

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

## Bayesian regularization of the length of memory in reversible sequences

### This is the author's manuscript

*Original Citation:*

*Availability:*

This version is available <http://hdl.handle.net/2318/1591251> since 2016-11-13T15:50:07Z

*Published version:*

DOI:10.1111/rssb.12140

*Terms of use:*

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)



*J. R. Statist. Soc. B* (2016)  
**78**, Part 4, pp. 933–946

# Bayesian regularization of the length of memory in reversible sequences

Sergio Bacallado and Vijay Pande,  
*Stanford University, USA*

Stefano Favaro  
*University of Turin and Collegio Carlo Alberto, Turin, Italy*

and Lorenzo Trippa  
*Dana–Farber Cancer Institute, Boston, USA*

[Received May 2014. Final revision August 2015]

**Summary.** Variable order Markov chains have been used to model discrete sequential data in a variety of fields. A host of methods exist to estimate the history-dependent lengths of memory which characterize these models and to predict new sequences. In several applications, the data-generating mechanism is known to be reversible, but combining this information with the procedures mentioned is far from trivial. We introduce a Bayesian analysis for reversible dynamics, which takes into account uncertainty in the lengths of memory. The model proposed is applied to the analysis of molecular dynamics simulations and compared with several popular algorithms.

**Keywords:** Bayesian analysis; Reinforced random walk; Reversibility; Variable order Markov model

## 1. Introduction

Time reversibility characterizes numerous stochastic models, from queuing networks (Kelly, 1979) to models of physical systems governed by reversible mechanics (Van Kampen, 1992). This property is commonly associated with Markov chains. However, in several applications the Markovian assumption is restrictive, and one needs to model higher order dependences. Higher order Markov models require the investigator to deal with the problem of a rapidly increasing number of unknown parameters. Here, we propose a Bayesian approach to this problem, tailored to reversible processes.

Our motivating application is the analysis of molecular dynamics simulations. These are computer experiments which mimic the structural transitions of macromolecules in time by using a physics-based Hamiltonian operator. The simulation is a reversible Markov chain representing the trajectory of hundreds of atoms in a three-dimensional space. This high dimensional process is typically projected onto a partition of the space of molecular structures, to produce a discrete time series which is still reversible, even though it is not necessarily a Markov chain. An accurate characterization of the projected dynamics provides the biologist with estimates of the rates of

*Address for correspondence:* Sergio Bacallado, Department of Statistics, Stanford University, 390 Serra Mall, Stanford, CA 94305, USA.  
E-mail: sergiobacallado@gmail.com

transitions between biologically relevant states. Such estimates, with associated uncertainty measures, are also useful in the design of adaptive simulations which can significantly reduce the computational burden of these experiments (Prinz *et al.*, 2011).

In variable order Markov models, the length of memory is a function of the previously observed states. If all possible sequences of states are arranged in a context tree, where every branch is truncated, then it is necessary to specify transition probabilities only at the nodes where truncations occur. The length of each branch is a context-specific length of memory. The literature on variable order Markov chains can be traced back to Rissanen (1983) and Weinberger *et al.* (1995), who developed algorithms based on context tree pruning. Statistical properties of these algorithms were developed by Bühlmann and Wyner (1999) and Bühlmann (2000). Begleiter *et al.* (2004) reviewed relevant algorithms, such as context tree weighting which is similar in spirit to Bayesian model averaging methods.

Bayesian non-parametric priors have been recently proposed for higher order Markov chains. In the hierarchical Dirichlet language model (MacKay and Peto, 1995), there are transition distributions  $G_u$  out of every point  $u$  in  $\cup_{i=1}^n \mathcal{X}^i$ , where  $\mathcal{X}$  is the state space and  $n$  is the length of memory. If, for instance,  $\mathcal{X}$  corresponds to the alphabet, then the prior for  $G_{\text{word}}$  is a Dirichlet distribution with mean identical to the ancestor  $G_{\text{ord}}$ . Teh (2006), Wood *et al.* (2009) and Mochihashi and Sumita (2008) developed successful non-parametric extensions of this construction for language modelling.

However, it would be difficult to combine the hierarchical structure of these Bayesian models with our goal of doing inference for a reversible process. Annis *et al.* (2010) have shown that, when the data-generating process is known to be reversible, models that enforce reversibility can have superior asymptotic properties. To develop a procedure that incorporates reversibility, we introduce the *random walk with amnesia*. This process generalizes the variable order Markov model and is related to *probabilistic suffix automata*, which is a construction that was introduced by Ron *et al.* (1996). Relying on the conjugate prior for reversible variable order Markov chains that was introduced by Bacallado (2011), we define a Bayesian analysis for the random walk with amnesia. Our main contribution is an efficient procedure for Bayesian inference of reversible dynamics when the context-specific lengths of memory are unknown.

We note that an alternative, and widely used, modelling approach for sequential data is offered by hidden Markov models and extensions such as tiered hidden Markov models; see Cappé *et al.* (2005) for a comprehensive treatment. These models introduce memory through a latent process instead of explicitly modelling the dependence on history, like the methods that are considered in this paper. Reversible versions have been applied recently (Palla *et al.*, 2014).

The data that are analysed in the paper and the programs that were used to analyse them can be obtained from

<http://wileyonlinelibrary.com/journal/rss-datasets>

## 2. Reversible dynamics

### 2.1. Random walk with amnesia

We define a random walk with finite memory taking values in a finite space  $\mathcal{X}$ . Our random walker remembers only the last portion of his trajectory, and the length of this vector of states changes in time. At every step, the walker can either

- (a) lose the first element of his memory or
- (b) proceed to a new state in  $\mathcal{X}$ .

For example, when  $\mathcal{X} = \{A, B, C\}$  a path could be

$$ABC \xrightarrow{(a)} BC \xrightarrow{(a)} C \xrightarrow{(b)} CA \xrightarrow{(b)} CAB, \tag{1}$$

where  $\xrightarrow{(a)}$  denotes a transition of the first kind and  $\xrightarrow{(b)}$  a transition of the second kind.

We use an  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  on the enriched space  $\mathcal{X}_s = \mathcal{X} \cup \{s_k; k = 1, \dots, n\}$  to represent this process. The maximum number of  $\mathcal{X}$ -valued states that the random walker can remember is  $n - 1$ . If, for some  $k > 0$ ,  $x_1 = \dots = x_k = s_k$  and  $(x_{k+1}, \dots, x_n) \in \mathcal{X}^{n-k}$  we say that  $(x_1, \dots, x_n)$  belongs to  $\mathcal{Z} \subset \mathcal{X}_s^n$ . This  $n$ -gram represents a state in which the random walker remembers only the sequence  $(x_{k+1}, x_{k+2}, \dots, x_n) \in \mathcal{X}^{n-k}$ ; the first  $k$  states  $x_1, \dots, x_k$  are equal to  $s_k$ , denoting a loss of memory. The sequence of  $n$ -grams belonging to  $\mathcal{Z}$  in the process  $(X_i)_{i \in \mathbb{N}}$  will mirror a trajectory like that in path (1). Definition 1 puts restrictions on the  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  which ensure that it represents the process that is described informally above. In the  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  the only random transitions occur from  $n$ -grams  $(X_j, \dots, X_{j+n-1})$  in  $\mathcal{Z}$ , whereas the remaining transitions are deterministic. The random transitions determine whether

- (a) the random walker forgets the first state in her memory or
- (b) adds a new  $\mathcal{X}$ -state to his trajectory.

