

Global multivariate model learning from hierarchically correlated data

Edwin Rodriguez Horta,^{1,2} Alejandro Lage,² Martin Weigt,¹ and Pierre Barrat-Charlaix^{3,*}

¹*Sorbonne Université, CNRS, Institut de Biologie Paris-Seine,
Laboratoire de Biologie Computationnelle et Quantitative – LCQB, Paris, France*

²*Group of Complex Systems and Statistical Physics,*

Department of Theoretical Physics,

University of Havana, Havana, Cuba

³*Biozentrum, Universität Basel, Basel, Switzerland*

Abstract

Inverse statistical physics aims at inferring models compatible with a set of empirical averages estimated from a high-dimensional dataset of independently distributed equilibrium configurations. However, in several applications such as biology, data result from stochastic evolutionary processes, and the relation between configurations is characterized by a hierarchical correlation structure, typically represented by a tree. In turn, empirical averages of observables superpose intrinsic signals related to the equilibrium distribution of the studied system, and spurious historical (or phylogenetic) signals characterizing the specific correlated data-generating process. The naive application of inverse statistical physics techniques therefore leads to systematic biases and an effective reduction of the sample size. To advance on the currently open task of extracting intrinsic signals from correlated data, we study a system described by a multivariate Ornstein-Uhlenbeck process defined on a finite tree. Using a Bayesian framework, we can disentangle covariances in the data, which are a result of their multivariate Gaussian equilibrium distribution, from those resulting from the historical correlations. Our approach leads to a clear gain in accuracy in the inferred equilibrium distribution, which corresponds to an effective two- to fourfold increase in sample size.

* Correspondence to: Pierre Barrat-Charlaix, pierre.barrat@unibas.ch

I. INTRODUCTION

With the emergence of large, high-dimensional datasets for complex systems across disciplines, methods of *inverse statistical physics* have seen rapidly growing interest during the last years [1]. In the most standard setting, the data provide observational samples of the “microscopic” degrees of freedom of the system under study – this can be biological sequences [2, 3], firing patterns of neurons [4, 5], individuals in animal ~~or human groups,~~ ~~stock prices etc. ...~~ groups [6, 7], stock markets [8, 9] etc. Within a static modeling approach, frequently based on the maximum-entropy principle [10], data \vec{x} are assumed to be generated independently from some unknown probability distribution $P(\vec{x})$. This distribution describes the underlying interaction patterns between the observed degrees of freedom, and has to be learned from data to unveil the rules governing the system. In more rare cases where data correspond to observed time series, theoretical and algorithmic development is much less advanced than in the case of independent static data ~~.[1].~~

~~One of the biggest application areas of inverse statistical mechanics is~~ Here we address a different case, the inverse problem for high-dimensional data, which show hierarchical correlations due to a data-generating process defined via a branching process. The motivation for this purely methodological study comes from the modeling of biological ~~processes.~~ ~~These applications are fuelled by the large amount of available data resulting from the impressive progress in experimental techniques in biology. This is especially visible in the case of biological sequences, with databases now harboring a vast amount of high-quality DNA or protein sequences [11, 12]. A common idea in this context is that it is possible to use characteristics of genes or organisms related by a common ancestry – called *homologous* – to construct models of the selection acting on them. A successful example in this regard is the representation of protein sequences by probabilistic models in the~~ sequence data, which is one of the major fields of application of inverse statistical mechanics [2, 3], but even more of other bioinformatic inference approaches, cf. [13, 14]. Sequence data are the result of natural evolution, i.e. a branching process initiating in some unknown common ancestor of extant and thus observable sequences [11, 12, 15]. The prototypical datasets in this case are multiple-sequence alignments (MSA), with lines being a so-called ~~DCA method [2, 3].~~ ~~We will use this type of application to motivate our otherwise purely methodological study.~~

~~When modeling homologous genes or organisms, one has to distinguish between two~~
~~complementary types~~ homologous, i.e. evolutionarily related sequences, and columns specific
positions deriving from some common ancestral position. Such MSA contain at least two
kinds of information:

- *Phylogenetic information:* the distances between ~~gene sequences or traits in organisms~~
sequences carry information about the evolutionary time since their common ancestor.
Using phylogenetic methods [14, 16] we may reconstruct the evolutionary history of ~~a~~
~~set of genes or species, and also the sequences or traits of their ancestors~~ our dataset,
represented by a phylogenetic tree.
- *Co-evolutionary information:* positions in a sequence ~~or traits in organisms~~ typically
do not evolve independently, but rather in a correlated way. This co-evolution carries
important information about the selection forces acting on evolving entities. This fact
has been extensively studied in the case of protein sequences, and used to predict
structure, mutational landscapes or networks of interacting proteins [2, 3].

These two types of information are contained in two complementary features of the
data: phylogenetic inference is based on the comparative analysis of different ~~species or~~
~~genes~~ sequences, while co-evolutionary information is contained in the correlation of different
~~trait values in the same gene or organism.~~

~~Models built using biological data usually ignore one of the two types of information.~~
~~For instance, columns of the MSA. Approaches aiming at one type of information typically~~
~~neglect the other one:~~ inference of phylogenies ~~typically~~ generally assumes that all positions
in a sequence evolve independently, while co-evolutionary models of proteins assume that ~~their~~
~~sequences~~ sequences in the MSA are independently distributed. This choice is motivated
by the fact that taking the two types of correlations into account, *i.e.* through time with
phylogeny and accross trait values for co-evolution, results in very hard inference procedures,
cf. [17, 18]. However, this can lead to biases in the model parameters: it has for instance been
shown that phylogenetic relations between protein sequences induce non-trivial correlations
that are not related to protein function [19, 20].

~~Accounting for both~~ We have used phylogeny and co-evolution ~~and phylogeny in models~~
~~requires one to understand how interdependent traits under selection evolve along a tree.~~
~~A of biological sequences to explain our motivation, but the underlying problem is much~~

more general. Instead of sequences of discrete characters, like amino acids or nucleotides, we may consider continuous phenotypic traits. The branching process is not necessarily the phylogeny of species, but it may be the genealogy of populations of the same species, or other branching processes like epidemics spreading or geographic migration.

Here we use a simple and very general model for the temporal evolution of correlated variables: a historically well-known way to represent such processes is to use Ornstein-Uhlenbeck dynamics (OU), which models ~~traits as a Gaussian vector configurations as~~ Gaussian vectors evolving in a quadratic potential that represents selection forces [16, 21, 22]. OU processes are commonly used in the field of phylogenetic comparative methods (PCM) [23, 24]. This modeling approach is *a priori* limited to continuous traits, but could potentially be used for protein sequences combined with a continuous-variable approximation, that has successfully been used in the past [25–27]. In this context, the equilibrium distribution reached by the OU process represents the probability distribution given by the DCA method, which can be used to predict non-trivial structural contacts in the protein fold, effects of amino-acid mutations or even designing novel functional sequences [28–30].

In this work, we are interested in constructing an inference method for parameters of an OU process from data correlated through a tree. Our approach is purely methodological, and the data can represent any set of continuous phenotypic traits, *e.g.* from different organisms, with the tree indicating the phylogenetic relations between data points. Inferred parameters then represent the selection forces without biasing effects from the phylogeny. The manuscript is divided as follows: we first review in section II the main characteristics of the multivariate OU process. We then describe the setting of the inference problem that we want to solve in section III A, propose a solution in sections III B and III C. Finally, we present results obtained on simulated data in section IV, with the context of pairwise models of protein sequences in mind.

II. THE MULTIVARIATE ORNSTEIN-UHLENBECK PROCESS

We consider a system characterized by L continuous degrees of freedom and whose state is fully described by an L -dimensional vector $\vec{x} \in \mathbb{R}^L$. These degrees of freedom can be continuous phenotypic traits of some living organism, or the sequence of a gene or a protein if

115 a continuous approximation is made. At equilibrium, \vec{x} is assumed to be normally distributed,

$$P_{eq}(\vec{x}) = \frac{1}{Z(\mathbf{J})} \exp\left\{-\frac{1}{2}\vec{x}^T \mathbf{J} \vec{x}\right\}, \quad (1)$$

116 where \mathbf{J} is the symmetric, positive definite *coupling matrix* and $Z(\mathbf{J}) = \sqrt{(2\pi)^L / \det \mathbf{J}}$ is
 117 the normalization constant; the means of all components of \vec{x} are set to zero without loss
 118 of generality. We are interested in inferring the coupling matrix from a given amount of
 119 observed states \vec{x} of the system. If these observations were independent from each other,
 120 due to the simple Gaussian form of Eq. (1), \mathbf{J} would simply be equal to the inverse of the
 121 empirical *covariance matrix* of the data, written $\mathbf{C} = \mathbf{J}^{-1}$.

122 However, we consider the case where observations are not independent. On the contrary,
 123 they result from a dynamical process taking place during a finite amount of time, and different
 124 data-points are therefore correlated to each other. This dynamical process is described below.

125 We suppose that the considered system evolves according to the following Langevin
 126 equation

$$\gamma^{-1} \frac{d\vec{x}}{dt} = -\mathbf{J}\vec{x} + \vec{\xi}(t). \quad (2)$$

127 Here, $\vec{\xi}(t)$ is a vector of uncorrelated white noise, and γ^{-1} is the characteristic timescale
 128 governing the dynamics. In short, Eq. (2) states that the system described by \vec{x} undergoes
 129 Brownian motion in a quadratic energy landscape characterized by the coupling matrix \mathbf{J} .

130 We are not interested in \vec{x} directly, but rather in its probability distribution $P(\vec{x}|\vec{x}_0, \Delta t)$,
 131 *i.e.* in the probability to find the system in state \vec{x} knowing it was in state \vec{x}_0 some time Δt
 132 in the past. The Fokker-Planck equation corresponding to Eq. (2) is straightforward to write,

$$\gamma^{-1} \partial_t P = \left(- \sum_{a,b=1}^L \frac{\partial}{\partial x_a} J_{ab} x_b + \sum_{a=1}^L \frac{\partial^2}{\partial x_a^2} \right) P, \quad (3)$$

133 where the parenthesized expression on the right hand side is understood as an operator acting
 134 on P . The solution to Eq. (3) is a multivariate normal distribution [31]:

$$P(\vec{x}|\vec{x}_0, \Delta t) = [(2\pi)^N \det \mathbf{\Sigma}]^{-1/2} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu})^T \mathbf{\Sigma}^{-1}(\vec{x} - \vec{\mu})\right\}, \quad (4)$$

135 where we introduce the matrices $\mathbf{\Sigma}$ and $\mathbf{\Lambda}$ as well as the vector $\vec{\mu}$ as

$$\mathbf{\Lambda} = e^{-\gamma \mathbf{J}}, \quad \vec{\mu} = \mathbf{\Lambda}^{\Delta t} \vec{x}_0, \quad \mathbf{\Sigma} = \mathbf{J}^{-1}(\mathbb{1} - \mathbf{\Lambda}^{2\Delta t}). \quad (5)$$

136 Eqs. (4) and (5) define a multivariate *Ornstein-Uhlenbeck* (OU) process.

Note that since matrix $\mathbf{\Lambda}$ is an exponential of \mathbf{J} , it is symmetric, has strictly positive eigenvalues and commutes with \mathbf{J} . We also underline that $\mathbf{\Sigma}$ and $\vec{\mu}$ depend on Δt , although this dependence is not explicitly written in our notation to make it less heavy. By taking $\gamma\Delta t \gg 1$ and using the fact that \mathbf{J} has strictly positive eigenvalues, one immediately recovers Eq. (1), meaning that the OU process converges to the desired equilibrium distribution.

We can compute the joint distribution of two configurations \vec{x}_1 and \vec{x}_2 separated by a time Δt by multiplying Eqs. (1) and (4),

$$P(\vec{x}_1, \vec{x}_2 | \Delta t) = P(\vec{x}_1 | \vec{x}_2, \Delta t) \times P_{eq}(\vec{x}_2) \propto \exp \left\{ -\frac{1}{2} (\vec{x}_1^T \mathbf{\Sigma}^{-1} \vec{x}_1 + \vec{x}_2^T \mathbf{\Sigma}^{-1} \vec{x}_2 - 2\vec{x}_1^T \mathbf{\Lambda}^{\Delta t} \mathbf{\Sigma}^{-1} \vec{x}_2) \right\}. \quad (6)$$

This equation illustrates the *time reversibility* of the OU process. Indeed, the distribution is symmetric in \vec{x}_1 or \vec{x}_2 and does not depend on which configuration came first.

