

PAPER • OPEN ACCESS

## Fluctuations and the limit of predictability in protein evolution

To cite this article: Saverio Rossi *et al* 2025 *Rep. Prog. Phys.* **88** 078102

View the [article online](#) for updates and enhancements.

You may also like

- [Quantum Fisher information and its dynamical nature](#)  
Matteo Scandi, Paolo Abiuso, Jacopo Surace et al.
- [Twistronics and moiré superlattice physics in 2D transition metal dichalcogenides](#)  
Dawei Zhai, Hongyi Yu and Wang Yao
- [Kicking Time Back in Black Hole Mergers: Ancestral Masses, Spins, Birth Recoils, and Hierarchical-formation Viability of GW190521](#)  
Carlos Araújo-Álvarez, Henry W. Y. Wong, Anna Liu et al.

# Fluctuations and the limit of predictability in protein evolution

Saverio Rossi<sup>1,2</sup> , Leonardo Di Bari<sup>3,4</sup> , Martin Weigt<sup>4</sup>  and Francesco Zamponi<sup>1,2,\*</sup> 

<sup>1</sup> Dipartimento di Fisica, Sapienza Università di Roma, Piazzale Aldo Moro 5, 00185 Rome, Italy

<sup>2</sup> Laboratoire de Physique de l'Ecole Normale Supérieure, ENS, Université PSL, CNRS, Sorbonne Université, Université Paris Cité, 75005 Paris, France

<sup>3</sup> DISAT, Politecnico di Torino, Corso Duca degli Abruzzi, 24, I-10129 Torino, Italy

<sup>4</sup> Sorbonne Université, CNRS, Laboratory of Computational and Quantitative Biology, 75005 Paris, France

E-mail: [francesco.zamponi@uniroma1.it](mailto:francesco.zamponi@uniroma1.it)

Received 17 December 2024, revised 17 June 2025

Accepted for publication 1 July 2025

Published 9 July 2025

Corresponding editor: Dr Paul Mabey



## Abstract

Protein evolution involves mutations occurring across a wide range of time scales. In analogy with disordered systems in statistical physics, this dynamical heterogeneity suggests strong correlations between mutations happening at distinct sites and times. To quantify these correlations, we examine the role of various fluctuation sources in protein evolution, simulated using a data-driven energy landscape as a proxy for protein fitness. By applying spatio-temporal correlation functions developed in the context of disordered physical systems, we disentangle fluctuations originating from the initial condition, i.e. the ancestral sequence from which the evolutionary process originated, from those driven by stochastic mutations along independent evolutionary paths. Our analysis shows that, in diverse protein families, fluctuations from the ancestral sequence predominate at shorter time scales. This allows us to identify a time scale over which ancestral sequence information persists, enabling its reconstruction. We link this persistence to the strength of epistatic interactions: ancestral sequences with stronger epistatic signatures impact evolutionary trajectories over extended periods. At longer time scales, however, ancestral influence fades as epistatically constrained sites evolve collectively. To confirm this idea, we apply a standard ancestral sequence reconstruction (ASR) algorithm and verify that the time-dependent recovery error is influenced by the properties of the ancestor itself. Overall, our results reveal that the properties of ancestral sequences—particularly their epistatic constraints—influence the initial evolutionary dynamics and the performance of standard ASR algorithms.

\* Author to whom any correspondence should be addressed.



Original Content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](https://creativecommons.org/licenses/by/4.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

Supplementary material for this article is available [online](#)

Keywords: protein evolution, sequence landscapes, epistasis, disordered systems

## 1. Introduction

Proteins play a key role in many biological processes that are essential for life. At the same time, they display a huge evolutionary flexibility, in that a large diversity of protein sequences can fold into the same structure and perform the same biological function. Such proteins are called ‘homologous’ and grouped into the same ‘protein family’ [1–4]. During evolution, a single ancestral protein belonging to a given family can diversify its amino acid sequence and, given enough time, explore a portion of the ‘neutral space’ of equivalently fit protein sequences [5]. Nevertheless, a random mutation in the amino acid sequence has a high probability of negatively impacting the structure and functionality of the protein [6–8]. Such kind of deleterious mutations are eliminated by natural selection, by which certain protein variants are favored over others due to their functional advantages: neutral (or even beneficial) mutations that maintain (or enhance) protein structure, stability, or function are more likely to persist in a population [9].

The picture gets more complex as one takes into account the concept of epistasis [10–18]: the effect of a mutation on the sequence fitness changes with the ‘background’ in which the mutation takes place, i.e. the amino acids that are present in the other sites of the protein. As a consequence, a mutation that would be deleterious in a particular background sequence can be instead beneficial in another (also called sign epistasis), allowing evolution to explore different pathways.

Understanding and characterizing the impact of epistasis in evolution requires careful experiments and modeling. In particular, recent developments have substantially increased the sequence divergence that can be reached by laboratory evolution experiments [19–23]. A large amount of data is therefore now becoming available with more expected to come soon. Yet, the sequence diversity of natural evolution still remains out of reach of such experiments, thus leaving an unexplored gap in evolutionary time scales. In order to fill this gap, one can simulate the evolution of protein sequences *in silico* [24–29], relying on the data-driven approach that goes under the name of direct coupling analysis (DCA) [15, 30]. DCA infers a fitness landscape (analogous to an energy function in the statistical physics vocabulary) starting from a multiple sequence alignment (MSA) of natural homologs constituting a given protein family [31–35]. The energy function that results from this inference procedure can then be used to assign a probability to each sequence. The resulting landscape is explored by means of a biologically motivated Monte Carlo Markov-Chain (MCMC) algorithm, providing *in silico* evolutionary trajectories that quantitatively mimic experimental results [24–29, 36].