After each random transition, a series of deterministic transitions will bring the process to the subsequent  $n$ -gram in  $\mathcal{Z}$ .

*Definition 1.* A random walk with amnesia is an  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  on the space  $\mathcal{X}_s$ . With probability 1,  $(X_1, \dots, X_n) \in \mathcal{Z}$ . Given  $(X_j, \dots, X_{j+n-1}) = (x_1, \dots, x_n) = x$ , for any  $j > 0$ , the next state  $X_{j+n}$  satisfies the following constraints. If  $x \in \mathcal{Z}$  and  $x_1 = \dots = x_k = s_k$ , then  $X_{j+n}$  can be

- (a)  $s_{k+1}$ , provided that  $k < n$ , or
- (b) a state in  $\mathcal{X}$ , provided that  $k > 1$ . If  $x \notin \mathcal{Z}$ , the state  $X_{j+n}$  is chosen deterministically according to the following rules:
  - (i) if  $x_1 = \dots = x_{k-1} = s_k$  and  $(x_k, \dots, x_n) \in \mathcal{X}^{n-k+1}$  for some  $1 < k \leq n$ , then  $X_{j+n} = s_{k-1}$ ;
  - (ii) if  $x_{n-m} = \dots = x_n = s_k$  and  $x_{n-m-1} \neq s_k$  for some  $m \leq k - 2 \leq n - 2$ , then  $X_{j+n} = s_k$ ;
  - (iii) if  $x_{m+1} = \dots = x_{m+k} = s_k$  and  $x_m \neq s_k$  for some  $0 < k < k + m \leq n$ , then  $X_{j+n} = x_1$ .

If we consider the path (1) and set  $n = 4$ , then, by definition 1, the corresponding transitions of the  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  on the enriched space  $\mathcal{X}_s$  are

$$\mathbf{s_1ABC} \xrightarrow{(a)} ABCs_2 \rightarrow BCs_2s_2 \rightarrow Cs_2s_2B \rightarrow \mathbf{s_2s_2BC} \xrightarrow{(a)} s_2BCs_3 \rightarrow BCs_3s_3 \rightarrow Cs_3s_3s_3 \rightarrow \mathbf{s_3s_3s_3C} \xrightarrow{(b)} s_3s_3CA \rightarrow s_3CAs_2 \rightarrow CAs_2s_2 \rightarrow As_2s_2C \rightarrow \mathbf{s_2s_2CA} \xrightarrow{(b)} s_2CAB \rightarrow CABS_1 \rightarrow ABS_1C \rightarrow BS_1CA \rightarrow \mathbf{s_1CAB},$$

where  $\xrightarrow{(a)}$  and  $\xrightarrow{(b)}$  denote random transitions from  $n$ -grams in  $\mathcal{Z}$ , displayed in bold, and  $\rightarrow$  denotes deterministic transitions. In this example  $(X_1, X_2, \dots) = (s_1, A, B, C, s_2, s_2, B, C, s_3, s_3, \dots)$ .

One can easily verify that the process  $(X_i)_{i \in \mathbb{N}}$  goes through only  $n$ -grams covered by definition 1. The law of the process is specified by the distribution of  $(X_1, \dots, X_n) \in \mathcal{Z}$  and the transition probabilities out of the  $n$ -grams in  $\mathcal{Z}$ .

Clearly, the definition of  $(X_i)_{i \in \mathbb{N}}$  does not allow every possible path. However, if  $(x_1, \dots, x_r) \in \mathcal{X}_s^r, r > n$ , is a realization of the process with  $(x_r, \dots, x_{r-n+1}) \in \mathcal{Z}$ , then its inverse  $(x_r, x_{r-1}, \dots, x_1)$  is also an allowed path. In other words, if  $(x_1, \dots, x_r)$  is consistent with the deterministic constraints (i)–(iii) in definition 1 and  $(x_r, \dots, x_{r-n+1}) \in \mathcal{Z}$ , then  $(x_r, x_{r-1}, \dots, x_1)$  is also consistent with these three requirements. For instance, with  $n = 3$  and  $\mathcal{X} = \{A, B, C\}$ , the path

$(x_1, \dots, x_r) = (s_1, A, B, s_2, s_2, B, C, s_1)$  corresponds to a random walker in  $B$  who forgets a previous visit to  $A$  and then goes to  $C$ , whereas the reverse sequence  $(s_1, C, B, s_2, s_2, B, A, s_1)$  corresponds to a random walker in  $B$  who forgets a previous visit to  $C$  and then goes to  $A$ . Every loss of memory in the forward path  $(x_1, \dots, x_r)$  is associated with a specular transition in  $(x_r, \dots, x_1)$  to an  $\mathcal{X}$ -valued state and vice versa. This will be a key property to introduce reversible random walks with amnesia in what follows. It will also be useful in Section 3 to describe our Bayesian approach to infer the transition probabilities of  $(X_i)_{i \in \mathbb{N}}$ , under the assumption of reversibility, as a model-based reinforcement learning procedure.

An irreducible  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  on  $\mathcal{X}_s$  is canonically represented by a *balanced* function  $w_{n+1} : \mathcal{X}_s^{n+1} \rightarrow [0, \infty)$ , which satisfies, for any  $(x_1, \dots, x_n) \in \mathcal{X}_s$ ,

$$\sum_{v \in \mathcal{X}_s} w_{n+1}(x_1, \dots, x_n, v) = \sum_{v \in \mathcal{X}_s} w_{n+1}(v, x_1, \dots, x_n) = w_n(x_1, \dots, x_n) \tag{2}$$

and

$$\text{pr}(X_{n+m+1} | X_1, \dots, X_{n+m}) = \frac{w_{n+1}(X_{1+m}, \dots, X_{n+m+1})}{w_n(X_{m+1}, \dots, X_{m+n})}. \tag{3}$$

Note that  $w_n$  is a stationary measure of this  $n$ th-order Markov chain. Without loss of generality, in what follows, we shall assume that  $\sum_x w_{n+1}(x) = 1$ .

The  $n$ th-order Markov chain  $(X_i)_{i \in \mathbb{N}}$  is reversible if and only if

$$w_{n+1}(x_1, \dots, x_{n+1}) = w_{n+1}(x_{n+1}, \dots, x_1) \quad \text{for every } (x_1, \dots, x_{n+1}) \in \mathcal{X}_s^{n+1}. \tag{4}$$

The following approach is a simple approach to define a function  $w_{n+1}$  that is consistent with this notion of reversibility. The resulting Markov chain is not necessarily a random walk with amnesia. Consider a cyclic sequence  $(x_1, \dots, x_r) \in \mathcal{X}_s^r$ , with  $r > n$  and  $(x_1, \dots, x_n) = (x_{r-n+1}, \dots, x_r)$ , and define

$$w_{n+1}(y_1, \dots, y_{n+1}) = \frac{1}{2(r-n)} \sum_{i=1}^{r-n} \mathbf{1}\{(x_i, \dots, x_{i+n}) = (y_1, \dots, y_{n+1})\} + \mathbf{1}\{(x_i, \dots, x_{i+n}) = (y_{n+1}, y_n, \dots, y_1)\}$$

for every  $n$ -gram  $(y_1, \dots, y_{n+1}) \in \mathcal{X}_s^{n+1}$ , where  $\mathbf{1}(\cdot)$  denotes an indicator function. The resulting function  $w_{n+1}$  satisfies equation (2) and allows us to specify the transition probabilities of a reversible  $n$ th-order Markov chain through the identity (3). This construction can be used to prove the following result.