Equation (6) allows for computing the joint covariance of the correlated equilibrium configurations \vec{x}_1 and \vec{x}_2 . The probability distribution in Eq. (6) is normal with an inverse covariance matrix defined by blocks: $\mathbf{\Sigma}$ on the diagonal and $-\mathbf{\Lambda}^{\Delta t} \mathbf{\Sigma}$ off-diagonal. By inverting this block matrix, given that $\mathbf{\Lambda}$ and $\mathbf{\Sigma}$ commute and are invertible, one obtains the following covariance:

$$\langle \vec{x}_1 \vec{x}_2^T \rangle_{\Delta t} = \mathbf{\Lambda}^{\Delta t} \mathbf{J}^{-1} = \mathbf{\Lambda}^{\Delta t} \mathbf{C}. \quad (7)$$

Eq. (7) allows us to readily distinguish two regimes. Let us call ρ_a the eigenvalues of \mathbf{J} . The eigenvalues of $\mathbf{\Lambda}^{\Delta t} \mathbf{C}$ are then equal to $\rho_a^{-1} e^{-\gamma \rho_a \Delta t}$. Since all ρ_a are positive, the eigenvalues of $\mathbf{\Lambda}^{\Delta t} \mathbf{C}$ vanish exponentially over time. The slowest timescale of exponential decay is set by $\tau_c^{-1} = \gamma \rho_{min}$, with ρ_{min} being the smallest eigenvalue of \mathbf{J} . Thus, for $\Delta t / \tau_c \gg 1$, \vec{x}_1 and \vec{x}_2 are uncorrelated. If this is verified for all pairs of observations \vec{x}_i and \vec{x}_j , the regime is that of *uncorrelated* data – the inference of \mathbf{J} can simply be performed by inverting the empirical covariance matrix extracted from the data. Inversely, for $\Delta t / \tau_c \ll 1$, \vec{x}_1 and \vec{x}_2 are highly correlated, defining a *strongly correlated* regime. It should be noted that for $\Delta t = 0$, the joint correlation matrix of \vec{x}_1 and \vec{x}_2 becomes non invertible, and Eq. (7) becomes irrelevant. Actually, \vec{x}_1 and \vec{x}_2 coincide at that point, *i.e.* we have $P(\vec{x}_1, \vec{x}_2 | \Delta t = 0) = P_{eq}(\vec{x}_1) \times \delta(\vec{x}_1 - \vec{x}_2)$ using the L -dimensional Dirac distribution.

III. METHODS

A. Statement of the problem

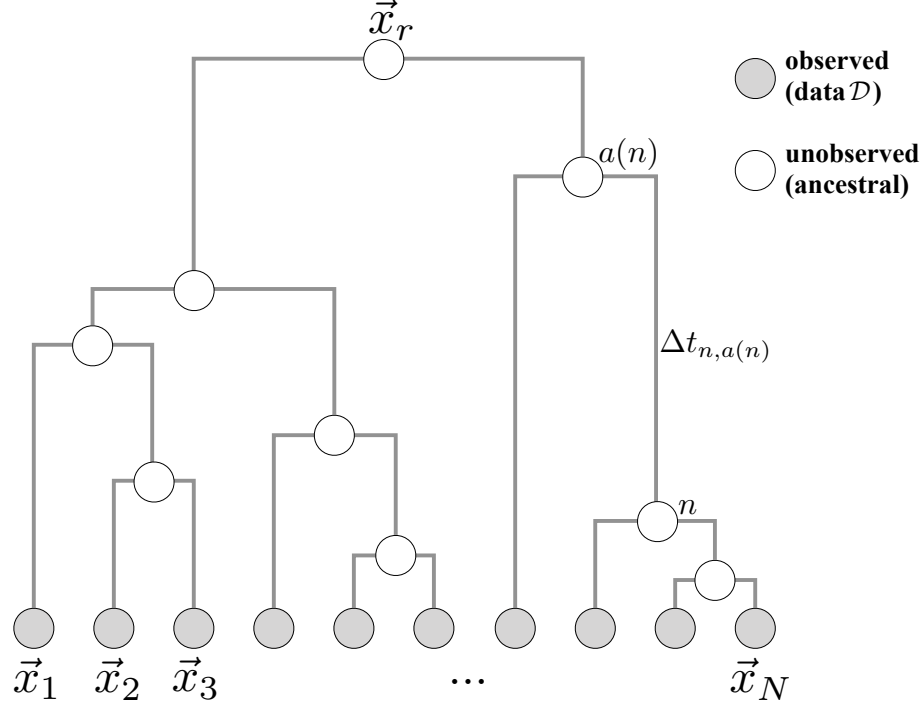


FIG. 1. Schematic representation of a tree \mathcal{T} underlying the data generating process. The process starts at the root node r with a configuration \vec{x}_r sampled from $P_{eq}(\vec{x}_r)$. The dynamics consist in independent realizations of the OU process on all branches from ancestral nodes $a(n)$ to child nodes n over times corresponding to the branch length $\Delta t_{n,a(n)}$, initialized in the ancestral configuration $\vec{x}_{a(n)}$. The observable data only consist of configurations of the leaf nodes (grey circles in the figure), while configurations of ancestral nodes remain unknown. There are no restrictions on the topology of tree \mathcal{T} and the length of the branches.

The problem discussed here is the inference of the probability distribution describing samples that are hierarchically correlated by a tree, cf. Fig. 1. Formally, we assume that the data consists of N real-valued vectors of length L , denoted $\{\vec{x}_i\} \in \mathbb{R}^L$ with $i = 1, \dots, N$. Taken individually, we assume that the \vec{x}_i are distributed according to Eq. (1), *i.e.* according to a multivariate Gaussian of zero mean and covariance \mathbf{C} . By construction, the equilibrium covariance between any pair of elements of a given vector $\vec{x} = (x^1, \dots, x^L)$ is given by the

inverse of the coupling matrix: $\langle x^a x^b \rangle - \langle x^a \rangle \langle x^b \rangle = \mathbf{C}_{ab} = (\mathbf{J}^{-1})_{ab}$ for all $a, b = 1, \dots, L$. This implies that inferring the coupling matrix defining the probability distribution amounts to finding the *equilibrium* covariance matrix \mathbf{C} .

However, this covariance cannot be directly measured as we consider observations that are not independently distributed. Instead, the set of measured configurations $\{\vec{x}_i\}_{i=1,\dots,N}$ is the result of an Ornstein-Uhlenbeck (OU) process taking place on a tree \mathcal{T} , as is illustrated in Fig. 1:

- The process starts at the root node r with a state vector \vec{x}_r drawn from the equilibrium distribution P_{eq} .
- On each branch $(n, a(n))$ of length $\Delta t_{n,a(n)}$ connecting node n with its ancestral node $a(n)$, the dynamics follow Eq. (2), starting from initial condition $\vec{x}_{a(n)}$, and running for time $\Delta t_{n,a(n)}$. In other words, given the state $\vec{x}_{a(n)}$ of the ancestral node, \vec{x}_n is sampled from $P(\vec{x}_n | \vec{x}_{a(n)}, \Delta t_{n,a(n)})$, see Eq. (4)
- As a consequence, OU processes on branches stemming from common ancestral node evolve independently, but from an identical initial condition.
- Observed data vectors correspond to the states of the leaves of the tree at the end of this process. The states of the internal nodes are not part of the observed data and remain unknown.

This process is thought to represent the evolution of biologic traits along a phylogenetic tree, with the leaf nodes corresponding to traits observed in today's species. Note that due to the reversible nature of our OU process, the joint probability of any pair of leaf configurations \vec{x}_i and \vec{x}_j , with $i, j \in \{1, \dots, N\}$, is given by $P(\vec{x}_i, \vec{x}_j | \Delta t_{ij})$ (Eq. (6)), with Δt_{ij} denoting the total branch length of the path connecting i and j in the tree.

The OU process is characterized by the quadratic potential $\mathbf{J} = \mathbf{C}^{-1}$ and the rate γ . Hence, the joint statistics of the leaf configurations $\{\vec{x}_i\}_{i=1,\dots,N}$ (*i.e.* the data) is fully determined by \mathbf{C} , γ , and the tree \mathcal{T} . The aim of this work is to derive a method for inferring the most likely values of \mathbf{C} and γ given the knowledge of the data $\mathcal{D} = \{\vec{x}_i\}_{i=1,\dots,N}$ and the underlying tree \mathcal{T} . We consider here that both the topology and the branch lengths of \mathcal{T} are known.

This problem shows two notable extreme cases: The first one is the case where the typical branch length of the tree is short compared to the timescales of the OU process. As a consequence, leaf configurations are close to identical to the root, *i.e.* $\vec{x}_i \simeq \vec{x}_r$, and the inference of \mathbf{C} becomes impossible. The second one is the opposite case where the typical branch length of the tree is long compared to the longest timescale of the OU process τ_c . In this case, the configuration of a child node is close to independent from that of its ancestor, and leaf configurations can be considered as independent samples from the equilibrium distribution P_{eq} . \mathbf{C} can then be readily estimated by computing the empirical covariance matrix. We are interested here in the intermediate regime where substantial tree-mediated correlations between data make it impossible to simply estimate \mathbf{C} with the empirical covariance, but the depth of the tree introduces enough variability in the data for one to hope of reconstructing the energy potential \mathbf{J} .

We adopt a Bayesian inference approach by writing the probability of a given set of parameters $\{\mathbf{C}, \gamma\}$ given the data $\{\mathcal{D}, \mathcal{T}\}$ using Bayes' equation

$$P(\mathbf{C}, \gamma | \mathcal{D}, \mathcal{T}) \propto P(\mathcal{D} | \mathbf{C}, \gamma, \mathcal{T}) \cdot P(\mathbf{C}, \gamma), \quad (8)$$

with the proportionality constant not depending on the parameters $\{\mathbf{C}, \gamma\}$. Here, $P(\mathbf{C}, \gamma)$ can be any arbitrarily chosen prior distribution. The difficulty in Eq. (8) lies in the estimation of the likelihood $P(\mathcal{D} | \mathbf{C}, \gamma, \mathcal{T})$, *i.e.* of the joint probability of the datapoints $\mathcal{D} = \{\vec{x}_i\}_{i=1, \dots, N}$ for an OU process given by its parameters $\{\mathbf{C}, \gamma\}$ and the tree \mathcal{T} . We detail the computation of this probability in the following section.

B. Calculation of the likelihood

The joint distribution of two configurations \vec{x}_1 and \vec{x}_2 separated by time Δt is given by Eq. (6) and corresponds to a joint normal distribution. This means that the vector $\vec{X} = [\vec{x}_1, \vec{x}_2]$, *i.e.* the concatenation of vectors \vec{x}_1 and \vec{x}_2 , follows a normal distribution with zero mean and variance described above in Eqs. (5). Of importance here is that this property of the OU process can be extended to the joint distribution of any subset of nodes in a tree. In other words, if we now define $\vec{X} = [\vec{x}_1, \dots, \vec{x}_N]$ to be the concatenation of all configurations in our dataset \mathcal{D} , we can write the distribution of \vec{X} as

$$P(\vec{X} | \mathbf{C}, \gamma, \mathcal{T}) = ((2\pi)^{LN} \det \mathbb{G})^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \vec{X}^T \mathbb{G}^{-1} \vec{X} \right\}, \quad (9)$$

226 where \mathbb{G} is the *joint covariance matrix* and depends on the tree as well as on \mathbf{C} and γ .