Within this framework, it has been recently shown [29] that different sites evolve with widely different time scales, which also depend on the background sequence, due to epistatic interactions. More precisely, sites that are more epistatically constrained need much longer times (or number of generations) to evolve (i.e. accumulate mutations) compared to less constrained sites [29]. Furthermore, whether a site is epistatically constrained or not depends on the rest of the sequence (the background) [37–39], which adds sequence-to-sequence heterogeneity on top of the site-to-site heterogeneity. This ‘heterogeneity of time scales’ is strongly reminiscent of similar phenomena observed in disordered physical systems [40–50]. There, it has been shown that a large heterogeneity in time scales implies the presence of strong and highly collective space-time correlations between sites that slow down the dynamical evolution. Furthermore, it has been shown that the structural disorder in the initial configuration partially encodes future correlations, and can be used (possibly by machine learning tools) to predict the future dynamics [51–54], see [55] for a recent review. These are the main observations that motivated this work.

In this paper, we measure space-time correlations to explore how the heterogeneity of evolutionary time scales relates to the epistatic interactions within the ancestral sequence. We follow the dynamics of the Hamming distance, i.e. the number of accepted mutations with respect to the ancestral sequence. We characterize how this quantity fluctuates (i) between different evolutionary trajectories originating from the same ancestor, and (ii) between different ancestors. We observe a strong dependence of the dynamics on the ancestral sequence for short enough times, over which the first source of fluctuations is found to be subdominant with respect to the second. On the other hand, at long times the stochasticity of evolutionary trajectories takes over, and any memory of the ancestor is lost. This has important consequences, as it allows us to properly quantify the time scale over which it should be possible to reconstruct the ancestral sequence of a certain set of evolutionary trajectories, and how this time scale depends on the epistatic interactions in the ancestor itself. In fact, we find that the amount of epistatically constrained sites in the ancestral sequence determines this time scale: more epistatic sequences leave their trace on evolutionary dynamics for longer times. We then measure the correlations between the evolution of all pairs of sites at the time scale at which such correlations are stronger, finding a different pattern for each ancestral sequence, which is then reflected in the evolutionary dynamics. Finally, we show that the magnitude of epistatic interactions, which controls the stochasticity of the evolutionary trajectories, is also related to the linear response

of the evolutionary dynamics to a change in selective pressure (controlled by the temperature in our Monte Carlo simulations). Hence, we show and quantify how more epistatically constrained ancestors lead to a more complex evolutionary dynamics over longer time scales, which is also more sensitive to perturbations of the environment.

## 2. Epistasis in biological data

Epistasis plays a central role in shaping protein evolution, influencing both the structure of fitness landscapes and the fate of evolutionary trajectories, and as such it has been the subject of theoretical, computational and experimental studies, see [10–18] for a few recent reviews. It has been shown that some epistatic effects can be accounted for by a global non-linearity of the phenotype-fitness relation [56–59], but epistasis is also due to a network of more specific pairwise and higher-order interactions. While some of these interactions are sparse [18, 60–66], multiple studies suggest that the most relevant effects tend to emerge collectively, from the accumulation of numerous weak interactions between a single site and multiple other residues across the protein sequence [14, 15, 32, 38, 39, 67–70]. This form of distributed epistasis suggests that the evolutionary constraints on a given mutation are shaped by the broader sequence context, rather than by a few dominant interactions, highlighting the complexity of protein fitness landscapes.

Experiments have provided valuable insights into how mutations interact, but each methodology comes with inherent limitations. While some studies focus on measuring epistasis in a restricted mutational space, others provide broader datasets but lack information on evolutionary dynamics. Here, we summarize the main categories of experimental data available and discuss their relevance to our theoretical framework.

- Combinatorial mutagenesis [18, 63, 65, 71–74] experiments involving a small number of residues have demonstrated epistatic effects. However, these studies typically do not allow the accumulation of sufficient mutational effects to observe the large-scale evolutionary patterns we investigate in this work.
- Deep mutational scanning experiments across homologous wild-type proteins [22, 39, 75, 76], reveal epistatic interactions and can be used to explore the concept of site variability within a small region of the landscape. However, these experiments lack direct observations of evolutionary dynamics over multiple generations.
- Experimental studies tracking evolutionary dynamics *in vitro* are available [19–21, 23], but only a few cases, such as TEM-1 [19] and PSE-1 [20]  $\beta$ -lactamases, involve multiple wild-type homologs belonging to the same protein family, whose evolutionary dynamics can be directly compared. Moreover, these experiments do not extend far enough in sequence space to capture the long-term epistatic effects central to our study.

Given these limitations, we anticipate that future advancements in experimental techniques will provide richer datasets capable of testing the predictions outlined in our study. In the meantime, our theoretical framework serves as a guide for interpreting existing data and shaping expectations for future experimental work.

## 3. Methods

### 3.1. Modeling evolution *in silico*

In order to mimic the evolution of protein sequences, we use the DCA model energy as a proxy for the fitness landscape. We start by building an MSA of naturally occurring sequences for each protein family we want to study. From the MSA we infer the parameters (fields and couplings) of the DCA model via Boltzmann machine learning (bmDCA) [77, 78]. The resulting model assigns a probability  $P(A) = \exp[-\mathcal{H}(A)]/Z$  to each sequence  $A = (a_1, \dots, a_L)$ , with  $L$  the common length of aligned sequences in the MSA and  $a_i$  being a symbol that takes 21 possible values corresponding to the 20 natural amino acids plus the gap symbol needed for alignment. According to the statistical physics language, low energy

$$\mathcal{H}(A) = - \sum_{i < j} J_{ij}(a_i, a_j) - \sum_i h_i(a_i) \quad (1)$$

corresponds to high probability, hence high fitness. In this expression, epistatic interactions between different amino acids are represented via the pairwise couplings  $J_{ij}(a_i, a_j)$  that have been found to be crucial in data-driven statistical models of biological sequences, see [15, 79]. Because training of these models has become a standard and well-documented procedure [77, 78], we do not give further details here.