*Proposition 1.* There exist random walks with amnesia with a reversible parameter  $w_{n+1}$ .

*Proof.* The random walk with amnesia is a class of  $n$ th-order Markov chains on  $\mathcal{X}_s$  where certain transitions are not allowed. Therefore, it is sufficient to follow the constructive approach that we described in the previous paragraph to specify a reversible  $w_{n+1}$ . In this case we define  $w_{n+1}$  through a cyclic sequence  $(x_1, \dots, x_r)$ , with  $(x_1, \dots, x_n) = (x_{r-n+1}, \dots, x_r)$ , which is in addition consistent with the constraints in definition 1. The remark that, if  $(x_1, \dots, x_r)$  is consistent with the constraints in definition 1, then so is the reversed sequence  $(x_r, \dots, x_1)$ , completes the proof.

We note that if  $w_{n+1} = \lambda w'_{n+1} + (1 - \lambda)w''_{n+1}$ ,  $\lambda \in (0, 1)$ , where both  $w'_{n+1}$  and  $w''_{n+1}$  are reversible parameters for random walks with amnesia, then  $w_{n+1}$  satisfies equations (2) and (4). This fact implies that reversible parameterizations for random walks with amnesia  $w_{n+1}$  constitute a convex space.

2.2. *Observable random sequence*

We use random walks with amnesia to model sequences with higher order dependences. Recall that the state space of the amnesia model is  $\mathcal{X}_s$ . We assume that we observe only the sequence  $(Z_j)_{j \in \mathbb{N}}$  of  $\mathcal{X}$ -valued states visited by our random walker. To make the definition precise, let  $\tau_0 = n$  and  $\tau_j = \min\{\tau; \tau > \tau_{j-1}, (X_{\tau-n}, \dots, X_{\tau-1}) \in \mathcal{Z}, X_\tau \in \mathcal{X}\}$ . The observable process is  $Z_j := X_{\tau_j}, j \geq 1$ . In the example displayed in path (1)  $Z_1 = A$  and  $Z_2 = B$ . The observable process does not need to be a Markov chain of any finite order, even though, as we show in what follows, reversible variable order Markov chains are a special case.

*Proposition 2.* Suppose that, given  $X_1, \dots, X_{n+1}$ , the process  $(X_i)_{i \in \mathbb{N}}$  is a reversible and irreducible random walk with amnesia with parameter  $w_{n+1}$ , and let

$$\text{pr}(X_1 = x_1, \dots, X_{n+1} = x_{n+1}) \propto w_{n+1}(x_1, \dots, x_{n+1}) \mathbf{1}\{(x_1, \dots, x_n) \in \mathcal{Z}, x_{n+1} \in \mathcal{X}\}. \tag{5}$$

Then the observable process  $(Z_i)_{i \in \mathbb{N}}$  is stationary and reversible, i.e.

$$\text{pr}(Z_1 = z_1, \dots, Z_m = z_m) = \text{pr}(Z_1 = z_m, \dots, Z_m = z_1), \tag{6}$$

for any  $m > 1$  and  $z = (z_1, \dots, z_m) \in \mathcal{X}^m$ .

In what follows, we neglect condition (5), as we are interested in estimating the transition probabilities in the latent process  $(X_i)_{i \in \mathbb{N}}$  from observations which are not stationary. Nonetheless, assumption (5) is unnecessary to verify the following notion of reversibility.

*Proposition 3.* Let  $(X_i)_{i \in \mathbb{N}}$  be a reversible and irreducible random walk with amnesia. For any  $m > 1$  and  $z = (z_1, \dots, z_m) \in \mathcal{X}^m$ , with probability 1,

$$\lim_{k \rightarrow \infty} \sum_{i=1}^k \frac{\mathbf{1}(Z_i = z_1, \dots, Z_{i+m-1} = z_m)}{k} = \lim_{k \rightarrow \infty} \sum_{i=1}^k \frac{\mathbf{1}(Z_i = z_m, \dots, Z_{i+m-1} = z_1)}{k}. \tag{7}$$

The proofs of propositions 2 and 3 are in the on-line supplementary materials. We conclude this section by showing that reversible variable order Markov chains are included within the family of observable processes  $(Z_i)_{i \in \mathbb{N}}$  defined above. A variable order Markov chain  $(U_i \in \mathcal{X})_{i \in \mathbb{N}}$  with maximum order  $n - 1$  and histories  $\mathcal{H} \subset \cup_{i=1}^{n-1} \mathcal{X}^i$  satisfies the equality

$$\text{pr}(U_j = x_j | U_1 = x_1, \dots, U_{j-1} = x_{j-1}) = \text{pr}(U_j = x_j | U_{j-k} = x_{j-k}, \dots, U_{j-1} = x_{j-1})$$

for every  $(x_1, \dots, x_j)$  whenever  $(x_{j-k}, \dots, x_{j-1}) \in \mathcal{H}$ . We shall assume without loss of generality that  $(U_1, \dots, U_{n-1})$  are fixed. We also assume that  $(U_i)_{i \in \mathbb{N}}$  is reversible.

*Proposition 4.* The sequence  $(U_i)_{i \geq n}$  is identical in distribution to a sequence  $(Z_i)_{i \in \mathbb{N}}$  of  $\mathcal{X}$ -valued states from a random walk with amnesia whose parameter  $w_{n+1} : \mathcal{X}_s^{n+1} \rightarrow [0, \infty)$  is equal to 0 in a subset  $\mathcal{S}_{\mathcal{H}}$  of  $\mathcal{X}_s^{n+1}$  specified by  $\mathcal{H}$ . A sequence  $(x_1, \dots, x_{n+1})$  with  $x_{n+1} \in \mathcal{X}$  belongs to  $\mathcal{S}_{\mathcal{H}}$  whenever  $(x_1, \dots, x_n) \in \mathcal{Z}$  has  $k$  elements in  $\mathcal{X}$  and contains a suffix of length shorter than  $k$  belonging to  $\mathcal{H}$ .

3. **Conjugate prior for reversible random walks with amnesia**

We shall infer the transition probabilities of a reversible random walk with amnesia  $(X_i)_{i \in \mathbb{N}}$  by using the sequence of  $\mathcal{X}$ -valued states  $(Z_1, \dots, Z_m)$ . Our model, like the special case of a variable order Markov chain, mitigates the need to estimate a large number of transition probabilities for dynamics with long memory in a subset of contexts. The major advantage of the approach that we propose is that it does not require model selection or model averaging over the possible sets of

histories  $\mathcal{H}$ . We propose to infer the law of a reversible sequence  $(Z_1, Z_2, \dots)$  by using a Bayesian model for random walks with amnesia. We specify a prior distribution that concentrates on amnesia processes with short memory. These are random walks with stationary measures  $w_n$  that tend to assign higher weights on the low memory states in  $\mathcal{Z}$ , such as  $(s_{n-1}, \dots, s_{n-1}, x_n)$ , compared with the high memory states, such as  $(s_1, x_2, \dots, x_n)$ , where  $(x_2, \dots, x_n) \in \mathcal{X}^{n-1}$ .

