford M Neal. Improving asymptotic variance of mcmc estimators: Non-reversible chains are better.  $arXiv\ preprint\ math/0407281$ , 2004.
- Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, 2(11):2, 2011.
- James R Norris. Markov chains. Number 2. Cambridge university press, 1998.
- Ari Pakman. Binary bouncy particle sampler. arXiv preprint arXiv:1711.00922, 2017.

- Ari Pakman, Dar Gilboa, David Carlson, and Liam Paninski. Stochastic bouncy particle sampler. In *International Conference on Machine Learning*, pages 2741–2750, 2017.
- Elias AJF Peters et al. Rejection-free monte carlo sampling for general potentials. *Physical Review E*, 85(2):026703, 2012.
- Vinayak Rao and Yee Whye Teh. Fast mcmc sampling for markov jump processes and extensions. The Journal of Machine Learning Research, 14(1):3295–3320, 2013.
- Yi Sun, Jürgen Schmidhuber, and Faustino J Gomez. Improving the asymptotic performance of markov chain monte-carlo by inserting vortices. In *Advances in Neural Information Processing Systems*, pages 2235–2243, 2010.
- Paula Tataru and Asger Hobolth. Comparison of methods for calculating conditional expectations of sufficient statistics for continuous time Markov chains. *BMC Bioinformatics*, 12:465–475, 2011.
- Paul Vanetti, Alexandre Bouchard-Côté, George Deligiannidis, and Arnaud Doucet. Piecewise deterministic markov chain monte carlo. arXiv preprint arXiv:1707.05296, 2017.
- Ziyu Wang, Shakir Mohamed, and Nando Freitas. Adaptive hamiltonian and riemann manifold monte carlo. In *International Conference on Machine Learning*, pages 1462–1470, 2013.
- G George Yin and Qing Zhang. Continuous-time Markov chains and applications: A two-time-scale approach, volume 37. Springer Science & Business Media, 2012.
- Tingting Zhao, Ziyu Wang, Alexander Cumberworth, Joerg Gsponer, Nando de Freitas, Alexandre Bouchard-Côté, et al. Bayesian analysis of continuous time markov chains with application to phylogenetic modelling. *Bayesian Analysis*, 11(4):1203–1237, 2016.