227 The joint covariance matrix is a matrix of dimension $(L \cdot N) \times (L \cdot N)$, built by $N \times N$
 228 blocks of size $L \times L$ with entries

$$\mathbb{G}_{ij}(a, b) = \langle x_i^a x_j^b \rangle - \langle x_i^a \rangle \langle x_j^b \rangle, \quad i, j \in \{1, \dots, N\}; a, b \in \{1, \dots, L\}, \quad (10)$$

229 where the (zero) marginals $\langle x_i^a \rangle$ and $\langle x_j^b \rangle$ are explicitly written for clarity. Each block \mathbb{G}_{ij} is
 230 describing the connected correlations between two data vectors \vec{x}_i and \vec{x}_j , which are separated
 231 by time Δt_{ij} , resulting as the sum of all branch lengths of the path connecting i and j on
 232 tree \mathcal{T} . Because the OU process is time reversible, we can directly apply Eq. (7) and give all
 233 blocks of \mathbb{G} in closed form,

$$\mathbb{G}_{ij} = \begin{cases} \mathbf{C} & \text{if } i = j \\ \mathbf{\Lambda}^{\Delta t_{ij}} \mathbf{C} & \text{otherwise,} \end{cases} \quad (11)$$

234 using the (currently unknown) covariance matrix \mathbf{C} of a single equilibrium vector \vec{x} . We
 235 remind here that $\mathbf{\Lambda} = e^{-\gamma \mathbf{C}^{-1}}$ depends only on γ and \mathbf{C} , and commutes with \mathbf{C} . As a direct
 236 consequence, all blocks \mathbb{G}_{ij} commute with each other and with \mathbf{C} .

237 Eq. (9) allows us to compute the log-likelihood of the data \vec{X} as a function of \vec{X} itself
 238 and of the joint covariance matrix. Indeed, taking its logarithm immediately gives

$$\mathcal{L}_{\mathcal{D}}(\mathbb{G}) = -\frac{1}{2} \log \det \mathbb{G} - \frac{1}{2} \vec{X}^T \mathbb{G}^{-1} \vec{X} + \text{const} , \quad (12)$$

239 but this expression is impractical for any numerical evaluation due to the large dimension of
 240 \mathbb{G} . However, the particular block structure of \mathbb{G} described in Eq. (11) allows us to simplify
 241 the expression. To do so, we first introduce the eigenvalues and eigenvectors $\{\rho_a, \vec{s}_a\}$ of \mathbf{C}^{-1} ,
 242 where the index a runs from 1 to L and vectors \vec{s}_a are of dimension L . By definition, we
 243 have $\rho_a > 0$ for all a . Using now Eq. (11), we immediately see that the vectors \vec{s}_a are also
 244 eigenvectors of the individual blocks \mathbb{G}_{ij} with eigenvalues $z(\rho_a, \Delta t_{ij})$ where we introduced

$$z(\rho_a, \Delta t_{ij}) = \rho_a^{-1} e^{-\gamma \rho_a \Delta t_{ij}} . \quad (13)$$

245 By convention, $\Delta t_{ii} = 0$ and the diagonal blocks are thus included via $z(\rho_a, \Delta t_{ii}) = \rho_a^{-1}$.

246 As the next step, we introduce $N \times N$ -dimensional matrices $\mathbf{G}^a, a = 1, \dots, L$, with elements

$$\mathbf{G}_{ij}^a = z(\rho_a, \Delta t_{ij}) , \quad 1 \leq i, j \leq N . \quad (14)$$

247 In other words, for a given index $1 \leq a \leq L$, \mathbf{G}^a is the matrix built by replacing all blocks
 248 of \mathbb{G} by their respective a th eigenvalue. Matrices \mathbf{G}^a are symmetric and have their own
 249 eigenmodes, that we denote by $\{\lambda_{ka}, \vec{u}_{ka}\}_{k=1, \dots, N}$.

250 To obtain the eigenmodes of the joint covariance matrix \mathbb{G} as a function of the \vec{s}_a and
 251 \vec{u}_{ka} , we construct the direct product of vectors \vec{s}_a and \vec{u}_{ka} , defining vectors \vec{S}_{ka} of dimension
 252 $L \times N$:

$$\begin{aligned} \vec{S}_{ka} &= \vec{u}_{ka} \otimes \vec{s}_a \\ &= [u_{ka}^1 \cdot \vec{s}_a, \dots, u_{ka}^N \cdot \vec{s}_a]. \end{aligned} \tag{15}$$

253 The i th block vector of \vec{S}_{ka} will thus be written as $\vec{S}_{ka}^i = u_{ka}^i \cdot \vec{s}_a$. We can now show that \vec{S}_{ka}
 254 are eigenvectors of matrix \mathbb{G} by considering the i th block vector of the product $\mathbb{G} \cdot \vec{S}_{ka}$:

$$\begin{aligned} (\mathbb{G} \cdot \vec{S}_{ka})^i &= \sum_{j=1}^N \mathbb{G}_{ij} u_{ka}^j \cdot \vec{s}_a \\ &= \sum_{j=1}^N z(\rho_a, \Delta t_{ij}) u_{ka}^j \cdot \vec{s}_a \\ &= (\mathbf{G}^a \cdot \vec{u}_{ka})^i \cdot \vec{s}_a \\ &= \lambda_{ka} (u_{ka}^i \cdot \vec{s}_a) \\ &= \lambda_{ka} \vec{S}_{ka}^i. \end{aligned} \tag{16}$$

255 We have first used the fact that \vec{s}_a is an eigenvector of \mathbb{G}_{ij} , then the definition of \mathbf{G}^a , and
 256 finally the fact that \vec{u}_{ka} is an eigenvector of \mathbf{G}^a . This demonstrates that the eigenmodes
 257 of \mathbb{G} are $\{\lambda_{ka}, \vec{S}_{ka}\}$ with $1 \leq k \leq N$ and $1 \leq a \leq L$. Since \mathbb{G} is the covariance matrix of a
 258 Gaussian distribution, we conclude the λ_{ka} to be strictly positive.

259 Note that this decomposition of the eigenvectors leads to a drastic decrease in computa-
 260 tional complexity for diagonalizing \mathbb{G} (at given \mathbf{C} , γ and \mathcal{T}), and in consequence also for
 261 calculating the likelihood according to Eq. (12), which depends on the inverse covariance
 262 matrix \mathbb{G}^{-1} . Matrix \mathbb{G} has linear dimension LN , so the numerical diagonalization or inversion
 263 takes time $\mathcal{O}((LN)^3)$. This is hardly achievable for systems of realistic length L of the state
 264 vector, and sufficient number N of data points for model learning. Following the above
 265 description, we need to first diagonalize \mathbf{C}^{-1} (or equivalently \mathbf{C}), which requires time of
 266 $\mathcal{O}(L^3)$, followed by inversion of the L matrices \mathbf{G}^a , each one having linear dimension N . The
 267 total time complexity therefore results in $\mathcal{O}(L^3) + \mathcal{O}(L \cdot N^3)$, and the calculation can be

268 easily achieved even on a standard PC. This observation is essential for inference, since we
 269 need to redo this calculation for many realizations of \mathbf{C} and γ , in order to find the ones
 270 maximizing the likelihood given the data \mathcal{D} and the tree \mathcal{T} . As is shown in section SA 4,
 271 this calculation simplifies even more when considering a fully balanced and homogeneous
 272 tree. In this case, the matrices \mathbf{G}^a commute and can be diagonalized simultaneously and
 273 analytically for any value of ρ^a .

274 For the case of arbitrary trees, Eq. (12) can now be rewritten using the eigen-decomposition
 275 of \mathbb{G} :

$$\begin{aligned}\mathcal{L}_{\mathcal{D}}(\mathbb{G}) &= -\frac{1}{2} \sum_{k=1}^N \sum_{a=1}^L \log \lambda_{ka} - \frac{1}{2} \sum_{k=1}^N \sum_{a=1}^L \lambda_{ka}^{-1} (\vec{X} \cdot \vec{S}_{ka})^2 \\ &= -\frac{1}{2} \sum_{k,a} \left(\log \lambda_{ka} + \lambda_{ka}^{-1} \left(\sum_{i=1}^N u_{ka}^i \vec{x}_i \cdot \vec{s}_a \right)^2 \right).\end{aligned}\tag{17}$$

276 Eq. (17) expresses the likelihood as a function of \vec{u}_{ka} , λ_{ka} (resulting from the tree \mathcal{T} and
 277 given ρ^a) and \vec{s}_a (resulting from \mathbf{C}). However, the definition of \mathbf{G}^a in Eq. (14) makes clear
 278 that its eigenmodes $\{\lambda_{ka}, \vec{u}_{ka}\}$ depend only of the eigenvalues ρ_a of \mathbf{C}^{-1} , on γ , as well as of
 279 the structure of the tree through the quantities Δt_{ij} , although this dependence cannot be
 280 analytically expressed in a simple manner. This means that the likelihood in equation (17)
 281 is in fact a function of $\{\rho_a, \vec{s}_a\}$, *i.e.* the eigenmodes of \mathbf{C}^{-1} , of the time scale parameter γ
 282 and of the pairwise distances on the tree Δt_{ij} .

283 C. Maximizing the likelihood

284 As stated at the beginning of this section, our main task is to find the equilibrium
 285 covariance matrix \mathbf{C} that maximizes the likelihood of the data. We also need to find the
 286 optimal time scale γ . In Eq. ((17)), the likelihood is expressed as a function of γ and
 287 $\{\rho_a, \vec{s}_a\}$, *i.e.* the eigenvalues and eigenvectors of \mathbf{C}^{-1} , either directly or through the quantities
 288 $\{\lambda_{ka}, \vec{u}_{ka}\}$. We now attempt to maximize the likelihood with respect to the eigenmodes
 289 $\{\rho_a, \vec{s}_a\}$ and to the time scale γ .

290 In order to perform this optimization, we need to compute the gradient of the likelihood
 291 with respect to the eigenvectors $\{\vec{s}_a\}$. Since \mathbf{C}^{-1} is a symmetric matrix, its eigenvectors
 292 form an orthogonal basis of the vector-space of dimension L and their components cannot be
 293 changed independently. One possible parametrization for the $\{\vec{s}_a\}$ consists in using $L(L-1)/2$

scalar *Eulerian angles* $\{\theta_{\alpha\beta}\}$ with $1 \leq \alpha < \beta \leq L$ [32, 33]. With the L eigenvalues ρ_a , this results in $L(L+1)/2$ independent values that fully parametrize the $L(L+1)/2$ values of \mathbf{C}^{-1} . A second possibility, that we have found faster in practice, is to express the matrix of the $\{\vec{s}_a\}$ as the exponential of a skew-symmetric matrix, see section A 3 of the SM. However, this parametrization does not allow a simple analytical expression of the gradient of the likelihood, and we use it along with automatic differentiation [34]. For this reason, we use the Eulerian angles below to express the gradient of the likelihood.

As a first step, we need to compute the gradient of the likelihood $\mathcal{L}_{\mathcal{D}}(\mathbb{G})$ with respect to all parameters $\{\rho_a, \theta_{\alpha\beta}\}$ and γ . To make explicit the dependences of eigenvalues and eigenvectors of the matrices \mathbf{G}^a on these parameters, we introduce the notation $\vec{u}_k(\rho_a, \gamma) = \vec{u}_{ka}$ and $\lambda_k(\rho_a, \gamma) = \lambda_{ka}$. Note that from the definition of \mathbf{G}^a in Eq. (14), its eigenvalues and vectors depend only on the eigenvalues of \mathbf{C}^{-1} and not on its eigenvectors. In the same way, we will now write $\mathbf{G}(\rho_a, \gamma)$ instead of \mathbf{G}^a .