Let  $w_{n+1}^0 : \mathcal{X}_s^{n+1} \rightarrow [0, \infty)$  be a balanced function, which determines the transition probabilities of a reversible recurrent random walk with amnesia. This function will specify a prior for  $w_{n+1}$ , the unknown parameter of the process  $(X_i)_{i \in \mathbb{N}}$ . We first define a reinforced process  $(Y_i)_{i \in \mathbb{N}}$  parameterized by  $w_{n+1}^0$ . We then show that it is a mixture of reversible random walks with amnesia. The mixing distribution is a conjugate prior for estimating, given a path  $X_1, \dots, X_N$ , with  $N > n$ , the transition probabilities of the  $n$ th-order random walk with amnesia. As in other reinforcement schemes, such as the Pólya urn and the edge-reinforced random walk (Diaconis and Rolles, 2006), the parameters of the prior, in our case  $w_{n+1}^0$ , and the unknown parameter  $w_{n+1}$  that we want to estimate, are functions defined on  $\mathcal{X}_s^{n+1}$  which share two characteristics: they are both balanced and reversible.

*Definition 2.* The process  $(Y_i)_{i \in \mathbb{N}}$  takes values in  $\mathcal{X}_s$ . With probability 1,  $(Y_1, \dots, Y_n)$  is a palindrome, i.e.  $(Y_1, \dots, Y_n) = (Y_n, \dots, Y_1)$ , and  $w_n^0(Y_1, \dots, Y_n) > 0$ . For  $i > n$ , the transition probabilities

$$\text{pr}(Y_i | Y_1, \dots, Y_{i-1}) = \frac{w_{n+1}^{i-n-1}(Y_{i-n}, \dots, Y_i)}{\sum_{v \in \mathcal{X}_s} w_{n+1}^{i-n-1}(Y_{i-n}, \dots, Y_{i-1}, v)}$$

are specified by the recursive reinforcement equations

$$w_{n+1}^{j+1}(u) = w_{n+1}^j(u) + \mathbf{1}(Y_{j+1} = u_1, \dots, Y_{j+n+1} = u_{n+1}) + \mathbf{1}(Y_{j+1} = u_{n+1}, \dots, Y_{j+n+1} = u_1), \quad (8)$$

where  $j \geq 0$  and  $u = (u_1, \dots, u_{n+1}) \in \mathcal{X}_s^{n+1}$ . The weights  $w_{n+1}^0$  parameterize the law of  $(Y_i)_{i \in \mathbb{N}}$ .

The reinforced processes  $(Y_i)_{i \in \mathbb{N}}$  in definition 2 constitute a subclass of the *reinforced random walks with memory* that was introduced in Bacallado (2011) (definition 3.1). Each process within our subclass is identified by a reversible balanced function  $w_{n+1}^0$  which parameterizes an  $\mathcal{X}_s$ -valued reversible random walk with amnesia. Similarly, reinforced random walks with memory are defined by using reversible balanced functions that parameterize reversible  $n$ th-order Markov chains. The reinforcement mechanisms in definition 2 and in reinforced random walks with memory are identical. The next proposition shows that  $(Y_i)_{i \in \mathbb{N}}$  can be used as a Bayesian model to infer the transition probabilities of  $(X_i)_{i \in \mathbb{N}}$ . The proof of the proposition, together with a few additional remarks on the process  $(Y_i)_{i \in \mathbb{N}}$ , is included in the on-line supplementary materials.

*Proposition 5.* The process  $(Y_i)_{i \in \mathbb{N}}$  is a mixture of reversible random walks with amnesia. For every  $N > n$

$$\text{pr}(Y_1, \dots, Y_N | Y_1, \dots, Y_n) = \int \prod_{i=1}^{N-n} \frac{w_{n+1}(Y_i, \dots, Y_{i+n})}{\sum_{v \in \mathcal{X}_s} w_{n+1}(Y_i, \dots, Y_{i+n-1}, v)} d\mu(w_{n+1}), \quad (9)$$

where  $\mu$  is a distribution on the space of balanced functions that parameterize reversible random walks with amnesia.

This result relies on a de-Finetti-type representation theorem for Markov chains, which relies on the notion of Markov exchangeability (Diaconis and Freedman, 1980; Fortini and Petrone, 2014). We say that a discrete process is Markov exchangeable if the probability of any finite

path can be expressed as a function of the initial state and the transition counts in the path between every pair of states in the state space. Important properties of  $(Y_i)_{i \in \mathbb{N}}$ , including both Markov exchangeability of the sequence of  $n$ -grams visited by the process and recurrence, follow directly from the study of reinforced random walks with memory in Bacallado (2011). It is not difficult to verify Markov exchangeability directly by using a simple closed form expression for the conditional probabilities

$$\begin{aligned} &\text{pr}(Y_{n+1} = y_{n+1}, \dots, Y_N = y_N | Y_1 = y_1, \dots, Y_n = y_n) \\ &= \prod_i \mathbf{1}\{w_{n+1}^0(y_i, \dots, y_{i+n}) > 0\} \times \prod_{x: w_{n+1}^0(x) > 0} g\{x, (y_1, \dots, y_N)\} \\ &\quad \times \prod_{x: w_n^0(x) > 0} g'\{x, (y_1, \dots, y_N)\}, \end{aligned} \tag{10}$$

where  $(y_1, \dots, y_n) = (y_n, \dots, y_1)$ ,  $w_n^0(y_1, \dots, y_n) > 0$ ,

$$g\{x, (y_1, \dots, y_N)\} = \begin{cases} \frac{\Gamma\{w_{n+1}^0(x) + \sum_i \mathbf{1}_x(y_i, \dots, y_{i+n}) + \mathbf{1}_x(y_{i+n}, \dots, y_i)\}^{1/2}}{\Gamma\{w_{n+1}^0(x)\}^{1/2}} & \text{if } (x_1, \dots, x_{n+1}) \neq (x_{n+1}, \dots, x_1), \\ \frac{\Gamma\{w_{n+1}^0(x)/2 + \sum_i \mathbf{1}_x(y_i, \dots, y_{i+n})\} \times 2^{\sum_i \mathbf{1}_x(y_i, \dots, y_{i+n})}}{\Gamma\{w_{n+1}^0(x)/2\}} & \text{if } (x_1, \dots, x_{n+1}) = (x_{n+1}, \dots, x_1), \end{cases}$$

for every  $x = (x_1, \dots, x_{n+1}) \in \mathcal{X}_s^{n+1}$  that satisfies  $w_{n+1}^0(x) > 0$ , and

$$g'\{x, (y_1, \dots, y_N)\} = \begin{cases} \frac{\Gamma\{w_n^0(x)\}^{1/2}}{\Gamma\{w_n^0(x) + \sum_i \mathbf{1}_x(y_i, \dots, y_{i+n-1}) + \mathbf{1}_x(y_{i+n-1}, \dots, y_i)\}^{1/2}} & \text{if } (x_1, \dots, x_n) \neq (x_n, \dots, x_1), \\ \frac{\Gamma\{[1 - \mathbf{1}_x(y_1, \dots, y_n) + w_n^0(x)]/2\} \times 2^{-\sum_i \mathbf{1}_x(y_i, \dots, y_{i+n-1})}}{\Gamma\{[1 - \mathbf{1}_x(y_1, \dots, y_n) + w_n^0(x)]/2 + \sum_i \mathbf{1}_x(y_i, \dots, y_{i+n-1})\}} & \text{if } (x_1, \dots, x_n) = (x_n, \dots, x_1), \end{cases}$$

for every  $x = (x_1, \dots, x_n) \in \mathcal{X}_s^n$  such that  $w_n^0(x) = \sum_{x_{n+1}} w_{n+1}^0(x_1, \dots, x_n, x_{n+1}) > 0$ . The above expression specifies the law of  $(Y_i)_{i \in \mathbb{N}}$  using the gamma function,  $\Gamma: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ , and has similarities to other well-known reinforcement schemes such as the edge-reinforced random walk (Diaconis and Rolles, 2006) and the Pólya urn model. The expression shows that, conditionally on  $(Y_1 = y_1, \dots, Y_n = y_n)$ , the probability of the event  $(Y_1 = y_1, \dots, Y_N = y_N)$ ,  $N > n$ , depends only on the counts

$$\sum_{i \leq N-n} \mathbf{1}_x(y_i, \dots, y_{i+n}), \quad x \in \mathcal{X}_s^{n+1}.$$

Therefore, the sequence of  $n$ -grams  $(Y_1, \dots, Y_n), (Y_2, \dots, Y_{n+1}), \dots$ , is Markov exchangeable.