Following [24, 25, 29], we consider a fixed initial sequence  $A_0 = (a_1^0, \dots, a_L^0)$  as the ancestor, and we let many trajectories evolve from it in parallel, hence restricting ourselves to a star phylogeny describing an ensemble of independent evolutionary trajectories of common initialization. Our goal is to characterize the statistical properties of the specific ancestral sequence and the sequence space accessible from it in a given evolutionary time. The sequence evolution is simulated by Monte Carlo dynamics. In this work, we use the Metropolis algorithm acting on amino acids for simplicity, and we verified that more refined sampling strategies taking into account amino-acid accessibility via the genetic code, insertions and deletions [24, 25, 29] produce qualitatively the same results, see the supplemental material (SM). As it was previously shown [79], at large times the generated sequences accurately reproduces many statistical features of the natural sequences used for the training, and have the same probability of being biologically functional in a given experimental platform; hence, the model is generative. Here, we are concerned with what happens at short and intermediate times, where the influence of the ancestral sequence is still important.

The results we report mainly concern the DNA-binding domain (DBD) protein family already studied in a similar setting in [29] and experimentally in [22], but for some results we also generalize to other protein families, in particular the WW domain (WW), chorismate mutase (CM), aminoglycoside 6-N-acetyltransferase (AAC6), dihydrofolate reductase (DHFR), beta lactamase (BL), and serine protease (SP) families. These families have been chosen because they span several chain lengths, and experimental data obtained either from Deep Mutational Scanning or by *in vitro* evolution are available. The procedures to construct the natural MSAs for these families are given in the SM.

### 3.2. Measures of epistasis

By using the MSA of a given protein family as input data, and inferring the fitness landscape parameters via the DCA model, the authors of [29, 38] have been able to roughly classify the sites of any specific protein sequence belonging to the family in three categories: conserved, mutable, and epistatically constrained. In order to do so, they defined two site-mutability metrics.

- The context-independent entropy (CIE) is obtained by computing the empirical frequency  $f_i(a)$  of occurrence of amino acid  $a$  on site  $i$ , for every  $a$  and  $i$ , in the MSA obtained from the natural sequences. This is then used to compute a standard Shannon entropy as

$$\text{CIE}_i = - \sum_{a=1}^{21} f_i(a) \log_2 f_i(a) . \quad (2)$$

This quantity measures to what extent a site is variable or conserved across the input MSA.

- The context-dependent entropy (CDE) is defined for site  $i$  in sequence  $A$  as

$$\text{CDE}_i^A = - \sum_{a=1}^{21} P_i(a|A_{\setminus i}) \log_2 P_i(a|A_{\setminus i}) . \quad (3)$$

Using the conditional probability  $P_i(a|A_{\setminus i}) = P(a_i = a|A_{\setminus i})$  of having amino acid  $a$  on site  $i$  given the rest of the sequence  $A_{\setminus i} = (a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_L)$ . This quantity cannot be extracted directly from the input MSA and needs to be obtained from the model parameters; due to the epistatic couplings in the energy  $\mathcal{H}(A)$  it actually differs from the CIE. As indicated explicitly, the  $\text{CDE}_i^A$  depends on the site  $i$  and the sequence  $A$ , hence on the context in which site  $i$  finds itself. As a matter of fact, this metric quantifies the local mutability of a site, i.e. its mutability within a certain reference-sequence context  $A_{\setminus i}$ .

A large value of CIE or CDE means that many mutations are allowed, while a small value means that only one or a few amino acids are tolerated.

In terms of these quantities, a site  $i$  in a given background sequence  $A$  can thus be classified as follows:

- *Mutable sites* have a large CDE and a large CIE. These sites can tolerate many mutations, both in the natural MSA and in the considered background.
- *Conserved sites* have a small CDE and a small CIE. These sites do not tolerate mutations, neither in the natural MSA nor in the considered background.
- *Epistatically constrained sites* have a small CDE and a large CIE. These sites display a large variability in the natural alignment, but only because the rest of the sequence is mutating at the same time. In fact, in the considered background, only one or a few amino acids are tolerated.

Note that it is very rare that a site has a CDE larger than the CIE, because generally speaking, fixing the background reduces the number of mutations that can be tolerated [38]. We also stress that this classification depends on the background  $A$ , and [29] has shown that epistatically constrained sites in a background can be mutable in another background and vice versa (while conserved sites tend to remain so in all backgrounds).

Furthermore, and most importantly for the present work, [29] has considered a given background  $A_0$  as the ancestral sequence, and starting from it has performed many parallel evolutions *in silico*, looking at how each site diversifies in the library of mutants obtained after a certain evolutionary time. It was found that mutable sites evolve very rapidly and quickly reach their asymptotic mutability, i.e. the CIE. Conserved sites remain so at any time during evolution, hence do not display interesting dynamics. Epistatically constrained sites, instead, are conserved at short evolutionary times, due to the epistatic constraints in the background of the ancestor, which cause a small CDE. But, as soon as the background sequence mutates enough, they can tolerate more mutations, asymptotically reaching the large CIE that they display in the natural alignment. Their mutation, however, is contingent on several mutations happening in the rest of the sequence, which often require a rather long evolutionary time to take place [29].