The gradient of the likelihood is obtained by differentiating Eq. (17) with respect to the parameters of interest. This gives us three equations:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \rho_a} = & -\frac{1}{2} \sum_{k=1}^N \left\{ \frac{\partial \lambda_k}{\partial \rho_a} \lambda_k^{-1} - \frac{\partial \lambda_k}{\partial \rho_a} \lambda_k^{-2} \left(\sum_{i=1}^N u_k^i \vec{x}_i \cdot \vec{s}_a \right)^2 \right. \\ & \left. + 2 \lambda_k^{-1} \left(\sum_{i=1}^N u_k^i \vec{x}_i \cdot \vec{s}_a \right) \left(\sum_{i=1}^N \frac{\partial u_k^i}{\partial \rho_a} \vec{x}_i \cdot \vec{s}_a \right) \right\}, \end{aligned} \quad (18)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_{\alpha\beta}} = \sum_{k=1}^N \lambda_k^{-1} \left(\sum_{i=1}^N u_k^i \vec{x}_i \cdot \vec{s}_a \right) \left(\sum_{i=1}^N u_k^i \vec{x}_i \cdot \frac{\partial \vec{s}_a}{\partial \theta_{\alpha\beta}} \right), \quad (19)$$

and

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \gamma} = & -\frac{1}{2} \sum_{k=1}^N \left\{ \frac{\partial \lambda_k}{\partial \gamma} \lambda_k^{-1} - \frac{\partial \lambda_k}{\partial \gamma} \lambda_k^{-2} \left(\sum_{i=1}^N u_k^i \vec{x}_i \cdot \vec{s}_a \right)^2 \right. \\ & \left. + 2 \lambda_k^{-1} \left(\sum_{i=1}^N u_k^i \vec{x}_i \cdot \vec{s}_a \right) \left(\sum_{i=1}^N \frac{\partial u_k^i}{\partial \gamma} \vec{x}_i \cdot \vec{s}_a \right) \right\}, \end{aligned} \quad (20)$$

The derivatives of $\vec{u}_k(\rho, \gamma)$ and $\lambda_k(\rho, \gamma)$ with respect to ρ can then be computed using the following equations [35]:

$$\frac{\partial \lambda_i(\rho, \gamma)}{\partial \rho} = \vec{u}_k(\rho, \gamma)^T \frac{\partial \mathbf{G}(\rho, \gamma)}{\partial \rho} \vec{u}_k(\rho, \gamma) \quad (21)$$

and

$$\frac{\partial \vec{u}_k(\rho, \gamma)}{\partial \rho} = \sum_{l \neq k} \left(\vec{u}_k(\rho, \gamma)^T \frac{\partial \mathbf{G}(\rho, \gamma)}{\partial \rho} \vec{u}_l(\rho, \gamma) \right) (\lambda_k(\rho, \gamma) - \lambda_l(\rho, \gamma))^{-1} \vec{u}_l(\rho, \gamma). \quad (22)$$

Equivalent equations can be written for their derivatives with respect to γ .

The computation of the gradient of \mathcal{L} for a given set of parameters $\{\rho_a, \theta_{\alpha\beta}\}$ then goes as follows. For each eigenvalue ρ_a , we compute and diagonalize matrix $\mathbf{G}(\rho_a)$ to obtain its eigenmodes $\vec{u}_k(\rho_a)$ and $\lambda_k(\rho_a)$. Using equations (21) and (22) and their equivalent form for γ , we also numerically compute their derivatives with respect to ρ_a and γ . This gives us all the quantities to estimate the gradient of \mathcal{L} with respect to ρ_a using equation (18).

The optimization is performed by a quasi-Newton method [36]. Details are presented in section A 5 of the SM.

IV. RESULTS

In order to evaluate our inference procedure, we generate artificial data corresponding to the process described in section III A. We first build a balanced binary tree \mathcal{T} with $2^9 = 512$ leaves. The length of each branch of \mathcal{T} is chosen from a uniform distribution in the interval $[0, 1]$. We also sample positive semi-definite coupling matrix \mathbf{J} of size $L \times L$ with $L = 4$ or $L = 10$, with entries normally distributed with mean $\mu_J = 0.8$ and $\sigma_J = 0.2$.

In the case of statistical models of protein sequence, a major achievement is the ability of pairwise models to predict contacts in the three-dimensional structure of the protein from an inferred coupling matrix. In order to replicate this setting and to perform interaction prediction, we randomly set to 0 off-diagonal elements of J with probability 0.7, resulting in a sparsified coupling matrix of approximate density 0.3. Zero elements of J correspond to variables that do not interact, in analogy to non-contacts in the case of an application to protein sequences.

In order to investigate the different regimes of tree-induced correlation, we vary the parameter γ around a reference timescale γ_d defined as follows:

$$\gamma_d = \frac{1}{\Delta t_{av} \rho_{min}} \quad (23)$$

where Δt_{av} is the average branch length separating two leaves of \mathcal{T} . For $\gamma \gg \gamma_d$, leaf configurations are on average well decorrelated, whereas for $\gamma \ll \gamma_d$ all leaves will be strongly correlated. By simulating data using different γ in the range $[10^{-2}, 2] \cdot \gamma_d$, we investigate

all relevant temporal regimes. For each value of γ , we then sample configurations of leaves of \mathcal{T} using the process described in section A 1 of the supplementary material. To avoid statistical noise when assessing the quality of our inference, we repeat the sampling of leaf configurations 100 times for each value of γ .

For each repetition of the sampling process, we perform our maximum likelihood procedure and obtain an inferred covariance matrix \mathbf{C}_{max} . As a means of comparison, we also compute the empirical covariance matrix \mathbf{C}_{emp} as if leaf configurations were independent. Fig. 2 shows the Pearson correlation between the real covariance matrix $\mathbf{C} = \mathbf{J}^{-1}$ and the empirical or inferred ones in the $L = 4$ case (similar figures for $L = 10$ are in Appendix A 6). As expected, both methods perform well in the large γ limit with a correlation close to 1, and worse in the low γ limit. In this latter case, correlations due to phylogeny are too strong for our maximum likelihood method to pick up signal, and both methods perform equally poorly. However, there exists an intermediate regime where \mathbf{C}_{max} is much closer to the actual correlation than \mathbf{C}_{emp} . In Fig. 3, we plot the relative l_2 -error between either covariance matrices in the left panel or coupling matrices in the right panel. In both cases, our maximum-likelihood method results in a consistent improvement over the empirical estimator. However, the relative error still reaches high values in the low γ regime, which is likely due to \mathbf{C}_{max} and \mathbf{C}_{emp} being close to singular in this case.

An interesting way to illustrate the benefits of reconstructing the covariance matrix using knowledge of the tree is to evaluate the gain in *effective sample size*. Intuitively, the use of correlated samples reduces the information contained in the data, as compared to an equally large dataset of *i.i.d.* configurations. It is therefore interesting to compare the accuracy of our inferences with the accuracy obtained on smaller but *i.i.d.* samples. To do so, we report in Fig. S5 the l_2 -error between true and empirical covariances computed from a *i.i.d.* samples of variable sizes N . As expected, the error increases with decreasing values of N . We can use this in turn to express values of the l_2 -error in correlated samples in terms of effective *i.i.d.* sample sizes. For example, the error reached by \mathbf{C}_{emp} for $\gamma/\gamma_d \in [0.5, 1]$ and $L = 4$ corresponds to the one obtained for an *i.i.d.* sample of size ~ 16 , whereas it corresponds to a sample of size $\sim 32 - 64$ for \mathbf{C}_{max} . Thus, our correction is equivalent to increasing by a factor 2-4 the number of effective samples.

Finally, we assess the performance of our method in improving the prediction of the network of interactions between the Gaussian variables $\{x_a\}$. We consider that two variables

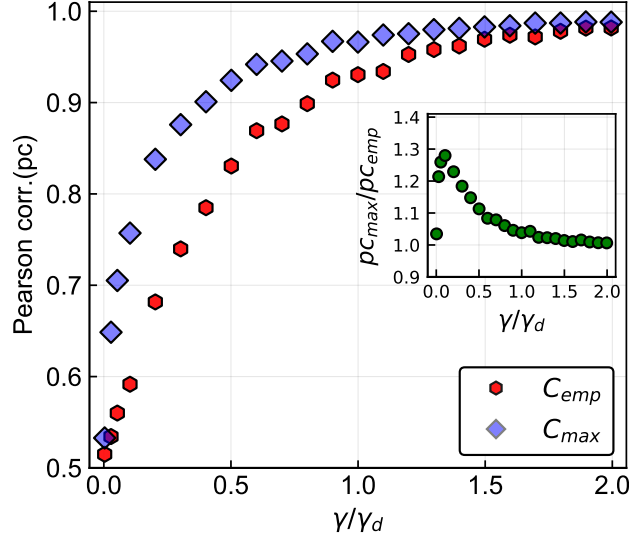


FIG. 2. Pearson correlation between empirical /maximum-likelihood covariance matrices and the true covariance matrix. The inset plot represents the ratio between the Pearson correlation for the maximum-likelihood covariance matrix and the one for the empirical covariance matrix. Simulations are performed for a tree of 512 leaves and system size $L = 4$.

x_a and x_b interact if the corresponding entry in the coupling matrix is non-zero, that is $J_{ab} \neq 0$. Using the data, we predict these interactions by taking the largest n elements (in absolute value) of the inferred coupling matrix, resulting in n predictions. The fraction TP/n of these n predictions that correspond to non-zero entries in the true matrix ($TP =$ true positives) defines the positive predictive value (PPV). This problem is equivalent to the one of predicting contacts in a protein structure

Fig. 4 shows the PPV as a function of the number of predictions for different values of γ and $L = 4$ (see Fig. S11 for the $L = 10$ case). In this case, the coupling matrix only has 6 independent non-diagonal elements, and only 6 predictions can be made. Our correction systematically outperforms the predictions from the empirical coupling matrix, with an always larger PPV. This gain is negligible in the extreme regimes of very high γ , where the prediction is close to identical to the one obtained with an *i.i.d.* sample, or very low γ , where it is essentially random. It is however much larger in the intermediate regime, with a significantly improved prediction in the region $\gamma/\gamma_d \in [0.5, 1]$.

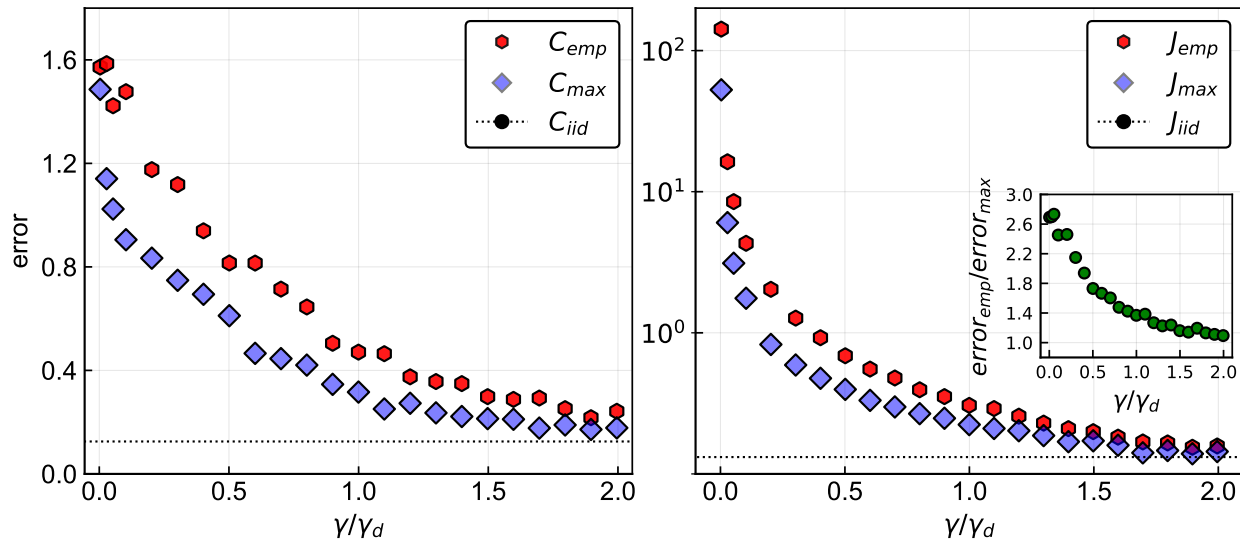


FIG. 3. **Left:**Relative l_2 -error between empirical or maximum-likelihood covariance matrices and the true covariance matrix. **Right:**Relative l_2 -error between empirical /maximum-likelihood coupling matrices and the true coupling matrix. Logarithmic scale is chosen for the y -axis because of large values of the error at low γ . The inset in both panels show the ratio between the two errors.

V. DISCUSSION

In this work, we proposed a method for inferring parameters of an Ornstein-Uhlenbeck process using data that is correlated through an evolutionary tree. We kept a very general setting in which data can in principle represent any set of continuous phenotypic traits or potentially discrete sequences if a continuous approximation is made. As such, our approach is purely methodological, and does not directly investigate any particular application.

