

Reconstructing ancestral protein sequences and structures

Michael Golden¹, Thomas Hamelryck², and Oliver Pybus¹

¹Department of Zoology, University of Oxford, UK

²Bioinformatics Centre, Section for Computational and RNA Biology, Department of Biology and Image Section, Department of Computer Science, University of Copenhagen

May 3, 2019

Abstract

We present a probabilistic model of protein evolution that captures several important features of protein sequence and local structure. The key feature being dependencies between neighbouring amino acid positions that are temporal in nature due to sequence mutations that occur during evolution. The model is trained on a large number of protein alignments and corresponding phylogenetic trees that represent the evolutionary history of the aligned proteins. This yields a model that acts as a rich prior distribution over protein evolution that can be used to perform several important inference tasks. One such task being Bayesian reconstruction of ancestral virus protein sequences, which we demonstrate to have better accuracy than competing methods. The model provides a complete probabilistic description of each protein’s backbone structure using an angle and bond length representation. Structure evolution is modelled jointly with sequence, permitting ancestral structures and sequences to be reconstructed in a phylogenetically rigorous manner. Likewise, the model can perform homology modelling to predict the unknown local backbone structure of a known protein sequence using additional information from potentially large numbers of homologous proteins. The model is highly flexible with respect to input, implying that arbitrary combinations of protein sequences and structures can be used when performing various inference tasks. The current model does not capture global features of protein structure that are necessary for accurate homology modelling or reconstruction of ancestral three-dimensional structures. However, it is ultimately expected to be combined with protein structure prediction models that account for such long-range dependencies, but that do not account for evolutionary information that can substantially enhance predictions of structures.

1 Introduction

2 Methods

2.1 Model

2.1.1 Model of a single protein

A single protein consisting of n amino acids:

$$P_a = (H_a, A_a, X_a) \\ = \langle (H_a^1, A_a^1, X_a^1), \dots, (H_a^n, A_a^n, X_a^n) \rangle$$

is a sequence of aligned sites where each site i is associated with a discrete-valued hidden state, H_a^i (taking on one of h possible values), a discrete-valued amino acid observation, A_a^i (representing one of the twenty possible amino acids), and a corresponding vector of continuous-valued structural observations X_a^i representing the backbone structure of the protein.

Structural observations The set of structural observations, X_a^i , at a particular site, i , consists of nine continuous-valued variables: three dihedral angles $(\phi_i, \psi_i, \omega_i)$, three additional bond angles $(\tau_i^{(1)} = \overrightarrow{C\alpha_{i-1}, C_{i-1}, N_i}, \tau_i^{(2)} = \overrightarrow{C_{i-1}, N_i, C\alpha_i}, \tau_i^{(3)} = \overrightarrow{N_i, C\alpha_i, C_i})$, and three bond lengths $(b_i^{(1)} = \overrightarrow{C_{i-1}, N_i}, b_i^{(2)} = \overrightarrow{N_i, C\alpha_i}, b_i^{(3)} = \overrightarrow{C\alpha_i, C_i})$.

Note that ϕ_1 , $\tau_1^{(1)}$, and $\tau_1^{(2)}$, are undefined at the first position in the peptide backbone of an unaligned protein. Similarly, ψ_n and ω_n are undefined for the last position, n , in each unaligned protein.

Given the structural observations, X_a , it is possible to exactly reconstruct the three-dimensional coordinates of each atom in a protein's backbone (Parsons *et al.*, 2005).

The ω_i dihedral angle (which determines the cis/trans conformation) at each site i is assumed to be distributed according a univariate von Mises (vM) distribution with mean μ_ω and concentration parameter κ_ω conditional on the hidden state H_a^i and the amino acid A_a^i :

$$\omega_i \sim \text{vM}(\mu_\omega(H_a^i, A_a^i), \kappa_\omega(H_a^i, A_a^i)). \quad (1)$$