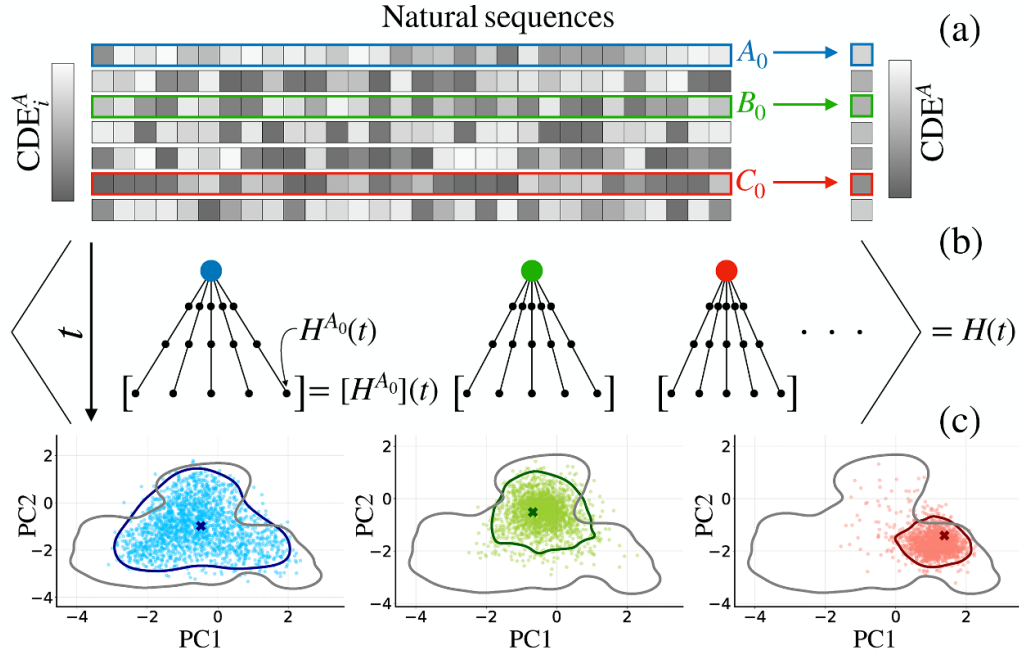
Along with the classification of sites discussed above, we also introduce a global measure of the level of epistasis in a given sequence  $A$ . Because each sequence has a different set of variable and epistatically constrained sites, we measure its overall level of epistatic constraints by averaging the CDE over sites,

$$\text{CDE}^A = \frac{1}{L} \sum_{i=1}^L \text{CDE}_i^A . \quad (4)$$

Hence, a more epistatically constrained sequence will display a lower  $\text{CDE}^A$ .

### 3.3. Dynamical fluctuations

Building on the previous work described above, and inspired by works on disordered systems [40–49], we devised a protocol (sketched in figure 1) to quantify how the amount of epistasis in the ancestral sequence impacts the subsequent evolution. We consider a series of ancestral sequences, taken from the natural MSA, let us call them  $A_0, B_0, C_0$ , etc. Each ancestor



**Figure 1.** Schematic representation of (a) an MSA in which each site is colored based on its mutability as measured by the  $CDE_i^A$  and (b) the protocol of our evolutionary dynamics with the relevant averages. (c) Depending on the epistatic constraints acting on the ancestral sequence and for equal evolutionary times, we observe a different level of diversity of the evolved sequences (colored), corresponding to a different exploration of the functional sequence space (gray).

has a different level of epistatic constraints, as measured by the average CDE introduced in equation (4), as illustrated in figure 1(a).

Next, we consider a star phylogeny of independent MCMC evolutionary trajectories, all starting from the same ancestor and evolving in parallel (figure 1(b)). Hence, for each ancestor, we construct, as a function of evolutionary time, an MSA of descendants that mimic those obtained through *in vitro* evolution experiments. As a measure of diversity, we focus on the evolution of the Hamming distance, between the ancestor  $A_0$  and the MSA of evolving sequences at time  $t$ ,  $A_t = (a_1^t, \dots, a_L^t)$ , defined as

$$H^{A_0}(t) = H(A_t, A_0) = \frac{1}{L} \sum_{i=1}^L (1 - \delta_{a_i^t, a_i^0}), \quad (5)$$

introducing the Kronecker delta symbol  $\delta_{a_i^t, a_i^0}$ , which is zero if site  $i$  is mutated between the two sequences  $A_t$  and  $A_0$ , and one otherwise. This quantity corresponds to the ‘overlap’ in the disordered systems literature and to the number of accepted mutations in evolution. Other measures could be considered as well, but we focus on this one for simplicity. For a given ancestor, simulating many parallel evolutionary trajectories, we thus obtain a set of realizations of the random variable  $H^{A_0}(t)$ .

As illustrated through a projection in PCA space in figure 1(c), we observe that the level of epistasis in the ancestor determines the size of the portion of the fitness landscape explored by evolution. For comparable evolutionary time, ancestors with more epistatically constrained sites (e.g. sequence  $C_0$  in figure 1) lead to a less diverse set of evolved

sequences. Conversely, ancestors with less epistatically constrained sites (e.g.  $A_0$  in figure 1) lead to a widely diverse set of evolved sequences.

To make this observation quantitative, we characterize the statistical properties of the resulting MSAs, and how they depend on the ancestor, by introducing two distinct averages and corresponding fluctuations, inspired from the disordered systems literature [45, 46, 48–50] and illustrated in figure 1(b). Recall that the number of accepted mutations  $H^{A_0}(t)$  at evolutionary time  $t$  for fixed ancestor  $A_0$  is a random variable, whose realizations depend on the stochasticity of the evolutionary trajectory.