We showed that due to the Gaussian and time reversible nature of the OU process, it is possible to write the joint covariance matrix of all data vectors in a simple way. The resulting matrix \mathbb{G} consists of block entries that represent covariances between pairs of leaves. The dependence of these blocks on the coupling matrix \mathbf{J} characterizing the OU process and on the tree structure can be written explicitly. Interestingly, \mathbb{G} only depends on the tree structure through the pairwise path length Δt_{ij} separating leaves along the tree.

We then proposed a way to compute the likelihood of the data given the tree and the parameters of the OU process, namely the coupling matrix \mathbf{J} and timescale γ . This method

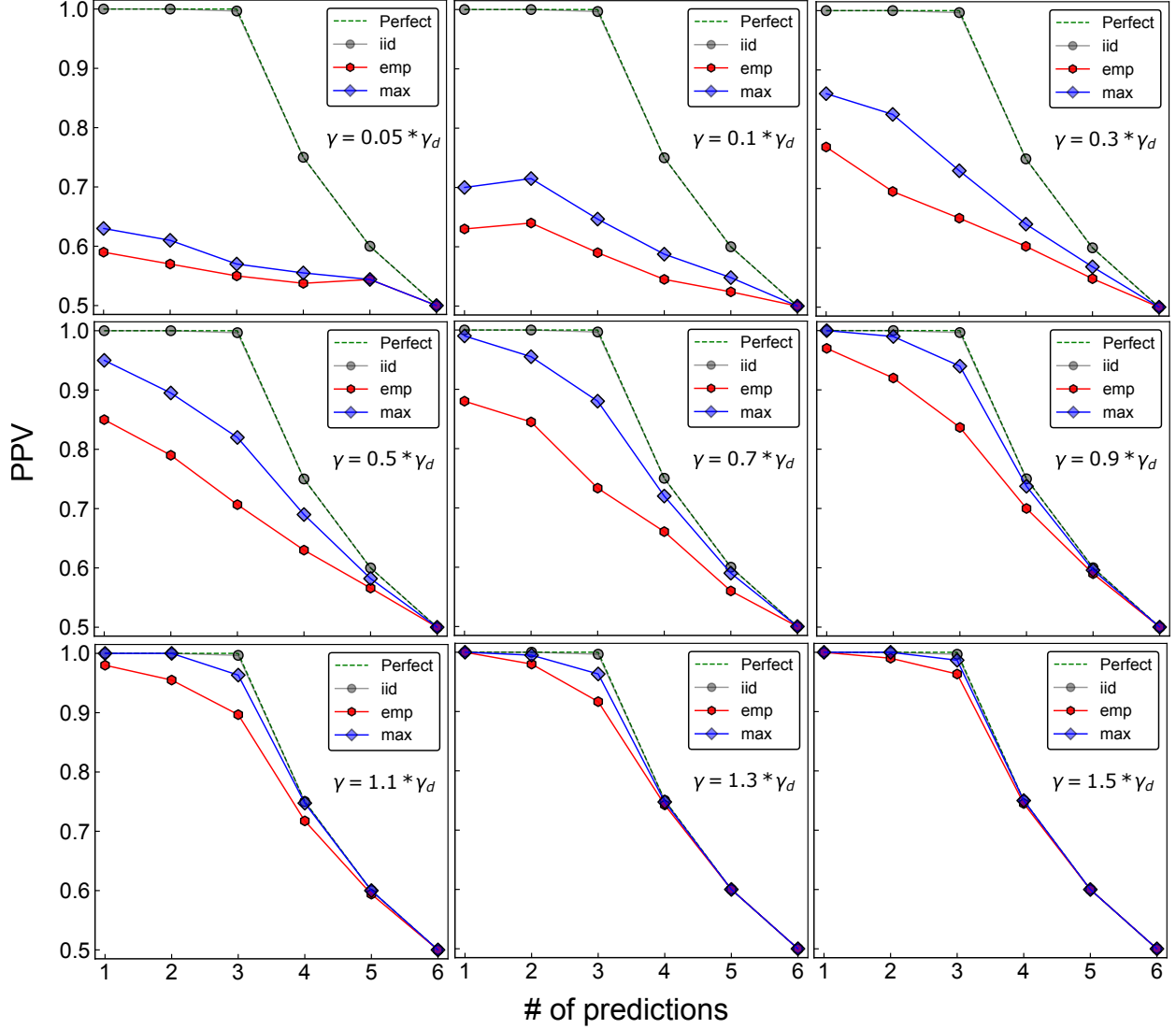


FIG. 4. Quality of prediction of interactions for different values of γ and system size $L = 4$. Interactions are defined as non-zero elements of the coupling matrix. In the $L = 4$ case, there are 6 possible interactions. Predictions are made by taking the largest elements (in absolute terms) of the inferred coupling matrix. The PPV is the fraction of correctly predicted contacts for a given number of predictions.

relies on computing the eigenvalues and vectors of the joint covariance matrix in an efficient manner. Indeed, it is possible to separate this calculation in two steps: the first in which we perform the eigen-decomposition of the matrix \mathbf{J} , and the second in which we compute eigenvalues and vectors of matrices \mathbf{G}^a that embed the tree structure. This reduces the computational complexity from $\mathcal{O}(L^3 N^3)$ for a naive inversion of \mathbb{G} to $\mathcal{O}(L^3) + \mathcal{O}(LN^3)$.

We also show that this method can be used to compute the gradient of the likelihood with respect to parameters with the same complexity. This makes the problem of inferring \mathbf{J} amenable to maximum likelihood methods using a gradient ascent approach.

Finally, we showed that this process gives encouraging results on simulated data, with a more accurate reconstruction of parameters than if empirical estimation was performed. These simulations highlight the fact that this method is only useful in the intermediate regime of phylogenetic correlations. If the timescale γ characterizing the branch lengths of the tree is too large, correlation of data points through the tree is weak and an empirical estimation performs well. On the other hand, a very low γ results in strong phylogenetic biases that make recovering \mathbf{J} impossible, basically due to a strong reduction of the information in a too redundant dataset. However, in an intermediate regime where intrinsic and historical correlations in the dataset coexist, our tree-aware re-construction of \mathbf{J} results in clear benefits over a tree-unaware empirical estimation.

A limitation of our approach remains the long computational time. Even with the efficient computation of the gradient, it was necessary to use small system sizes, $L = 10$ at most, to repeat our inference process many times with simulated data in a reasonable time. For this reason, the framework proposed here is limited to a small number of variables. In this respect, it is interesting to note that a different manner of computing the likelihood developed in [24] and based on Gaussian integrations on every branch of the tree results in an asymptotic complexity of $\mathcal{O}(NL^3)$.

Although our method can in principle be used for any set of traits, a major motivation in developing it is its potential application to model of proteins sequences. Several results in the last years have shown that selection forces shaping the evolution of protein sequences are well described by a pairwise potential. The estimation of this potential is performed using homologous sequences, and is therefore biased by the phylogenetic relations between these sequences. Results presented here are a first step in disentangling effects due to phylogeny from effects due to selection in a principled way.

However, there remain several challenges in using this framework for protein sequences. First, the computational power required to process actual sequences is much larger than what was needed for the small simulated systems presented here. As an example, a protein of length $L = 100$ will be represented by $q \times 100 = 2000$ Gaussian variables, where $q = 20$ is the number of amino acids. This is of course much larger than the $L = 10$ system used as an

example to test our approach.

A second question is the capacity of a continuous variable approximation, necessary when using Ornstein-Uhlenbeck dynamics, to represent dynamical properties of the landscape protein sequences evolve in. This type of approximation has been successfully used before, but in quite different contexts [25–27]. Its use in the context of modelling the evolutionary dynamics of protein sequences remains an open question.

Acknowledgments: We acknowledge interesting discussions with Roberto Mulet. PBC and MW acknowledge the hospitality of the Department of Theoretical Physics of University of Havana, where part of this work was done. Our work was partially funded by the EU H2020 Research and Innovation Programme MSCA-RISE-2016 under Grant Agreement No. 734439 InferNet.

-
- [1] Nguyen H. Chau, Zecchina R. N, and Berg J. Inverse statistical problems: from the inverse
ising problem to data science. *Adv. Phys*, 2017,66,197-261.
- [2] Ronald M Levy, Allan Haldane, and William F Flynn. Potts Hamiltonian models of protein
co-variation, free energy landscapes, and evolutionary fitness. *Current Opinion in Structural
Biology*, 43:55–62, April 2017.
- [3] Simona Cocco, Christoph Feinauer, Matteo Figliuzzi, Remi Monasson, and Martin Weigt.
Inverse Statistical Physics of Protein Sequences: A Key Issues Review. *Reports on Progress in
Physics*, 81(3):032601, March 2018. arXiv: 1703.01222.
- [4] Elad Schneidman, Michael J Berry, Ronen Segev, and William Bialek. Weak pairwise correla-
tions imply strongly correlated network states in a neural population. *Nature*, 440(7087):1007–
1012, 2006.
- [5] Yasser Roudi, Joanna Tyrcha, and John Hertz. Ising model for neural data: model quality and
approximate methods for extracting functional connectivity. *Physical Review E*, 79(5):051915,
2009.
- [6] William Bialek, Andrea Cavagna, Irene Giardina, Thierry Mora, Edmondo Silvestri, Massi-
miliano Viale, and Aleksandra M Walczak. Statistical mechanics for natural flocks of birds.
Proceedings of the National Academy of Sciences, 109(13):4786–4791, 2012.
- [7] Andrea Cavagna, Irene Giardina, and Tomás S Grigera. The physics of flocking: Correlation
as a compass from experiments to theory. *Physics Reports*, 728:1–62, 2018.
- [8] Thomas Bury. Market structure explained by pairwise interactions. *Physica A: Statistical
Mechanics and its Applications*, 392(6):1375–1385, 2013.
- [9] Stanislav S Borysov, Yasser Roudi, and Alexander V Balatsky. Us stock market interaction
network as learned by the boltzmann machine. *The European Physical Journal B*, 88(12):1–14,
2015.
- [10] Edwin T Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620,
1957.
- [11] Eric W Sayers, Mark Cavanaugh, Karen Clark, James Ostell, Kim D Pruitt, and Ilene
Karsch-Mizrachi. GenBank. *Nucleic Acids Research*, 47(D1):D94–D99, January 2019.
- [12] The UniProtConsortium. UniProt: the universal protein knowledgebase. *Nucleic Acids*

- Research*, 46(5):2699–2699, March 2018.
- [13] Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998.
- [14] Joseph Felsenstein and Joseph Felsenstein. *Inferring phylogenies*, volume 2. Sinauer associates Sunderland, MA, 2004.
- [15] Sara El-Gebali, Jaina Mistry, Alex Bateman, Sean R Eddy, Aurélien Luciani, Simon C Potter, Matloob Qureshi, Lorna J Richardson, Gustavo A Salazar, Alfredo Smart, et al. The pfam protein families database in 2019. *Nucleic acids research*, 47(D1):D427–D432, 2019.
- [16] Joseph Felsenstein. Phylogenies and quantitative characters. *Annual Review of Ecology and Systematics*, 19(1):445–471, November 1988. Publisher: Annual Reviews.
- [17] Benedikt Obermayer and Erel Levine. Inverse ising inference with correlated samples. *New Journal of Physics*, 16(12):123017, 2014.
- [18] Edwin Rodriguez Horta, Pierre Barrat-Charlaix, and Martin Weigt. Toward inferring potts models for phylogenetically correlated sequence data. *Entropy*, 21(11):1090, 2019.
- [19] Chongli Qin and Lucy J. Colwell. Power law tails in phylogenetic systems. *Proceedings of the National Academy of Sciences*, 115(4):690–695, January 2018.
- [20] Edwin Rodriguez Horta and Martin Weigt. Phylogenetic correlations have limited effect on coevolution-based contact prediction in proteins. *bioRxiv*, 2020.
- [21] G. E. Uhlenbeck and L. S. Ornstein. On the Theory of the Brownian Motion. *Physical Review*, 36(5):823–841, September 1930. Publisher: American Physical Society.
- [22] Thomas F. Hansen. Stabilizing Selection and the Comparative Analysis of Adaptation. *Evolution*, 51(5):1341–1351, 1997. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1558-5646.1997.tb01457.x>.
- [23] Krzysztof Bartoszek, Jason Pienaar, Petter Mostad, Staffan Andersson, and Thomas F. Hansen. A phylogenetic comparative method for studying multivariate adaptation. *Journal of Theoretical Biology*, 314:204–215, December 2012.
- [24] Venelin Mitov, Krzysztof Bartoszek, Georgios Asimomitis, and Tanja Stadler. Fast likelihood calculation for multivariate Gaussian phylogenetic models with shifts. *Theoretical Population Biology*, 131:66–78, February 2020.
- [25] David T. Jones, Daniel W. A. Buchan, Domenico Cozzetto, and Massimiliano Pontil. PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple

- sequence alignments. *Bioinformatics*, 28(2):184–190, January 2012.
- [26] J. P. Barton, S. Cocco, E. De Leonardis, and R. Monasson. Large pseudocounts and L₂-norm penalties are necessary for the mean-field inference of Ising and Potts models. *Physical Review E*, 90(1), July 2014.
- [27] Carlo Baldassi, Marco Zamparo, Christoph Feinauer, Andrea Procaccini, Riccardo Zecchina, Martin Weigt, and Andrea Pagnani. Fast and Accurate Multivariate Gaussian Modeling of Protein Families: Predicting Residue Contacts and Protein-Interaction Partners. *PLoS ONE*, 9(3), March 2014.
- [28] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, and M. Weigt. Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proceedings of the National Academy of Sciences*, 108(49):E1293–E1301, December 2011.
- [29] Matteo Figliuzzi, Herv Jacquier, Alexander Schug, Oliver Tenaillon, and Martin Weigt. Coevolutionary Landscape Inference and the Context-Dependence of Mutations in Beta-Lactamase TEM-1. *Molecular Biology and Evolution*, 33(1):268–280, January 2016.
- [30] William P. Russ, Matteo Figliuzzi, Christian Stocker, Pierre Barrat-Charlaix, Michael Socolich, Peter Kast, Donald Hilvert, Remi Monasson, Simona Cocco, Martin Weigt, and Rama Ranganathan. An evolution-based model for designing chorisate mutase enzymes. *Science*, 369(6502):440–445, July 2020. Publisher: American Association for the Advancement of Science Section: Report.
- [31] Rajesh Singh, Dipanjan Ghosh, and R. Adhikari. Fast bayesian inference of the multivariate ornstein-uhlenbeck process. *arxiv:1706.04961*, 2017.
- [32] Richard C. Raffenetti and Klaus. Ruedenberg. Parametrization of an orthogonal matrix in terms of generalized eulerian angles. *International Journal of Quantum Chemistry, Vol.III*, 625-634, 1970.
- [33] Ron Shepard, Scott R. Brozell, and Gergely Gidofalvi. The representation and parametrization of orthogonal matrices. *Journal of Physical Chemistry A*, 119, 7924-7939, 2015.
- [34] Michael Innes. Don’t unroll adjoint: Differentiating ssa-form programs. *CoRR*, abs/1810.07951, 2018.
- [35] Kaare Brandt Petersen and Michael Syskind Pedersen. *The Matrix Cookbook*. 2015.
- [36] Steven G. Johnson. The nlopt nonlinear-optimization package.

Appendix A: ~~Supplementary material~~ Description of technical details

1. Generating artificial data

We are interested in the case where the ~~described~~ dynamics of the L -dimensional Ornstein-Uhlenbeck process takes place on a tree. For example, if configurations \vec{x} represent quantitative traits of some organisms, the tree can represent the genealogy or phylogeny of these organisms. Therefore, to generate our datasets, we have to be able to simulate the OU process on a tree. In practice, given a rooted tree such as the one shown in ~~figure 1 of the main text~~ Fig. 1 above, we want to sample a configuration \vec{x} for every node in such a way that ~~equation Eq. (6)~~ holds for every pair of nodes, the time Δt ~~then being the branch length connecting them.~~ being the path length connecting the nodes along the tree.

We use a simple methodology to achieve this. First, note that given an arbitrary configuration \vec{x}_0 and a time Δt , we can generate a new configuration \vec{x} distributed according to the propagator ~~in equation Eq. (4)~~ by ~~the transformation exploiting the transformation~~

$$\vec{x} = \mathbf{\Lambda}^{\Delta t} \vec{x}_0 + \mathbf{\Sigma}^{1/2} \vec{\eta}, \quad (\text{A1})$$

where $\mathbf{\Lambda}$ and $\mathbf{\Sigma}$ are defined in ~~equation Eq. (5)~~, and $\vec{\eta}$ is a vector of uncorrelated variables ~~with distributions drawn individually from the normal distribution~~ $\mathcal{N}(0, 1)$. Moreover, if \vec{x}_0 is distributed according to ~~equation the equilibrium distribution Eq. (1)~~, then \vec{x} and \vec{x}_0 are distributed according to ~~equation the joint distribution Eq. (6) describing two equilibrium configurations at finite time difference.~~ Note that ~~equation Eq. (A1)~~ is quite different from the Langevin ~~equation Eq. (2)~~, ~~even though they have similar forms. While the Langevin equation describes the motion which describes the instantaneous dynamics~~ of \vec{x} in the potential ~~given by \mathbf{J} , eq. directly samples from the OU process. Given an and which could also be simulated in more complicated situation where no analytical expression for the propagator can be derived.~~

Given any already sampled internal node in the tree, ~~equation Eq. (A1)~~ allows to ~~sample~~ emit a configuration for each of its ~~children~~ child nodes. To sample the whole tree, we first ~~sample the root node draw the root configuration~~ \vec{x}_0 from the equilibrium distribution Eq. (1). By recursive applications of Eq. (A1), we then simply work our way down the tree until all leaves are sampled. Only the configurations in the leaves form the data set, the internal configuration remain hidden to our model-learning task.

2. Initializing parameters

a. Eigenvalues and eigenvectors of \mathbf{C}^{-1}

~~The initial value that we take for~~ We initialize the covariance matrix ~~is using~~ the empirical one:

$$\mathbf{C}^{emp} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i \cdot \vec{x}_i^T.$$

Its eigenmodes $\{\rho_a^0, \vec{s}_a^0\}$ determine the starting point of the optimization.

~~Insert here how we go from~~ A suitable parametrization of \vec{s}_a^0 to the corresponding set of
in terms of generalized Eulerian angles or ~~to the corresponding a~~ skew symmetric matrix ~~is~~
described below in Sec. A 3.

b. Time scale parameter γ

The optimization also requires that we initialize the time scale γ . For ~~this, we try~~
coherence with the last section, we need to find the optimal γ given the data \mathbf{X} , the tree,
and the OU process defined by the empirical covariance matrix.

The probability distribution P for the configurations of two leaves \vec{x}_i and \vec{x}_j separated by time Δt_{ij} is given by ~~equation~~ Eq. (6) of the main text. With this distribution we can ~~calculate analytically~~ analytically calculate the average of the scalar product ~~$\vec{x}_i^T \vec{x}_j$~~ :

$$\begin{aligned} \underline{\langle \vec{x}_i^T \vec{x}_j \rangle_P} &= \underline{\int d\vec{x}_i d\vec{x}_j P(\vec{x}_i, \vec{x}_j | \Delta t_{ij}) \sum_{a=1}^L x_i^a x_j^a} \\ &= \underline{\sum_{i=a}^L \langle x_i^a x_j^a \rangle_P}. \end{aligned}$$

$\vec{x}_i^T \cdot \vec{x}_j$ of two equilibrium configurations at given time separation:

$$\begin{aligned} \underline{\langle \vec{x}_i^T \cdot \vec{x}_j \rangle_P} &= \underline{\int d\vec{x}_i d\vec{x}_j P(\vec{x}_i, \vec{x}_j | \Delta t_{ij}) \sum_{a=1}^L x_i^a x_j^a} \\ &= \underline{\sum_{a=1}^L \langle x_i^a x_j^a \rangle_P}. \end{aligned} \tag{A2}$$

The covariance $\langle x_i^a x_j^a \rangle_P$ ~~between of two~~ observations separated by a time Δt_{ij} is given by ~~equation Eq. (7) of the main text~~. Using this, we ~~now have find~~

$$\begin{aligned}\langle \vec{x}_i^T \cdot \vec{x}_j \rangle_P &= \sum_{a=1}^L (\Lambda^{\Delta t_{ij}} \mathbf{C})_{aa} \\ &= \text{Tr}(\Lambda^{\Delta t_{ij}} \mathbf{C}) \\ &= \sum_{a=1}^L \rho_a^{-1} e^{-\gamma \rho_a \Delta t_{ij}}.\end{aligned}\tag{A3}$$

Having initialized the covariance matrix \mathbf{C} ~~at with~~ its empirical value, we know the values of all members of the r.h.s. of ~~equation Eq. (A3)~~ except ~~that of the one of~~ γ . ~~On the other hand, equivalent versions of equation can be written for all pairs of configurations i and j .~~ To find an initial value of γ which is consistent with the data and the empirical covariance matrix ~~for all pairs of data configurations $i < j$~~ , we search for one that best explains the observed scalar products between configurations. We thus define γ^0 to be the argument minimizing the functional $F(\gamma)$:

$$F(\gamma) = \sum_{1 \leq i < j \leq N} \left[\vec{x}_i^T \cdot \vec{x}_j - \sum_{a=1}^L \rho_a^{-1} e^{-\gamma \rho_a \Delta t_{ij}} \right].\tag{A4}$$

~~As Since~~ F depends on ~~one a single~~ scalar parameter, it is straightforward to minimize it ~~;~~ ~~allowing us and thereby~~ to initialize γ to ~~a an empirically~~ reasonable value.

3. Parametrizations of eigenvectors

a. Parametrization ~~in term of using~~ generalized Eulerian angles

The ~~aim is to parameterize each base vector~~ ~~idea is to write the base vectors~~ \vec{s}_a ~~according to the a -th column of an arbitrary as columns of an~~ orthogonal matrix ~~S parameterized, and to parameterize this matrix~~ in terms of $L(L-1)/2$ independent variables $\theta_{pq}, p=1, 2, \dots, L; q=1, 2, \dots, L; p < q$, ~~named Generalized Eulerian angles (Eulerian angles is for $L=3$)~~ θ_{pq} with $1 \leq p < q \leq L$. These parameters are called generalized Eulerian angles, since they generalize Eulerian angles to $L > 3$.

To construct this matrix we start from a ~~rotational~~ transformation in a two-dimensional subspace of an ~~L -dimensional space which~~ L -dimensional space. It is given by an ~~L -dimensional~~ L -dimensional matrix of the form \div

$$T_{pq} = \begin{pmatrix} 1 & & & & \\ & 1 & & & \\ & & \cos \theta_{pq} & \sin \theta_{pq} & \\ & & & 1 & \\ & & -\sin \theta_{pq} & \cos \theta_{pq} & \\ & & & & 1 \end{pmatrix} \quad (A5)$$

where all diagonal elements are unity except for the diagonal elements in the p th column and the q th row, which are equal $\cos \theta_{pq}$ and all off-diagonal elements are zero except for the one corresponding to the intersection of the p th row and the q th column, which is $\sin \theta_{pq}$, and that on the intersection of the q th row and the p th column, which is $-\sin \theta_{pq}$. There are $L(L-1)/2$ matrices of the form indicated in Equation (1), corresponding to all choices of p and q with $1 \leq p < q \leq L$.

An arbitrary L -dimensional orthogonal matrix \mathbf{S} can be represented as a product of these $L(L-1)/2$ orthogonal matrices with appropriate values of the $L(L-1)/2$ independent parameters θ_{pq} .

$$\mathbf{S}(\boldsymbol{\theta}) = \prod_{\{pq\}} T_{pq}(\theta_{pq}) \quad (A6)$$

This matrix is in charge of rotations of vectors in L -dimensional space. The basic idea is to express a general rotation in L dimensions as the result of $L(L-1)/2$ successive rotations by angles θ_{pq} in the two-dimensional subspaces given by the components p, q . For $L = 3$ we get the classical Eulerian matrix.