The ϕ_i and ψ_i dihedral angles are assumed to be drawn from a bivariate von Mises (bvM) distribution with mean vector $\mu_{\phi, \psi} = \langle \mu_\phi, \mu_\psi \rangle$ and covariance parameters $\kappa_{\phi, \psi} = \langle \kappa_1, \kappa_2, \kappa_3 \rangle$:

$$(\phi_i, \psi_i) \sim \text{bvM}(\mu_{\phi, \psi}(H_a^i, A_a^i), \kappa_{\phi, \psi}(H_a^i, A_a^i)), \quad (2)$$

where κ_1 is the variance associated with ϕ , κ_2 is the variance associated with ψ , and κ_3 is the correlation between ϕ and ψ ,

The three additional bond angles $(\tau_i^{(1)}, \tau_i^{(2)}, \tau_i^{(3)})$ are each distributed according a univariate (vM) distribution conditional on the hidden state H_a^i only:

$$\begin{aligned} \tau_i^{(1)} &\sim \text{vM}(\mu_{\tau^{(1)}}(H_a^i), \kappa_{\tau^{(1)}}(H_a^i)) \\ \tau_i^{(2)} &\sim \text{vM}(\mu_{\tau^{(2)}}(H_a^i), \kappa_{\tau^{(2)}}(H_a^i)) \\ \tau_i^{(3)} &\sim \text{vM}(\mu_{\tau^{(3)}}(H_a^i), \kappa_{\tau^{(3)}}(H_a^i)). \end{aligned} \quad (3)$$

The three bond lengths $(b_i^{(1)} > 0, b_i^{(2)} > 0, b_i^{(3)} > 0)$ are distributed according a truncated Multivariate Normal (MVN) with mean vector, μ , of length 3 and a 3×3 covariance matrix, Σ , conditional on the hidden state H_a^i :

$$(b_i^{(1)}, b_i^{(2)}, b_i^{(3)}) \sim \text{MVN}(\mu(H_a^i), \Sigma(H_a^i)). \quad (4)$$

Note that the parameters in (3) and (4) are no longer conditional upon the amino acid observations A_a^i . This was done to reduce the number of model parameters, as the values of the three bond angles and three bond lengths are all largely invariant, implying that the values can be reasonably fixed. Despite this, we still opted to treat them as random variables so that the model gives a complete probabilistic description of a protein backbone structure.

Site likelihood The likelihood of a structural observation, $p(X_a^i | H_a^i, A_a^i, \hat{\theta})$, at site i conditional on the hidden state, H_a^i , and amino acid, S_a^i , is given by a product of the densities in Equations 1, 2, 3, and 4.

Hidden Markov model A sequence of structural observations representing a single protein backbone structure is modelled using a Hidden Markov Model