- First, we consider averaging over many evolutionary trajectories for fixed ancestor. We denote this average as  $\langle \dots \rangle$ . The variance of  $H^{A_0}(t)$  can then be defined as

$$\chi_{\text{dyn}}^{A_0}(t) = \langle [H^2] \rangle - [H]^2, \quad (6)$$

where the dependence of  $H^{A_0}(t)$  on  $A_0$  and  $t$  is omitted to simplify the notation. The quantity  $\chi_{\text{dyn}}^{A_0}(t)$  depends on the ancestor  $A_0$  and on time  $t$ .

- Second, we consider the average over different ancestors, which we define as  $\langle \dots \rangle$ . In particular, we consider the average number of mutations  $\langle [H^{A_0}(t)] \rangle$  for a given ancestor, that measures the average diversity of the evolved MSA, and we measure how this quantity fluctuates from ancestor to ancestor via the variance

$$\chi_{\text{bg}}(t) = \langle \langle [H^2] \rangle \rangle - \langle [H] \rangle^2, \quad (7)$$

where again the dependence of  $[H] = [H^{A_0}(t)]$  on  $A_0$  and  $t$  is omitted to simplify the notation. This quantity is the variance

of  $[H]$  associated to the fluctuations in the background of the ancestral sequence  $A_0$ , hence the suffix ‘bg’.

Note that the total variance of  $H^{A_0}(t)$  over both sources of randomness, i.e. the random choice of ancestor and the stochasticity of mutations along the evolutionary trajectory, can be decomposed as

$$\begin{aligned}\chi_{\text{tot}}(t) &= \langle [H^2] \rangle - \langle [H] \rangle^2 \\ &= \langle [H^2] \rangle - \langle [H^2] \rangle + \langle [H^2] \rangle - \langle [H] \rangle^2 \\ &= \chi_{\text{bg}}(t) + \chi_{\text{dyn}}(t),\end{aligned}\quad (8)$$

where

$$\chi_{\text{dyn}}(t) = \langle \chi_{\text{dyn}}^{A_0}(t) \rangle \quad (9)$$

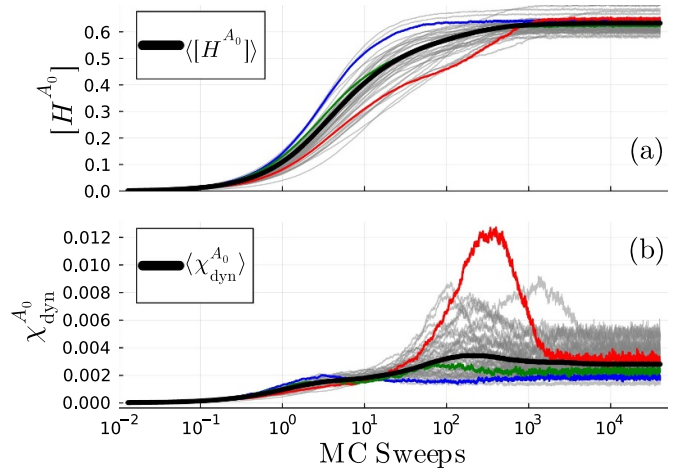
is the average of the ancestor-dependent  $\chi_{\text{dyn}}^{A_0}(t)$  over  $A_0$ . Using these quantities, which are called ‘dynamical susceptibilities’ in the physics of disordered systems, we can thus quantify the relative importance of the ancestor and the evolutionary noise in determining the diversity of the resulting MSAs at a fixed evolutionary time. In the SM we also study some alternative definitions of susceptibilities, based for example on the Hamming distance between two chains starting from the same ancestor. However, because our MCMC algorithm is time-reversible, these two quantities are related.

## 4. Results

### 4.1. Mutational dynamics and its fluctuations

The dynamical evolution of the Hamming distance  $H^{A_0}(t)$  (number of accepted mutations with respect to the ancestor) is displayed as a function of the number of Monte Carlo sweeps in figure 2, with a sweep corresponding, on average, to one attempted mutation per site. More specifically, figure 2(a) shows the average  $\langle [H^{A_0}(t)] \rangle$  over a star phylogeny of  $10^3$  parallel evolutions, for 200 choices of the ancestor  $A_0$  taken at random (with weights, see SM) from the DBD protein family. Figure 2(b) reports the variance  $\chi_{\text{dyn}}^{A_0}(t) = [H^2] - [H]^2$  defined in equation (6) over the same phylogeny, and for the same initial sequences as in figure 2(a).

Before describing three interesting cases that we highlighted with colors, we discuss the general traits of such dynamics. The average Hamming distance and its variance are both initially null because all chains start from the same ancestral sequence. Because the DCA model is generative and the evolutionary dynamics respects detailed balance, at large times we expect the simulated sequences to be independent samples from the DCA model, which then reproduce statistical features of the natural ones used for training. This means that the average of  $H^{A_0}(t)$  and its variance converge, respectively, to the average and the variance of the Hamming distance between the chosen ancestral sequence and the rest of the natural ones. These values can vary significantly (the final average Hamming distance in figure 2(a) varies from roughly 0.55–0.65) depending on how close  $A_0$  is to the other sequences



**Figure 2.** Average (a) and variance (b) of the Hamming distance between the evolving sequence and the ancestor estimated using  $10^3$  independent trajectories, for many different ancestors (gray lines). The thick black line represents the average over ancestors, i.e.  $\langle [H^{A_0}(t)] \rangle$  and  $\chi_{\text{dyn}}(t) = \langle \chi_{\text{dyn}}^{A_0}(t) \rangle$ , over 200 ancestors. The green, blue and red lines highlights the same specific choices of ancestor used in figure 1.