In reference [7] is explained an algorithm to efficiently perform this multiplication, as well as the construction of the matrix derivatives respect to parameters θ_{pq} .

For the first, equation (2) is transformed in the recurrence equations. To start, Eq. (A6) is transformed into a set of recursive equations, with \mathbf{S} being equal to the final step

$$\mathbf{S}(\boldsymbol{\theta}) = \mathbf{T}^{(L)} \quad (A7)$$

For all $1 \leq k < l \leq L$ we iterate over n :

$$\mathbf{S}(\vec{\theta}) = \mathbf{T}^{(L)}$$

~~with $\theta_{pq}, p = 1, 2, \dots, L; q = 1, 2, \dots, L; p < q$ and~~

[7] What is the initialization? Please check carefully, also the correct ranges of the kl which seem to go to n only. It would also be better to avoid s , which we use for the eigenvectors.

$$T_{kl}^{(n)} = \cos \theta_{kn} * t_{kl}^{(n)} - \sin \theta_{kn} * s_{kl}^{(n)} \quad (\text{A8})$$

with

$$\begin{aligned} s_{k+1,l}^{(n)} &= \sin \theta_{kn} * t_{kl}^{(n)} + \cos \theta_{kn} * s_{kl}^{(n)} \\ s_{k+1,n}^{(n)} &= \cos \theta_{kn} * s_{kn}^{(n)} \\ s_{1,l}^{(n)} &= -\delta_{ln} \\ \theta_{nn} &= \pi/2 \end{aligned} \quad (\text{A9})$$

and

$$\mathbf{t}^{(n)} = \begin{pmatrix} \mathbf{T}^{(n-1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \quad (\text{A10})$$

~~This mean that~~ As a consequence, the eigenvectors \vec{s}_a of the matrix \mathbf{C} with dimensions $L \times L$ can be written as:

$$s_a^k = T_{.,a}^{(L)} = \cos \theta_{kL} * t_{ka}^{(L)} - \sin \theta_{kL} * s_{ka}^{(L)} \quad \text{for } k = 1, \dots, L \quad (\text{A11})$$

~~Then there is another important problem: how~~ To use this expression, we still need to determine parameters $\boldsymbol{\theta}$ given an orthogonal matrix ~~\mathbf{Q}~~ \mathbf{S} , such that all equations

$$\underline{S_{ij}(\boldsymbol{\theta}) = Q_{ij}}$$

$$\underline{S_{ij}(\boldsymbol{\theta}) = S_{ij}} \quad (\text{A12})$$

are satisfied. This system of nonlinear transcendental equations ~~can't be algebraically solved~~, ~~however it's possible~~ cannot be solved algebraically. However, it is possible to overcome this issue finding the set of $\boldsymbol{\theta}$ which minimize the ~~function:~~ square distance between the target and the parametrized matrices:

$$\underline{f(\boldsymbol{\theta})} \hat{\boldsymbol{\theta}} = \underline{\text{argmin}_{\boldsymbol{\theta}}} \sum_{\underline{i,j} \leq j} \left[\underline{Q} \underline{S_{ij}} - S_{ij}(\boldsymbol{\theta}) \right]^2 \quad (\text{A13})$$

This is useful when we initialize parameters $\boldsymbol{\theta}$ ~~from~~ for the matrix formed by the eigenvectors of the empirical covariance matrix.

b. Parametrization in term of the exponential of a skew-symmetric matrix

~~The~~

Any orthogonal matrix can be written as the exponential of a skew-symmetric matrix

$$\mathbf{X} = -\mathbf{X}^T;$$

$$\mathbf{S} = \exp(\mathbf{X}), \quad (\text{A14})$$

~~is an orthogonal matrix: $\exp(\mathbf{X})^T = \exp(\mathbf{X}^T) = \exp(-\mathbf{X}) = \exp(\mathbf{X})^{-1}$~~ since $\exp(\mathbf{X})^T = \exp(\mathbf{X}^T) = \exp(-\mathbf{X})$

The derivatives ~~on~~ with respect to the $L(L-1)/2$ independent variables of \mathbf{X} is formally defined by

$$\frac{\partial \mathbf{S}}{\partial X_{jk}} = \lim_{h \rightarrow 0} \frac{1}{h} (\exp(\mathbf{X} + h \mathbf{E}^{jk}) - \exp(\mathbf{X})) \quad (\text{A15})$$

where \mathbf{E}^{jk} for $j > k$ is defined as a skew-symmetric matrix that has only two nonzero elements entries in positions (j, k) and (k, j) :

$$E_{jk}^{pq} = \delta_{pj} \delta_{qk} - \delta_{pk} \delta_{qj} \quad (\text{A16})$$

To obtain a skew-symmetric matrix \mathbf{X} from an orthogonal matrix \mathbf{Q} ~~is enough~~, it is sufficient to invert the exponential relation:-

$$\underline{\mathbf{X} = \log \mathbf{Q}}$$

~~$\mathbf{X} = \log \mathbf{Q}$~~ . If we write eigenvectors of the ~~likelihood function covariance matrix~~ as the columns of the exponential of the skew-symmetric matrix \mathbf{X} , ~~then~~ we are able to perform the optimization over ~~it's $L(L-1)$ independent variables of \mathbf{X}~~ its $L(L-1)/2$ independent entries.

Furthermore, as the eigenvectors ~~were parameterized~~ are parametrized by an algebraic function, the complete likelihood is described in terms of ~~arithmetics operation~~ arithmetic operations and elementary functions. This allow to compute the gradient via automatic differentiation (AD) [CITE](#), making computations in our optimization process faster. The idea is to represent a function as a computational graph, also called ~~a Wengert list~~ Wengert list [CITE](#), where each node in ~~this list will represent the list represents~~ an intermediate result of the computation. The intermediate results can ~~then~~ be assembled using the chain rule to get the final derivative we ~~re~~ are looking for. There ~~is~~ are two main algorithms to ~~traverses~~ traverse the chain rule in AD: a forward mode and ~~reverse mode, here we used reverse mode~~

algorithm which is a reverse mode. Here, we used the reverse-mode algorithm, which is also the one of choice for back-propagation in deep learning. In particular we implemented it using the Julia package Zygote.jl.

4. Homogeneous and fully balanced tree

Let's assume that the tree is binary, symmetric and completely homogeneous with all branches having the same length Δt . As an example, the covariance matrix for such a tree with $K = 2$ branching events and four leaves is

$$\mathbb{G} = \begin{pmatrix} \mathbf{C} & \mathbf{C}\Lambda^{2\Delta t} & \mathbf{C}\Lambda^{4\Delta t} & \mathbf{C}\Lambda^{4\Delta t} \\ \mathbf{C}\Lambda^{2\Delta t} & \mathbf{C} & \mathbf{C}\Lambda^{4\Delta t} & \mathbf{C}\Lambda^{4\Delta t} \\ \mathbf{C}\Lambda^{4\Delta t} & \mathbf{C}\Lambda^{4\Delta t} & \mathbf{C} & \mathbf{C}\Lambda^{2\Delta t} \\ \mathbf{C}\Lambda^{4\Delta t} & \mathbf{C}\Lambda^{4\Delta t} & \mathbf{C}\Lambda^{2\Delta t} & \mathbf{C} \end{pmatrix}. \quad (\text{A17})$$

The associated matrix $\mathbf{G}^a = z(\rho_a, \gamma, \Delta t)$ defined in Eq. (A17) becomes

$$\mathbf{G}^a = \begin{pmatrix} \rho_a^{-1} & \rho_a^{-1}e^{-2\gamma\rho_a\Delta t} & \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} & e^{-4\gamma\rho_a\Delta t} \\ \rho_a^{-1}e^{-2\gamma\rho_a\Delta t} & \rho_a^{-1} & \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} & \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} \\ \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} & \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} & \rho_a^{-1} & \rho_a^{-1}e^{-2\gamma\rho_a\Delta t} \\ \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} & \rho_a^{-1}e^{-4\gamma\rho_a\Delta t} & \rho_a^{-1}e^{-2\gamma\rho_a\Delta t} & \rho_a^{-1} \end{pmatrix}. \quad (\text{A18})$$

For hyper-geometric matrices as (A18) of dimension 2^K , there are $K+1$ different eigenvalues given by:

$$\lambda_k(\rho_a, \gamma) = \rho_a^{-1} * \begin{cases} (1 + \sum_{l=1}^{k-1} 2^{l-1} e^{-2l\gamma\rho_a\Delta t} - 2^{k-1} e^{-2k\gamma\rho_a\Delta t}) & k \in [1, K] \\ (1 + \sum_{l=1}^K 2^{l-1} e^{-2l\gamma\rho_a\Delta t}) & k = K + 1 \end{cases} \quad (\text{A19})$$

where $\lambda_{K+1} \geq \lambda_K \cdots \geq \lambda_1$. For $k < K + 1$, the degeneracy of eigenvalue λ_k is $d_k = 2^{K-k}$. The associated eigenvectors are independent of the parameter ρ_a and reflect the events in the phylogenetic tree. Each eigenvector \vec{u}_k captures the duplication events in the $(K + 1 - k)st$ generation:

$$\vec{u}_k = \begin{cases} \underbrace{(1, \dots, 1)}_{2^{k-1}}, \underbrace{(-1, \dots, -1)}_{2^{k-1}}, 0, \dots, 0) \cup \Gamma(u_k) & k \in [1, K] \\ (1, 1, 1, \dots, 1, 1, 1) & k = K + 1 \end{cases}$$

where $\Gamma(\vec{u}_k)$ represents the d_k combinations obtained by shifting the block of length Q , generating all eigenvectors corresponding to the eigenvalue λ_k . The eigenvectors are orthogonal to each other, and can be normalized and arranged horizontally into a matrix U .

To compute the gradient of the likelihood, we need derivatives of $\lambda_k(\rho_a, \gamma)$ with respect to ρ_a and γ , which can be directly obtained from expression (A19).

5. Optimization scheme

The proposed inference scheme was transformed into a multidimensional nonlinear optimization problem for which a variant of quasi-Newton Methods (QNMs) was used. The main feature in QNMs is that Hessian matrix \hat{H} does not need to be computed instead its approximated. The Hessian approximation \hat{H} is chosen to satisfy the secant equation:

$$\nabla L(\vec{x}_k + \Delta \vec{x}) = \nabla L(\vec{x}_k) + \hat{H} \Delta \vec{x}$$

for $n - dimensions$ \hat{H} is undetermined, the various QNMs differ in their choice of the solution to the secant equation. We used LBFGS variant (limited memory version of BFGS), this particular method is based in choosing \hat{H} as a positive definite matrix where

$$\hat{H}_{k+1} = \hat{H}_k + \frac{y_k y_k^T}{y_k^T \Delta x_k} - \frac{\hat{H}_k \Delta x_k (\hat{H}_k \Delta x_k)^T}{\Delta x_k^T \hat{H}_k \Delta x_k} \quad \text{and} \quad y_k = \nabla L(\vec{x}_{k+1}) - \nabla L(\vec{x}_k)$$

The method implementation was done using NLOpt package [36].

6. ~~Supplementary~~ Supporting figures

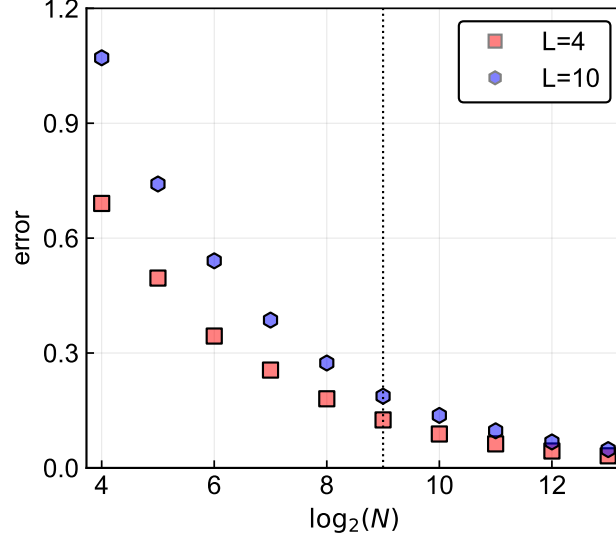


FIG. 5. Relative l_2 -error between the empirical covariance matrix calculated from an *i.i.d.* sample and the true covariance matrix, for system sizes $L = 4$ and $L = 10$. The dashed vertical line corresponds to the number of leaves of the tree used in the simulations.