In what follows we use the mixing distribution  $\mu$  in proposition 5 as Bayesian prior for the unknown parameter  $w_{n+1}$  of our amnesia process  $(X_i)_{i \in \mathbb{N}}$ . Under this Bayesian model we can sample from the predictive distributions  $\text{pr}(Y_{N+1} | Y_1, \dots, Y_N)$ . The computation of these predictive probabilities reduces to the evaluation of the reinforced weights  $w_{n+1}^N$  in definition 2.

The fact that the initial  $n$ -gram in  $(Y_i)_{i \in \mathbb{N}}$  is a palindrome is a necessary hypothesis for proposition 5. However, we can define priors for a random walk with amnesia with initial  $n$ -gram  $(y_1, \dots, y_n) \neq (y_n, \dots, y_1)$  by using the distribution of  $(Y_i)_{i > m}$  conditional on  $Y_1 = y_1, \dots, Y_m = y_m$ , with  $(y_1, \dots, y_m)$  specified up front together with  $w_{n+1}^0$ .



*Remark 1.* Consider a prior parameter  $w_{n+1}^0$  with  $w_{n+1}^0(x) > 0$  for every  $x = (x_1, \dots, x_{n+1})$  such that the transition from  $(x_1, \dots, x_n) \in \mathcal{X}_s^n$  to  $x_{n+1}$  is allowed by the definition (1) of a random walk with amnesia. Assume that the investigator observes the trajectory  $(X_1, \dots, X_s)$  of a recurrent random walk with amnesia with unknown transition probabilities, and fix a recurrent  $n$ -gram  $(x_1, \dots, x_n) \in \mathcal{X}_s^n$ . The linear reinforcement (8) indicates that the Bayesian estimate of the transition probability  $w_{n+1}(x_1, \dots, x_{n+1}) / \sum_v w_{n+1}(x_1, \dots, x_n, v)$ , defined as the posterior mean, converges in probability to the empirical estimate

$$\frac{\sum_{j=1}^s \mathbf{1}(X_j = x_1, \dots, X_{j+n} = x_{n+1}) + \mathbf{1}(X_j = x_{n+1}, \dots, X_{j+n} = x_1)}{\sum_{j=1}^s \mathbf{1}(X_j = x_1, \dots, X_{j+n-1} = x_n) + \mathbf{1}(X_j = x_n, \dots, X_{j+n-1} = x_1)}$$

when the length of the trajectory  $s$  diverges. This fact follows from the convergence in probability, for any  $(x_1, \dots, x_n) \in \mathcal{X}_s^n$  and  $(x_1, \dots, x_{n+1}) \in \mathcal{X}_s^{n+1}$ , of the ratios

$$\frac{\sum_{j=1}^{s-n} \mathbf{1}(X_j = x_1, \dots, X_{j+n} = x_{n+1})}{s},$$

and

$$\frac{\sum_{j=1}^{s-n+1} \mathbf{1}(X_j = x_1, \dots, X_{j+n} = x_n)}{s}$$

to the parameters  $w_{n+1}(x_1, \dots, x_{n+1})$  and  $w_{n+1}(x_1, \dots, x_n)$  respectively.

#### 4. Posterior simulations

The process  $(Y_i)_{i \in \mathbb{N}}$  in definition 2 is used as a prior distribution for the random walk with amnesia. This process is driven by a random  $n$ th-order transition matrix (proposition 5). Recall that the observable sequence  $(Z_j)_{j \in \mathbb{N}}$  takes values in  $\mathcal{X}$ . Let  $Z_j = Y_{\tau_j}$ , where  $\tau_j = \min\{\tau; \tau > \tau_{j-1}, (Y_{\tau-n}, \dots, Y_{\tau-1}) \in \mathcal{Z}, Y_\tau \in \mathcal{X}\}$  and  $\tau_0 = n$ . The length of memory of the process  $(Z_j)_{j \in \mathbb{N}}$  at each transition is captured by a sequence of latent variables  $T_j := \max\{t : \tau_j - t > 0, (Y_{\tau_{j-1}}, \dots, Y_{\tau_j-t}) \in \mathcal{X}^t\}$ ,  $j \geq 1$ . Without loss of generality, assume that  $Y_1 = \dots = Y_n = s_n$ .

The observed states  $Z_1, \dots, Z_{r-1}$  and the lengths of memory  $T_1, \dots, T_{r-1}$  identify the path  $Y_1, \dots, Y_{\tau_{r-1}}$ , which in turn can be used to obtain the reinforced weights  $\{w_{n+1}^{\tau_{r-1}-n}(u); u \in \mathcal{X}_s^{n+1}\}$  and to compute the predictive distribution

$$\text{pr}(Y_{1+\tau_{r-1}}, \dots, Y_{\tau_r} | T_1, \dots, T_{r-1}, Z_1, \dots, Z_{r-1}) = \text{pr}(Y_{1+\tau_{r-1}}, \dots, Y_{\tau_r} | Y_1, \dots, Y_{\tau_{r-1}}).$$

This makes it straightforward to sample  $T_r$  conditionally on  $(T_1, \dots, T_{r-1}, Z_1, \dots, Z_r)$ , which will be useful in the sequential importance sampling algorithm (Gordan *et al.*, 1993) that is proposed below.

The goal of the algorithm is to infer the lengths of memory  $T_1, \dots, T_m$  given a sequence of observations  $Z_1, \dots, Z_m$ . The algorithm is initialized with particles  $t_1^{(i)} = 0$  and importance weights  $v_1^{(i)} = 1$  for all  $i = 1, \dots, N$ . For each  $r = 2, \dots, m$ , the particle  $t_r^{(i)}$  is sampled from the conditional distribution of  $T_r$  given the observed random variables  $(Z_1 = z_1, \dots, Z_r = z_r)$  and  $(T_1 = t_1^{(i)}, \dots, T_{r-1} = t_{r-1}^{(i)})$  for all  $i = 1, \dots, N$ . At each step, the importance weights are updated by

$$\mathbf{v}_r^{(i)} = \mathbf{v}_{r-1}^{(i)} \text{pr}(Z_r = z_r | Z_1 = z_1, \dots, Z_{r-1} = z_{r-1}, T_1 = t_1^{(i)}, \dots, T_{r-1} = t_{r-1}^{(i)}).$$

Finally, we approximate the posterior distribution of the lengths of memory by using the weighted particles,

$$\text{pr}(T_1, \dots, T_m | Z_1, \dots, Z_m) \approx \sum_{i=1}^N \tilde{\mathbf{v}}_m^{(i)} \mathbf{1}(T_1 = t_1^{(i)}, \dots, T_m = t_m^{(i)}), \tag{11}$$

where

$$\tilde{\mathbf{v}}_r^{(i)} = \frac{\mathbf{v}_r^{(i)}}{\sum_{j=1}^N \mathbf{v}_r^{(j)}}.$$