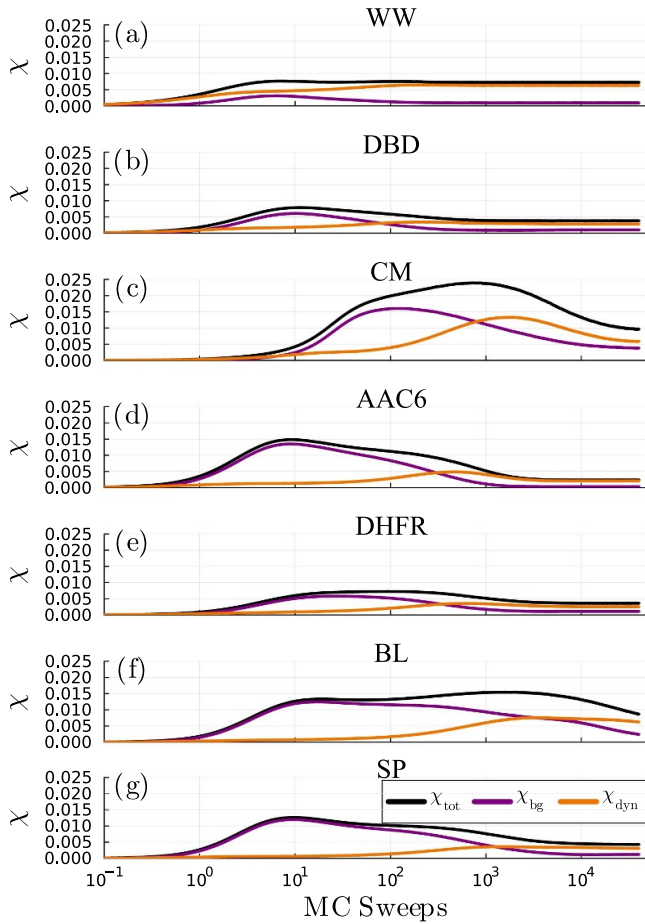
in the natural MSA. The different curves, corresponding to distinct ancestors  $A_0$ , show a wide range of behaviors and time scales: some trajectories reach the steady state rapidly while others take much longer, displaying intermediate plateaus as a hallmark of epistasis. In some cases  $\chi_{\text{dyn}}^{A_0}(t)$  displays a peak and then decreases, while in others the equilibrium value is reached in a monotonic way. The large value reached by  $\chi_{\text{dyn}}^{A_0}(t)$  for some initial sequences implies strong fluctuations between distinct evolutionary trajectories, making it hard to predict the dynamics.

In both panels of figure 2, we highlighted with colors some representative curves. In particular, the green curve corresponds to an ancestor that behaves in a rather ‘typical’ way, close to the average. The blue curve corresponds to an ancestor that has less epistatically constrained sites, hence the dynamics is faster and less heterogeneous. Finally, the red curve corresponds to a highly epistatically constrained ancestor, which leads to a slower dynamics with an intermediate plateau in  $[H^{A_0}(t)]$ , and a large peak in  $\chi_{\text{dyn}}^{A_0}(t)$ . The same three ancestors have been used to construct the PCA plots in figure 1 that correspond to time  $t = 50$  MC Sweeps.

### 4.2. Dynamical heterogeneity across families

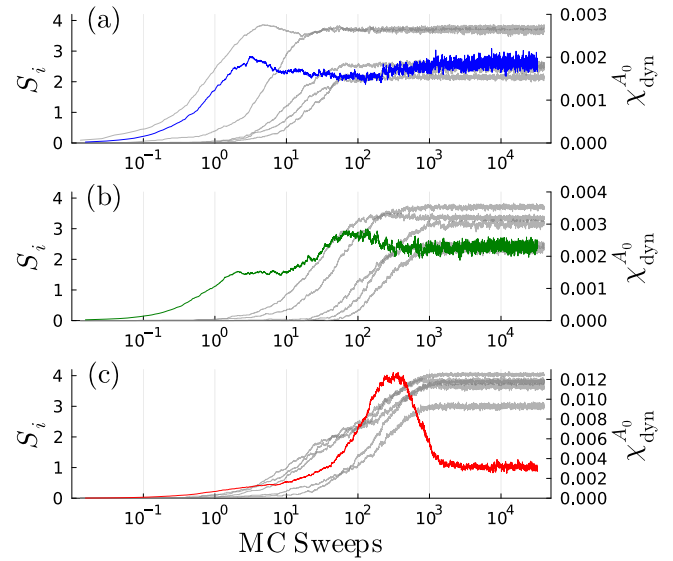
In figure 3 we generalize our results to many protein families, chosen to have different ranges of sequence length, MSA depth (number of natural sequences) and equilibration time scales. These families are also interesting because for some of them, experimental data from Deep Mutational Scans and *in vitro* evolution are available.

For compactness, we only display the average dynamical fluctuations, with  $\chi_{\text{tot}}$  corresponding to the total variance of  $H$ ,  $\chi_{\text{dyn}}$  to the variance due to the stochasticity of the evolution



**Figure 3.** Evolution of the different susceptibilities  $\chi_{\text{tot}}$ ,  $\chi_{\text{bg}}$ , and  $\chi_{\text{dyn}}$ , for different protein families. (a) WW domain, (b) DNA-binding domain (DBD), (c) chorismate mutase (CM), (d) AAC6, (e) DHFR enzyme, (f) beta-lactamase (BL), and (g) serine protease (SP).