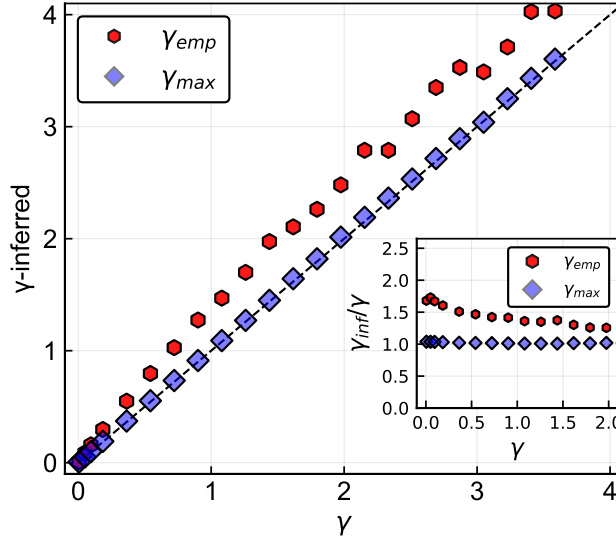


FIG. 6. Inferred γ values as a function of real γ , for system size $L = 4$. γ_{emp} is the value obtained by the process described in section A 2 of the SM. γ_{max} is the value inferred by the maximum-likelihood calculation. The inset represents the ratio of both inferred parameters γ_{emp} or γ_{max} to the real γ .

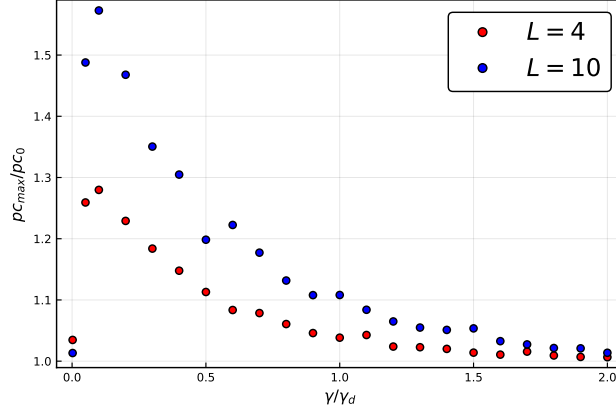


FIG. 7. Ratio between Pearson correlation of the maximum-likelihood and true covariance matrix to the Pearson correlation of the empirical and true covariances matrices for the two system sizes $L = 4$ and $L = 10$.

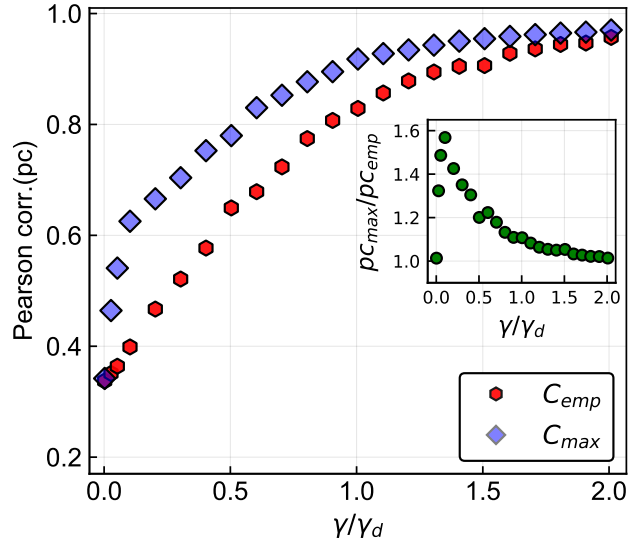


FIG. 8. Pearson correlation between empirical / maximum-likelihood covariance matrices and the true covariance matrix, the inset plot represent the ratio between the person correlation for the maximum-likelihood covariance matrix and the one for the empirical covariance matrix.

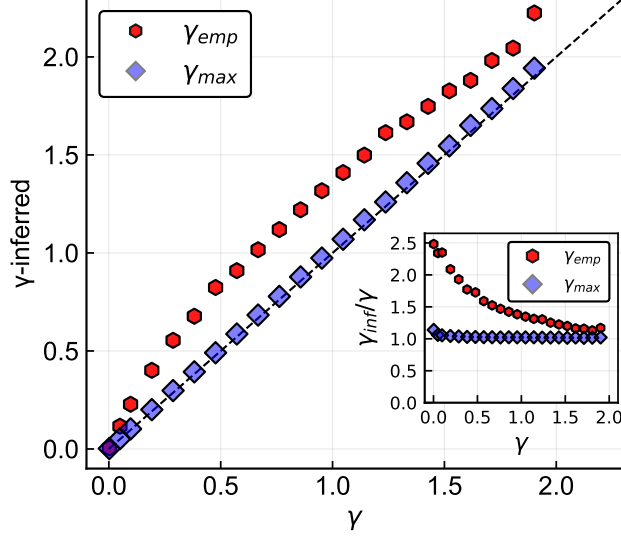


FIG. 9. Inferred γ values as a function of real γ , for system size $L = 10$. γ_{emp} is the value obtained by the process described in section A 2 of the SM. γ_{max} is the value inferred by the maximum-likelihood calculation. The inset represents the ratio of both inferred parameters γ_{emp} or γ_{max} to the real γ .

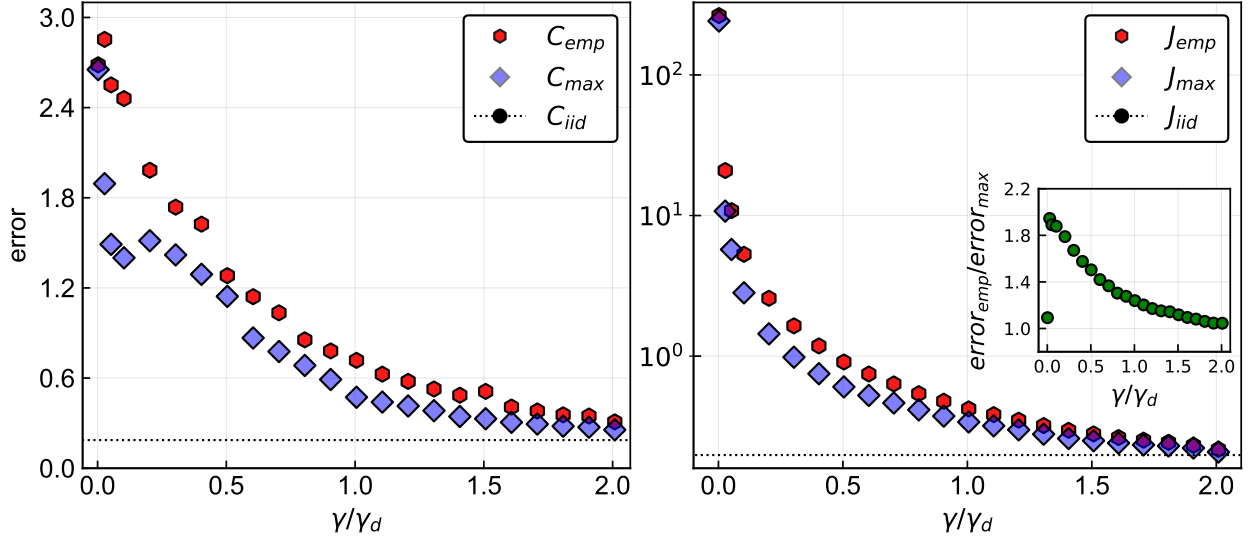


FIG. 10. **Left:** Relative l_2 -error between empirical or maximum-likelihood covariance matrices and the true covariance matrix. **Right:** Relative l_2 -error between empirical /maximum-likelihood coupling matrices and the true coupling matrix. Logarithmic scale is chosen for the y-axis because of large values of the error at low γ . The inset in both panels show the ratio between the two errors. For system size $L = 10$.

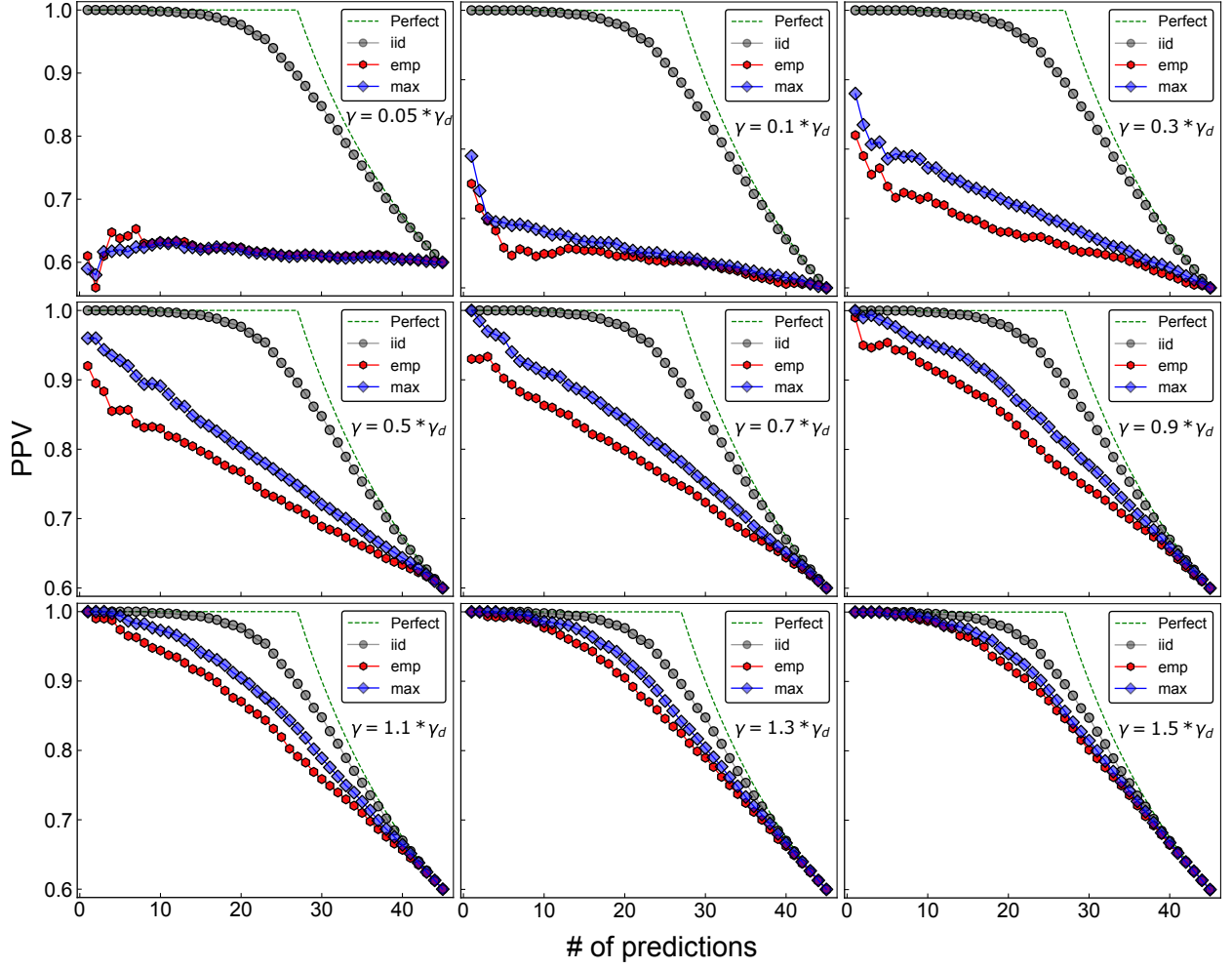


FIG. 11. Quality of prediction of interactions for different values of γ and system size $L = 4$. Interactions are defined as non-zero elements of the coupling matrix. In the $L = 4$ case, there are 6 possible interactions. Predictions are made by taking the largest elements (in absolute terms) of the inferred coupling matrix. The PPV is the fraction of correctly predicted contacts for a given number of predictions.