It is well understood that, in many cases, a better approximation is obtained by combining the structure of sequential importance sampling with resampling operations (Gordon *et al.*, 1993). We shall apply a simple strategy known as sequential importance resampling, or the bootstrap particle filter. A resampling operation at time  $r$  consists of the following two steps. First, after the weights  $\tilde{\mathbf{v}}_r^{(i)}$ ,  $i = 1, \dots, N$ , have been computed,  $N$  random variables  $a_i$ ,  $i = 1, \dots, N$ , taking values in  $\{1, \dots, N\}$  are independently generated with  $\text{pr}(a_i = j) = \tilde{\mathbf{v}}_r^{(j)}$ ,  $j = 1, \dots, N$ . Second, each particle  $(t_1^{(i)}, \dots, t_r^{(i)})$  is replaced by  $(t_1^{(a_i)}, \dots, t_r^{(a_i)})$ , and the weights  $\mathbf{v}_r^{(i)}$  are set equal to 1, for  $i = 1, \dots, N$ . The algorithm is then ready to proceed as described in the previous paragraph.

Let  $r_1, \dots, r_k$  be the times at which resampling steps are made, and let  $\hat{\mathbf{v}}_r$  be the average of  $\mathbf{v}_r^{(1)}, \dots, \mathbf{v}_r^{(N)}$  before the resampling operation at time  $r$ . Then, it can be shown that  $\prod_{j=1}^k \hat{\mathbf{v}}_{r_j}$  is an unbiased estimate for the marginal probability of the observation  $\text{pr}(Z_1, \dots, Z_{r_k})$ ; see proposition 7.4.1 in del Moral (2004). This will be convenient in the following sections for estimating predictive probabilities of the form  $\text{pr}(Z_{r+1}, \dots, Z_m | Z_1, \dots, Z_r)$ .

The particles  $(t_1^{(i)}, \dots, t_m^{(i)}; i = 1, \dots, N)$  approximate the posterior distribution of the length of memories  $T_1, \dots, T_m$ . They can be used to sample approximately from the predictive distribution of the process  $(Z_i)_{i \in \mathbb{N}}$ . A long simulation from the predictive distribution yields, by de Finetti representation (proposition 5) arguments, an approximate posterior sample of the transition probabilities that drive the random walk with amnesia. Posterior samples of  $w_{n+1}$  can be used, among other things, to calculate predictive probabilities of the form  $\text{pr}(Z_{r+1}, \dots, Z_m | Z_1, \dots, Z_r)$  via a forward-backward algorithm; this is an alternative to the estimates that were mentioned above.

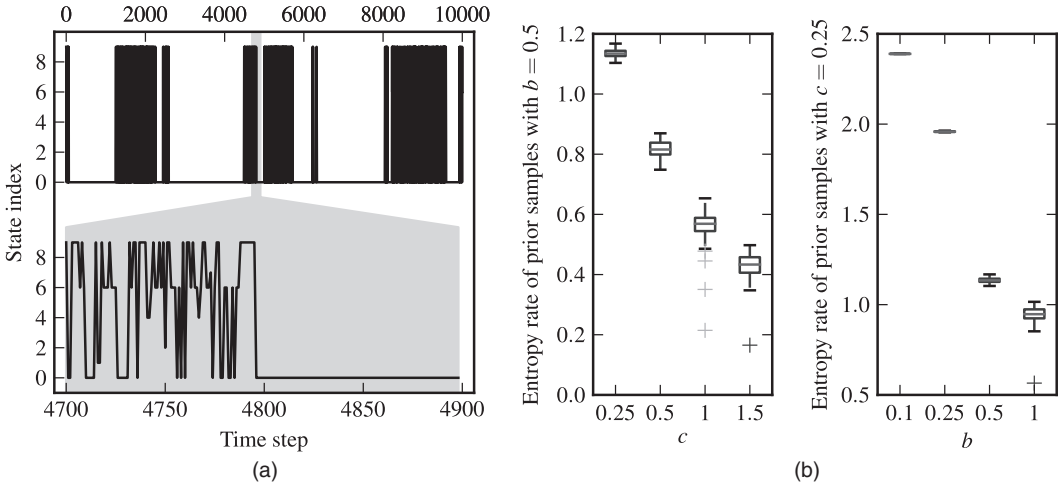
### 5. Simulation study

We specify the prior model choosing, for every  $(x_1, \dots, x_n) \in \mathcal{Z}$ ,

$$w_{n+1}^0(x_1, \dots, x_{n+1}) = c(|\mathcal{X}|b)^{-\sum_i \mathbf{1}(x_i \in \mathcal{X})}, \tag{12}$$

for some  $b, c > 0$ , if the random walk with amnesia allows the transition  $(x_1, \dots, x_n) \rightarrow x_{n+1}$ , and  $b, c = 0$  otherwise. This reduces the set of hyperparameters that we need to tune to the pair  $(b, c)$ . The initial function  $w_{n+1}^0$  corresponds to a random walker who, from any  $n$ -gram in  $\mathcal{Z}$ , with some fixed probability forgets the first  $\mathcal{X}$ -valued element in his memory. Also, under  $w_{n+1}^0$ , transitions to the  $\mathcal{X}$ -valued states from any  $n$ -gram in  $\mathcal{Z}$  are all equally likely.

The parameter  $b$  tunes the length of memory. The greater  $b$ , the shorter the length of memory tends to be during the reinforced walk  $(Y_i)_{i \geq 1}$ . The parameter  $c$  has an interpretation that is similar to the total initial mass of the edge-weighted graph in Diaconis and Rolles (2006) and tunes the



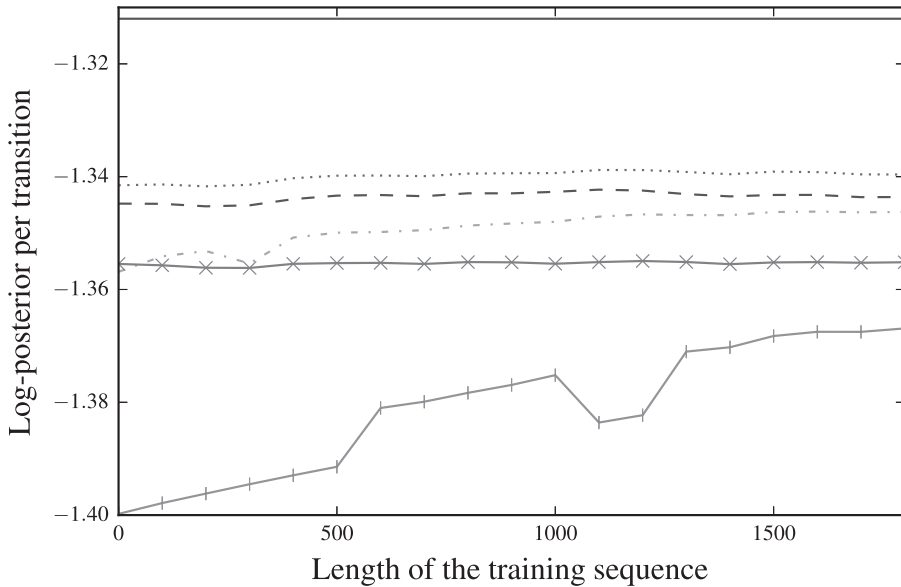
**Fig. 1.** (a) Conformational states of the WW domain, indexed by  $0, \dots, 9$ , observed in a molecular dynamics simulation (the top panel shows the full sequence of 10000 steps whereas the bottom panel shows a magnified subsequence) and (b) entropy rates  $\lim_{n \rightarrow \infty} n^{-1} H(X_1, \dots, X_n)$  of 100 random walks with amnesia sampled from the prior distribution under eight choices of  $(b, c)$

concentration of  $w_{n+1}$  around the prior mean. Fig. 1(b) illustrates how the choice of the parameters  $(b, c)$  tunes the complexity of a typical random walk with amnesia with 10 states sampled from the prior. As a measure of complexity, we use the entropy rate  $\lim_{n \rightarrow \infty} n^{-1} H(X_1, \dots, X_n)$  of the random walk with amnesia, where  $H$  is the Shannon entropy.