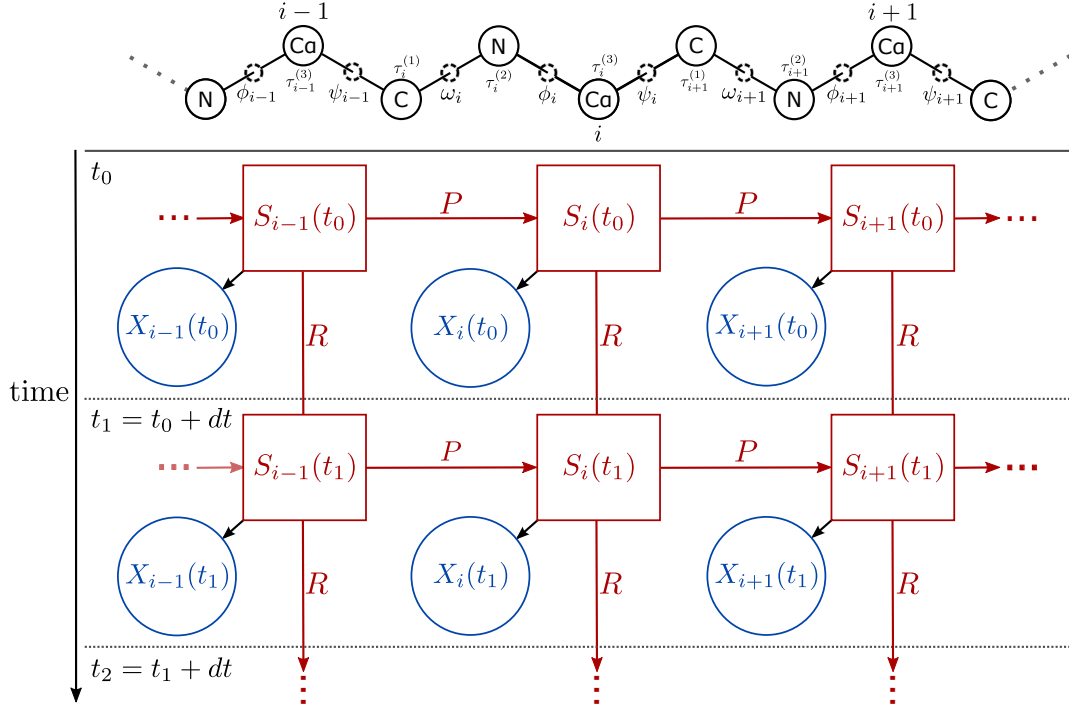


Figure 1: Above: a depiction of a protein backbone (three amino acids long) with the ω , ϕ and ψ dihedral angles and the three additional bond angles ($\tau_i^{(1)}$, $\tau_i^{(2)}$, $\tau_i^{(3)}$) shown. Bond lengths are implicit. Bond angles and bond lengths are not to scale. Also shown are C_α atoms which attach to the amino acid side-chains. Each amino acid side-chain determines the characteristic nature of each amino acid. Every amino acid position corresponds to a hidden node in the model below.

Below: Graphical depiction of the model architecture showing three amino acid positions ($i-1$, i , and $i+1$) at two time instants (t_0 and t_1) along a single branch of a phylogenetic tree. Note that $S_i(t) = (H_i(t), A_i(t))$.

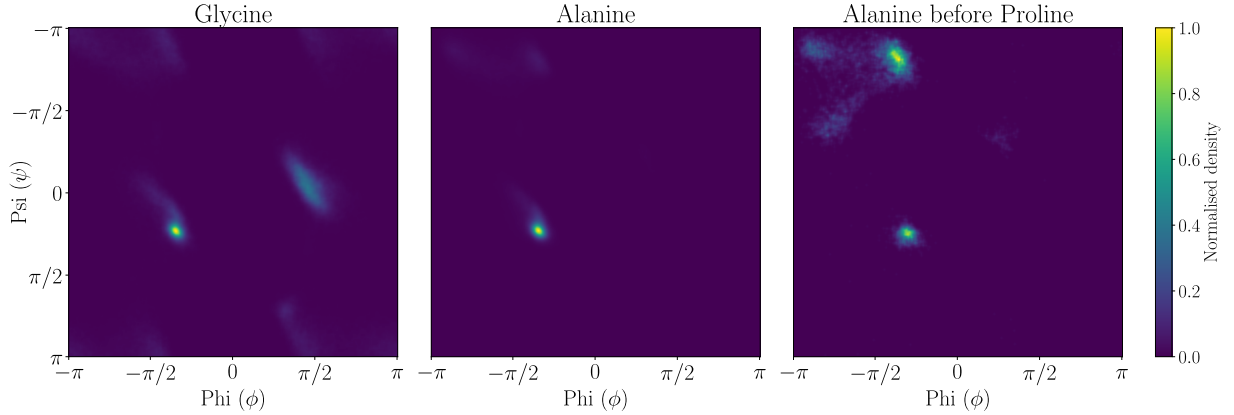


Figure 2: Ramachandran plots depicting empirical distributions of ϕ and ψ angle combinations under three different amino acid contexts. Glycine has the smallest amino acid side-chain and is therefore the least constrained. Alanine has a larger side-chain than glycine, constraining most ϕ and ψ angle to lie near a single peak. When alanine precedes proline in the peptide backbone, the ϕ and ψ angle combinations previously favoured become sterically hindered **MG: is this the correct terminology?**, favouring angle combinations at multiple peaks.

(HMM). Hidden states in the HMM are primarily intended to encode the angle and bond lengths distributions and their association with the different amino acids as specified in Equations 1-4. However, the HMM is also critically to capture neighbouring dependencies, such as steric effects on dihedral angle conformations (Figure 2). These neighbouring dependencies are captured using a $h \times h$ transition probability matrix $P = p(H_a^i | H_a^{i-1}, \hat{\theta})$.