(averaged over the ancestor), and  $\chi_{\text{bg}}$  to the variance due to the choice of ancestor. All the  $\chi_{\text{tot}}$  curves display a similar behavior, increasing from zero to a maximum value that is maintained over one or more decades, before finally decreasing to a smaller value corresponding to equilibrium. What matters most to us, however, is the relative importance of the two terms  $\chi_{\text{dyn}}$  and  $\chi_{\text{bg}}$  in which  $\chi_{\text{tot}}$  can be decomposed. The key observation is that background-related fluctuations dominate over dynamical ones at short evolutionary times, while the inverse is true at larger times. This means that, at least over time scales for which  $\chi_{\text{bg}} \gg \chi_{\text{dyn}}$ , we can hope to reconstruct with good accuracy the ancestral sequence from which a set of evolutionary trajectories started. Conversely, at larger time scales, the dynamical noise contribution dominates and the trajectory-to-trajectory fluctuations are large enough to hide the signal coming from the ancestral sequence, precluding the possibility to reconstruct it. This behavior is observed in all the protein families that we tested, with the exception of the WW domain in which  $\chi_{\text{bg}}$  never grows to larger values than  $\chi_{\text{dyn}}$ . We attribute this difference to the small size of the WW domain, which prevents the sequences to accumulate enough



**Figure 4.** Evolution of the entropy of the most epistatically constrained sites for a specific ancestral sequence of DBD plotted together with  $\chi_{\text{dyn}}^{A_0}$  of the same sequence. (a)–(c) are for the blue, green, and red sequences in figure 2, respectively.

epistatic interactions. However, both the time scales and the relative importance of the two terms contributing to the total susceptibility vary greatly from family to family (figure 3). In DBD, AAC6, DHFR, and SP the peak of  $\chi_{\text{tot}}$  is reached quite early in the dynamics and the two contributions  $\chi_{\text{dyn}}$  and  $\chi_{\text{bg}}$  are almost non-overlapping. On the other hand, for CM and BL, the peak occurs much later due to the two contributions having a significant overlap.

For a given family and a typical ancestor, the time scale at which  $\chi_{\text{dyn}}$  and  $\chi_{\text{bg}}$  cross, the former becoming larger than the latter, defines a characteristic evolutionary time scale, around which the memory of the ancestor is lost.

#### 4.3. Epistatic constraints and dynamical fluctuations

Up to this point, we discussed the time dependence of the different contributions to the total fluctuations. We argued that, as long as the dominant contribution to  $\chi_{\text{tot}}$  comes from  $\chi_{\text{bg}}$ , the ancestral sequence is strongly related to the evolving ones. We now want to better understand the transition to the final regime in which dynamical fluctuations, i.e.  $\chi_{\text{dyn}}$ , dominate.

In figure 4 we show, for the three ancestral sequences highlighted in figure 2, how the dynamical part of the susceptibility is strongly related to the epistatically constrained sites (defined as in section 3.2 and in [29]). To make this relation clear, we first checked that the time scale at which the dynamical susceptibility reaches its peak is compatible with the time scale at which epistatically constrained sites evolve. For a given ancestor  $A_0$ , we consider the five most epistatically constrained sites, i.e. those with the largest  $\text{CIE}_i - \text{CDE}_i^{A_0}$ . For these sites, we consider at time  $t$  the frequency of appearance of amino acid  $a$ ,  $f_i^t(a)$ , in the set of sequences that evolved from  $A_0$ , and from it we compute a time-dependent entropy

$$S_i(t) = - \sum_{a=1}^{21} f_i(a) \log_2 f_i(a) . \quad (10)$$

In figure 4,  $\chi_{\text{dyn}}^{A_0}(t)$  (colored curve) is superposed to  $S_i(t)$  of the five sites (gray curves), as a function of evolutionary time  $t$ . We observe that the peak of  $\chi_{\text{dyn}}^{A_0}$  is reached just before the equilibration of the epistatically constrained sites, i.e. the time at which  $S_i(t)$  approaches the CIE $_i$ . The time of equilibration of epistatically constrained sites also increases from figure 4(a) to figure 4(c), similarly to what happens for the equilibration of the Hamming distance in the same three sequences in figure 2. Epistasis does not only affect the time scale at which the peak of  $\chi_{\text{dyn}}^{A_0}$  is reached, but also the intensity of the peak, which means that more epistatic ancestors lead to a more heterogeneous dynamics at intermediate times, when the epistatic sites mutate.

To quantify these observations, in figure 5 we consider the same set of families as in figure 3, namely WW, DBD, CM, AAC6, DHFR, BL, and SP, spanning a wide range of sequence length (from  $L = 31$  in WW to  $L = 220$  in SP). For each family, we consider a set of ancestors  $A_0$  and we report a scatter plot of the following quantities, each  $A_0$  being a point:

- the value of  $\text{CDE}^{A_0}$  defined as in equation (4),
- the maximum value  $\max_t \chi_{\text{dyn}}^{A_0}(t)$ , indicated as  $\max(\chi_{\text{dyn}}^{A_0})$  for simplicity,
- and the time  $t_{90}$  at which the average  $[H^{A_0}]$  over the trajectories reaches 90% of the equilibrium ( $t \rightarrow \infty$ ) value.

We observe that both  $\max(\chi_{\text{dyn}}^{A_0})$  and  $t_{90}$  markedly increase as  $\text{CDE}^{A_0}$  decreases. Hence, a ‘highly epistatic’ sequence with a small value of  $\text{CDE}^{A_0}$  (with respect to the typical value of the family) has many sites with low context-dependent entropy, which cannot mutate at the beginning of evolution, leading to a slower and more heterogeneous overall dynamics. This is for example the case of the red sequence in figure 2, for which we have  $\text{CDE}^{A_0} = 1.45$ , while for the blue one  $\text{CDE}^{A_0} = 2.11$ . The three sequences considered in Figure 2 are indicated with stars of the corresponding colors in figure 5(b).

Epistatically constrained sites also carry the specific signature of the ancestral sequence, because the conserved sites are roughly the same for every sequence of the family and the variable ones carry little to no information. These results thus suggest that tracing back an evolutionary trajectory to its ancestral sequence becomes more difficult as one approaches the peak of  $\chi_{\text{dyn}}$ , because this is when the epistatic sites, i.e. the sites that carry information about that initial sequence, start to mutate.