We tested the proposed model by using data simulated from a reversible fourth-order Markov chain in the space of nucleic bases  $\{a, c, g, t\}$ . The transition probabilities of the Markov chain were the frequencies of every fourth-order transition in a cyclic deoxyribonucleic acid (DNA) sequence of 37243 bases read in both directions (reference sequence NC 021042). The random walk with amnesia was fitted to a simulated sequence by sequential importance resampling, and the predictive probability of a different simulated sequence was approximated by using the importance weights as discussed at the end of Section 4. We used  $N = 1000$  particles and performed resampling operations every three steps. Fig. 2 shows the log-predictive probabilities per step given training sequences of lengths ranging between 100 and 1900. The hyperparameters  $b$  and  $c$  were selected to optimize the probability of the training sequence provided by the sequential importance sampling algorithm over a grid with  $b \in \{0.5, 1, 1.5, 2\}$  and  $c \in \{0.5, 1, 1.5, 2\}$ .

For comparison, Fig. 2 shows the same statistics for a Bayesian analysis of first-order and fourth-order reversible Markov models (Bacallado, 2011) For each model, the concentration parameter, i.e. the initial weight of the palindrome in the model of Bacallado (2011), is optimized to maximize the probability of the training sequence. Even though the sequence was generated from a fourth-order Markov chain, the fourth-order model performs poorly compared with the random walk with amnesia. We observe, as expected, that the inferred random walk with amnesia adapts to the complexity of the training sequence, with longer lengths of memory associated with the more frequent contexts.

Bacallado (2011) also defined a prior for reversible variable order Markov models. We compared the random walk with amnesia with a Bayesian analysis based on variable order Markov models, with and without reversibility. In the non-reversible case, we apply an independent Dirichlet prior to the transition probabilities out of every context. Before applying these models, it is necessary to estimate the appropriate lengths of memory. For this, we applied a context



**Fig. 2.** Reversible simulation example—log-posterior predictive probabilities per transition for a left-out test set, as a function of the length of the training sequence: — — —, first-order reversible Markov model; +, fourth-order reversible Markov model; ······, random walk with amnesia; - · - · - ·, reversible variable order model selection; ×, variable order Markov model selection; ———, expected log-likelihood of a transition per step in the model used to simulate the sequence

tree pruning algorithm that was similar in structure to the algorithm Context of Rissanen (1983).

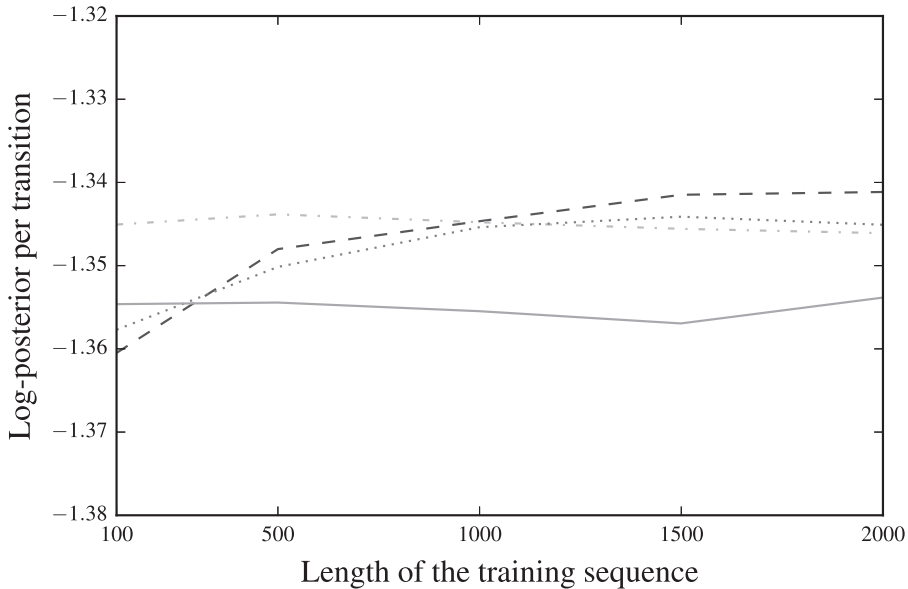
We describe the main characteristics of the pruning algorithm. First, the algorithm grows the largest possible context tree in which every context appears at least five times in the training sequence and is no longer than six. Then, we employ a backward tree pruning procedure with a local criterion. The criterion is a Bayes factor comparing the variable order model before and after pruning the tree. We loop through the leaves in alphabetical order; in Context, which uses a slightly different local criterion, the order is irrelevant. The threshold for the Bayes factor and the concentration parameter of the prior for the variable order Markov models are used as tuning parameters and are selected to optimize the probability of the training data under the model. The pruning strategy is slightly modified in the reversible case, as the set of contexts must satisfy the closure properties in proposition 4.2 of Bacallado (2011).

Fig. 2 shows that, as expected, the reversible models outperform non-reversible models. The random walk with amnesia outperforms both of the model selection schemes in this example.

To evaluate the performance of the methods above on data which are not reversible, we repeated the analysis by using the original cyclic DNA sequence instead of the reversible simulated data. Fig. 3 shows that, as we would expect, models that enforce reversibility pay a price in bias and perform worse than an appropriately pruned variable order Markov model for longer training sets. This is not so when reversibility is given, as in the application of the following section.

## 6. Molecular dynamics

We analyse a molecular dynamics simulation of a protein known as the WW domain (Shaw



**Fig. 3.** Non-reversible example—log-posterior predictive probabilities per transition for a left-out test set, as a function of the length of the training sequence: — — —, variable order Markov model selection; - · - · -, first-order reversible Markov model; ·····, reversible variable order model selection; ———, random walk with amnesia

*et al.*, 2010). Markov models have become, in recent years, essential tools for the analysis of this type of dynamical data (Prinz *et al.*, 2011). The data consist of a sequence of states that are representative of distinct conformations of the protein observed at regular time intervals. In this application, the stochastic law of the process is reversible owing to the nature of Hamiltonian mechanics. We have 10 states in  $\mathcal{X}$  and a trajectory of length 10 000, with one observation every 20 ns. We use the first half of the trajectory as a training data set and the second as a test data set.

Fig. 1(a) plots the trajectory. The plot shows that the process alternates phases during which it remains stable at state 0 with phases during which it rapidly moves across states. These phases correspond to periods during which the protein is folded and unfolded respectively.