We let $\Pi(a)$ denote the joint likelihood of a sequence of hidden states (H_a) and structural observations (X_a) conditional on a sequence of amino acids (A_a) defined as follows:

$$\begin{aligned} \Pi(a) &= p(H_a, A_a, X_a | \hat{\theta}) = \\ &= p(A_a^1, X_a^1 | H_a^1, \hat{\theta}) p(H_a^1 | \hat{\theta}) \\ &\times \prod_{i=2}^n p(A_a^i, X_a^i | H_a^i, \hat{\theta}) p(H_a^i | H_a^{i-1}, \hat{\theta}), \end{aligned} \quad (5)$$

where $p(H_a^1 | \hat{\theta})$ is the initial probability of starting in state H_a^1 at the first site. Whilst $\Pi(a)$ describes a single protein, it is used in the next section to ‘weight’ an evolutionary model such that protein evolutionary trajectories are visited with probability proportional to $\Pi \times \pi$, where π is weighting corresponding to the amino acid sequences.

2.1.2 Evolutionary model

Thus far we have only considered a single protein. In this section we outline how multiple phylogenetically related proteins are modelled evolutionarily. Following Choi *et al.* (2008) we construct a rate matrix that represents changes between two sequences, a and b , instead of character states, such as amino acids, as is typical of substitution models. Furthermore, each sequence position combines a hidden state (h_i), an amino acid (aa_i), and a set of structural observations (x_i), into a joint character state (h_i^s, aa_i^s, x_i^s), where s refers to a particular sequence. The rate matrix is

given as follows:

$$R_{ab} = \begin{cases} \sqrt{\frac{\Pi(b)}{\Pi(a)}} U_{a_i b_i} \pi_{aa_i^b}^{h_i^b} & \text{Single amino acid} \\ & \text{difference at site } i. \\ \sqrt{\frac{\Pi(b)}{\Pi(a)}} V_{a_i b_i} \pi_{aa_i^b}^{h_i^b} & \text{Single hidden state} \\ & \text{difference at site } i. \\ 0 & \text{Both hidden state and} \\ & \text{amino acid differences} \\ & \text{at site } i. \\ 0 & \text{Differences at two} \\ & \text{or more sites.} \\ -\sum_{k \neq a} R_{ak} & a = b \end{cases} \quad (6)$$

where i is the position that differs between sequence a and b , U is a symmetric 20×20 amino acid exchangeability matrix, and V is a symmetric $h \times h$ hidden state exchangeability matrix. The term, $\frac{\Pi(b)}{\Pi(a)}$, weights each hidden state or amino acid change, such that sequences are visited with probability given by its probability under the hidden Markov model multiplied by the amino acid sequence probability. This gives an evolutionary model on amino acid sequences and protein structures that accounts for neighbouring dependencies between adjacent sites and introduces temporal evolutionary dependencies between proteins.

The evolutionary dependencies between structures is introduced via the hidden states, thus avoiding having to directly implement an evolutionary process on structure, which is cumbersome given the continuous nature of the structural observations. We have previously developed a continuous diffusion process on angles for modelling protein dihedral angles, (García-Portugués *et al.*, 2018; Golden *et al.*, 2017), however, protein backbone angles and bond lengths do not evolve in a continuous fashion, rather they are expected to ‘jump’ when changes occur. The hidden states capture this jump behaviour.

Note that proteins P_a and P_b referred to in the ratio $\frac{\Pi(b)}{\Pi(a)}$ in Equation 6 always differ at exactly one site, implying that at most three terms in Equation 5

need to be considered when computing the ratio.

Additionally note, although the summation in Equation 6 appears to involve an exponential number of terms, most terms are equal to zero, except those that differ from P_a at one position. Furthermore, an amino acid transition and a hidden state transition are not permitted to occur simultaneously, further reducing the number of terms that are summed ($19n$ amino acid terms plus $(h - 1)n$ hidden state terms).