We trained a reversible random walk with amnesia with  $n = 6$  by sequential importance resampling. We used  $N = 1000$  particles and performed resampling operations every three steps. The hyperparameters  $b$  and  $c$  were chosen to maximize the probability of the training sequence under the model. Table 1 shows the parameters selected and the log-likelihood of the test sequence. To evaluate the efficacy of the sequential importance resampling algorithm, we replicated the computations 10 times at every hyperparameter setting. The unbiased estimates of marginal likelihood varied little in these iterations; for example, the estimate that is shown in the last row of Table 1 had a standard deviation of 67.

To compare our procedure with some common alternatives, we considered five methods that were reviewed by Begleiter *et al.* (2004), none of which takes into account reversibility. Following the recommendations in Begleiter *et al.* (2004), we tuned the parameters of each method via threefold cross-validation using only the training sequence. We then computed the log-likelihood of the test data set at the optimal parameters. Table 1 lists the methodologies, the tuning parameters that were selected, in the notation of Begleiter *et al.* (2004), and the log-likelihood of the test trajectory. The best of the methods in this example is decomposed context

**Table 1.** Predictive performance of the random walk with amnesia compared with the methods reviewed in Begleiter *et al.* (2004) and two model selection approaches for variable order Markov models†

| Method                                 | Test log-likelihood | Tuning parameters   |
|--|---------------------|---|
| Decomposed context tree weighting      | -3210.4             | $D \in \{1, 3, 5, 7, 9, 15, 20\}$   |
| Prediction by partial matching         | -3359.8             | $D \in \{1, 3, 5, 7, 9, 15, 20\}$   |
| Lempel Ziv 78                          | -4041.4             | None  |
| Lempel Ziv MS                          | -3862.1             | $M \in \{0, 2, 4, 6, 8\}$<br>$S \in \{0, 2, 4, 6, 8\}$  |
| Probabilistic suffix trees             | -3280.9             | $P_{\min} \in \{0.0001, 0.001, 0.01, 0.1\}$<br>$\gamma \in \{0.0001, 0.001, 0.01, 0.1\}$<br>$D \in \{1, 3, 5, 7, 9, 15, 20\}, r = 1.05, \alpha = 0$ |
| Variable order Markov model            | -3308.4             | $c \in \{0.0001, 0.001, 0.01, 0.1, 1\}, t \in \{-10, 0, 10\}$   |
| Reversible variable order Markov model | -3268.8             | $c \in \{0.0001, 0.001, 0.01, 0.1, 1\}, t \in \{-10, 0, 10\}$   |
| Random walk with amnesia               | -3015.4             | $b \in \{0.1, 0.4, 0.6, 1, 1.5\}, c \in \{0.5, 1, 1.5, 2\}$   |

†The tuning parameters of each method from Begleiter *et al.* (2004) were chosen via cross-validation using 50% of the data. We report the range of values explored for each tuning parameter and the optimal value in italics. We also report the corresponding log-likelihood of the test data set.

tree weighting. The random walk with amnesia has a higher log-likelihood than all five methods have.

Finally, we applied the Bayesian analyses of variable order Markov models that were described in Section 6, selecting the model selection threshold  $t$  and the concentration parameter  $c$  which optimize the likelihood of the training data.

The results are summarized in Table 1. The reversible model outperforms the non-reversible model, but both of the methods are outperformed by decomposed context tree weighting and the random walk with amnesia. The tuning of the random walk with amnesia was repeated by using a threefold cross-validation strategy, and this had a negligible effect on the results in Table 1.

### Acknowledgements

SF is supported by the European Research Council through grant StG N-BNP 306406, LT has been supported by the Claudia Adams Barr Program in Innovative Cancer Research and SB received funding from the Stein Fellowship.

### References

Annis, D., Kiessler, P., Lund, R. and Steuber, T. (2010) Estimation in reversible markov chains. *Am. Statistn.*, **64**, 116–120.

Bacallado, S. (2011) Bayesian analysis of variable-order, reversible markov chains. *Ann. Statist.*, **39**, 838–864.

Begleiter, R., El-Yaniv, R. and Yona, G. (2004) On prediction using variable order Markov models. *J. Artif. Intell. Res.*, **22**, 385–421.

Bühlmann, P. (2000) Model selection for variable length Markov chains and tuning the context algorithm. *Ann. Inst. Statist. Math.*, **52**, 287–315.

Bühlmann, P. and Wyner, A. J. (1999) Variable length Markov chains. *Ann. Statist.*, **27**, 480–513.

Cappé, O., Moulines, E. and Ryden, T. (2005) *Inference in Hidden Markov Models*. New York: Springer.

Diaconis, P. and Freedman, D. (1980) de Finetti’s theorem for Markov chains. *Ann. Probab.*, **8**, 115–130.

Diaconis, P. and Rolles, S. (2006) Bayesian analysis for reversible Markov chains. *Ann. Statist.*, **34**, 1270–1292.

Fortini, S. and Petrone, S. (2014) Predictive characterizations of mixtures of Markov chains. *Preprint arXiv:1406.5421*.

- Gordon, N. J., Salmond, D. J., Smith, A. F. M. and Steuber, T. (1993) Novel approach to nonlinear/non-Gaussian Bayesian state estimation. *IEEE Proc. F*, **140**, 107–113.
- Kelly, F. P. (1979) *Reversibility and Stochastic Networks*. Chichester: Wiley.
- MacKay, D. and Peto, L. (1995) A hierarchical Dirichlet language model. *Nat. Lang. Eng.*, **1**, 289–308.
- Mochihashi, D. and Sumita, E. (2008) The infinite Markov model. *Adv. Neur. Inf. Process. Syst.*, **20**, 1017–1024.
- del Moral, P. (2004) *Feynman-Kac Formulae*. New York: Springer.
- Palla, K., Knowles, D. and Ghahramani, Z. (2014) A reversible infinite HMM using normalised random measures. In *Proc. 31st Int. Conf. Machine Learning*. Montreal: International Machine Learning Society.
- Prinz, J., Wu, H., Sarich, M., Keller, B., Senne, M., Held, M., Chodera, J., Schütte, C. and Noé, F. (2011) Markov models of molecular kinetics: generation and validation. *J. Chem. Phys.*, **134**, article 174105.
- Rissanen, J. (1983) A universal data compression system. *IEEE Trans. Inform. Theor.*, **29**, 656–664.
- Ron, D., Singer, Y. and Tishby, N. (1996) The power of amnesia: learning probabilistic automata with variable memory length. *Mach. Learn.*, **25**, 117–149.
- Shaw, D., Maragakis, P., Lindorff-Larsen, K., Piana, S., Dror, R., Eastwood, M., Bank, J., Jumper, J., Salmon, J., Shan, Y. and Wriggers, W. (2010) Atomic-level characterization of the structural dynamics of proteins. *Science*, **330**, 341–346.
- Teh, Y. (2006) A hierarchical Bayesian language model based on Pitman-Yor processes. In *Proc. 21st Int. Conf. Computational Linguistics*, pp. 985–992. Stroudsburg: Association for Computational Linguistics.
- Van Kampen, N. G. (1992) *Stochastic Processes in Physics and Chemistry*. Amsterdam: North-Holland.
- Weinberger, M., Rissanen, J. and Feder, M. (1995) A universal finite memory source. *IEEE Trans. Inform. Theor.*, **41**, 643–652.
- Wood, F., Archambeau, C., Gasthaus, J., James, L. F. and Teh, Y. W. (2009) A stochastic memorizer for sequence data. In *Proc. 26th Int. Conf. Machine Learning*, pp. 1129–1136. Montreal: International Machine Learning Society.

#### Supporting information

Additional ‘supporting information’ may be found in the on-line version of this article:

‘Supplementary information’.