Stationary probability of proteins Following Choi *et al.* (2008), and by construction, the stationary probability of a protein a is given by:

$$p(P_a|\hat{\theta}) = \frac{\Pi(a) \prod_i \pi_{aa_i}^{h_i^a}}{\sum_k \Pi(k) \prod_i \pi_{aa_i}^{h_i^k}} \quad (7)$$

Time-reversibility Since U and V in (6) are symmetric matrices, i.e. $U_{a_i b_i} = U_{b_i a_i}$ and $V_{a_i b_i} = V_{b_i a_i}$, time-reversibility of the model holds, in other words:

$$p(P_a|\hat{\theta})M_{ab} = p(P_b|\hat{\theta})M_{ba}. \quad (8)$$

Time-reversibility implies that at any rooting of the tree can be used if the equilibrium probabilities are taken to be the initial probabilities (Felsenstein, 1981), which is indeed the case for our model.

Dataset likelihood The likelihood of a given dataset \mathcal{D}_d of proteins related by a tree \mathcal{T}_d consisting of a set of branch paths \mathcal{B}_d is given as follows:

$$p(\mathcal{D}_d|\mathcal{T}_d, \mathcal{B}_d, \hat{\theta}) = p(P_{root}) \prod_{b \in \mathcal{B}_d} p(X_s(t_{end}) | H_b(t_{end}), A_b(t_{end}), \hat{\theta}) \times [e^{R_{b_n}(t_{end}-t_n)} \prod_{k=1}^n e^{-R_{b_k b_k}(t_k-t_{k-1})} R_{b_{k-1} b_k}]. \quad (9)$$

The first term, $p(P_{root})$, is the probability of the protein, P_{root} , at the root of the tree. The outermost product is a product over branches in \mathcal{B}_d , where the first term is the likelihood of any structural observations at the tip of each branch. The terms in square brackets represent the likelihood of the the hidden

states and amino acids along a branch paths, as specified the rate matrix R . The first term in square parentheses is the probability that no events occur after the last event in a given branch path, whereas the second term is the probability of the events in a branch path and the waiting times between them.

2.2 Inference

2.2.1 Branch path inference: a phylogeny

Inference for a given dataset \mathcal{D}_d consists of sampling the set of branch paths, \mathcal{B}_d , conditional on the tree topology and branch lengths, \mathcal{T}_d , and model parameters $\hat{\theta}$:

$$\mathcal{B}_d \sim p(\mathcal{B}_d|\mathcal{D}_d, \mathcal{T}_d, \hat{\theta}). \quad (10)$$

To sample this distribution, for each site i , Felsenstein's algorithm was used to calculate likelihoods in a forward pass up the tree, followed by a backwards sampling pass down the tree to propose new hidden node states at the tip of each branch. The hidden state rate matrices used were conditional on the amino acid branch paths at site i and hidden state branch paths at site $i - 1$ and $i + 1$.

Conditional upon proposed internal node states, modified rejection sampling was used to sample hidden state branch paths using the parent branch's proposed hidden node tip state as the start state and the current branches proposed hidden node tip state as the end state.

The proposed branch paths for site i were then accepted or rejected using the Metropolis-Hastings ratio together with proposal ratio.

An analogous algorithm was used for sampling the amino acid branch paths, where the amino acid rate matrices used were conditional on the hidden states branch paths at each site i .

2.2.2 Inference: a single protein

Inference for a dataset consisting of a single protein P_a is much simpler, given that the protein is assumed to be drawn from the stationary distribution of the

model, which is given by (7)

$$H_a \sim p(H_a|A_a, \hat{\theta}). \quad (11)$$

This distribution can be sampled exactly using the forward-filtering backward-sampling algorithm for HMMs (Frühwirth-Schnatter, 1994) in $\mathcal{O}(h^2|P_a|)$ computational time. Note that the amino acid sequence, A_a , is typically observed. However, regardless of which combinations of A_a and X_a are observed it remains possible to use the forward-filtering backward-sampling algorithm to efficiently sample H_a .

2.2.3 Backbone structure inference

Branch path inference gives the distribution of $S_b(t) = (H_b(t), A_b(t))$ at every point in time t along a branch b . Conditioned on $H_b(t)$ and $A_b(t)$ the angles and bond lengths, $X_b(t)$, comprising the backbone structure can be trivially sampled using Equations 1-4. The posterior marginal $p(X_b(t)|\mathcal{D}_b, \mathcal{T}_b)$ is therefore obtained by first sampling $H_b(t)$ and $A_b(t)$:

$$(H_b(t), A_b(t)) \sim p(H_b(t), A_b(t)|\mathcal{D}_b, \mathcal{T}_b), \quad (12)$$

followed by sampling $X_b(t)$ conditional on $H_b(t), A_b(t)$:

$$X_b(t) \sim p(X_b(t)|H_b(t), A_b(t)). \quad (13)$$

2.3 Model training

2.3.1 Datasets

2.3.2 Model estimation

Stochastic EM (StEM, Gilks *et al.* (1995)) was used to train the model. StEM is a stochastic version of the well known Expectation-Maximization algorithm (Gilks *et al.*, 1995). Its distinguishing feature is that the E-step consists of filling in the values of the latent variables using sampling. Only a single value is sampled. StEM is attractive due to its computational efficiency and its tendency to avoid getting stuck in local minima (Gilks *et al.*, 1995).

Sampling was used in the E-step to sample branch paths and times. In other words, at iteration k for each dataset, d , consisting of an aligned set of proteins, \mathcal{D}_d , and a corresponding set of branch paths, \mathcal{B}_d , we draw samples, from the following joint-distribution:

$$Z_d^{(k)} \sim p(\mathcal{B}_d|\mathcal{D}_d, \Psi^{(k)}).$$

In the M-step the samples from the previous E-step, were used to update the hidden node parameters ($\hat{\Psi}$) using efficient sufficient statistics (ESSs).

3 Results and Discussion

3.1 Benchmarks of ancestral sequence reconstruction

4 Conclusions

5 Software availability

Julia code (compatible with Windows and Linux) is available at: <https://github.com/michaelgoldendev/protein-evolution>

6 Acknowledgements

MG is supported by the ERC under the European Union’s Seventh Framework Programme (FP7/2007-2013)/ERC grant agreement no. 614725-PATHPHYLODYN.

References

- Choi, S. C., Redelings, B. D., and Thorne, J. L. 2008. Basing population genetic inferences and models of molecular evolution upon desired stationary distributions of DNA or protein sequences. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 363(1512): 3931–3939.

Table 1: Composition of training datasets

Category	Number in category	Amino acid observations	Structural observations
One sequence (no structure)	310	84,806	0
Two sequences (no structures)	208	138,892	0
Three or more sequences (no structures)	118	390,913	0
One structure (one corresponding sequence)	4,565	1,259,635	1,259,342
One structure (two or more sequences)	192	147,679	58,887
Two structures (two or more sequences)	81	72,331	50,982
Three or more structures (three or more sequences)	18	44,273	20,745
Total	5,492	2,138,529	1,389,956

Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of molecular evolution*, 17(6): 368–376.

Frühwirth-Schnatter, S. 1994. Data augmentation and dynamic linear models. *Journal of time series analysis*, 15(2): 183–202.

García-Portugués, E., Golden, M., Sørensen, M., Mardia, K. V., Hamelryck, T., and Hein, J. 2018. Toroidal diffusions and protein structure inference. In C. Ley and T. Verdebout, editors, *Applied Directional Statistics: Modern Methods and Case Studies*. Chapman and Hall/CRC.

Gilks, W. R., Richardson, S., and Spiegelhalter, D. 1995. *Markov chain Monte Carlo in practice*. CRC press.

Golden, M., García-Portugués, E., Sørensen, M., Mardia, K. V., Hamelryck, T., and Hein, J. 2017. A generative angular model of protein structure evolution. *Molecular Biology and Evolution*, 34: msx137.

Parsons, J., Holmes, J. B., Rojas, J. M., Tsai, J., and Strauss, C. E. 2005. Practical conversion from torsion space to cartesian space for in silico protein synthesis. *Journal of computational chemistry*, 26(10): 1063–1068.