#### 4.4. Cooperative mutational dynamics

We have established that, when epistatic sites start to mutate,  $\chi_{\text{dyn}}^{A_0}$  reaches its peak, and that the intensity of this peak is correlated with the amount of epistatically correlated sites in the ancestral sequence. We now want to show that the peak

of  $\chi_{\text{dyn}}^{A_0}$  is due to dynamical correlations between different sites, caused by the strong interaction between those sites and the context. These correlations result in a cooperative mutational process, in which epistatically constrained sites can only mutate because other such sites mutate, leading to an avalanche of mutations.

To precisely quantify this effect, one can interpret the dynamical susceptibility as the sum of a site-site dynamical correlation function [47]. The Hamming distance is defined in equation (5). Inserting its expression into equation (6), we obtain

$$\chi_{\text{dyn}}^{A_0}(t) = \frac{1}{L^2} \sum_{ij} G_{ij}^{A_0}(t) , \quad (11)$$

with

$$G_{ij}^{A_0}(t) = \left[ \delta_{a_i^t, a_i^0} \delta_{a_j^t, a_j^0} \right] - \left[ \delta_{a_i^t, a_i^0} \right] \left[ \delta_{a_j^t, a_j^0} \right] , \quad (12)$$

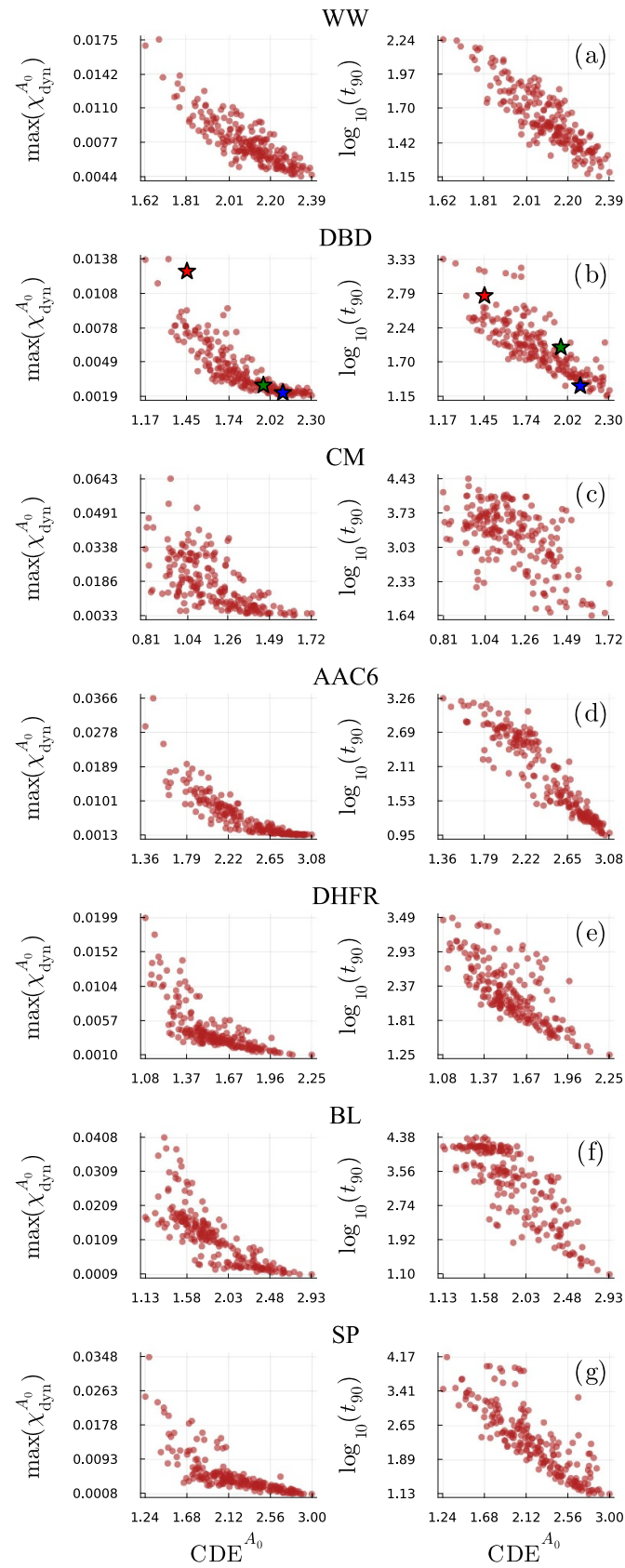
recalling that the Kronecker delta symbol  $\delta_{a_i^t, a_i^0}$  is zero if site  $i$  is mutated between the two sequences  $A_t$  and  $A_0$ , and one otherwise. The matrix  $G_{ij}^{A_0}$  is a time-dependent correlation matrix between sites  $i$  and  $j$ . It is high when the sites are dynamically correlated, i.e. if  $a_i^t$  is different from  $a_i^0$ , then  $a_j^t$  is likely different from  $a_j^0$  as well, and vice versa. It is low when the sites mutate independently. Because  $\chi_{\text{dyn}}^{A_0}$  is the sum of  $G_{ij}^{A_0}$  over all pairs  $(i, j)$ , a large value of  $\chi_{\text{dyn}}^{A_0}$  implies that many  $(i, j)$  are strongly correlated. Furthermore, by looking at this correlation matrix one can infer which pairs of sites mutate in a correlated way during evolution.

The value of the off-diagonal part of  $G_{ij}^{A_0}(t^*)$ , computed at the time  $t^*$  at which a peak in the dynamical susceptibility is observed, is shown in figure 6 for six different starting sequences in the DBD family.

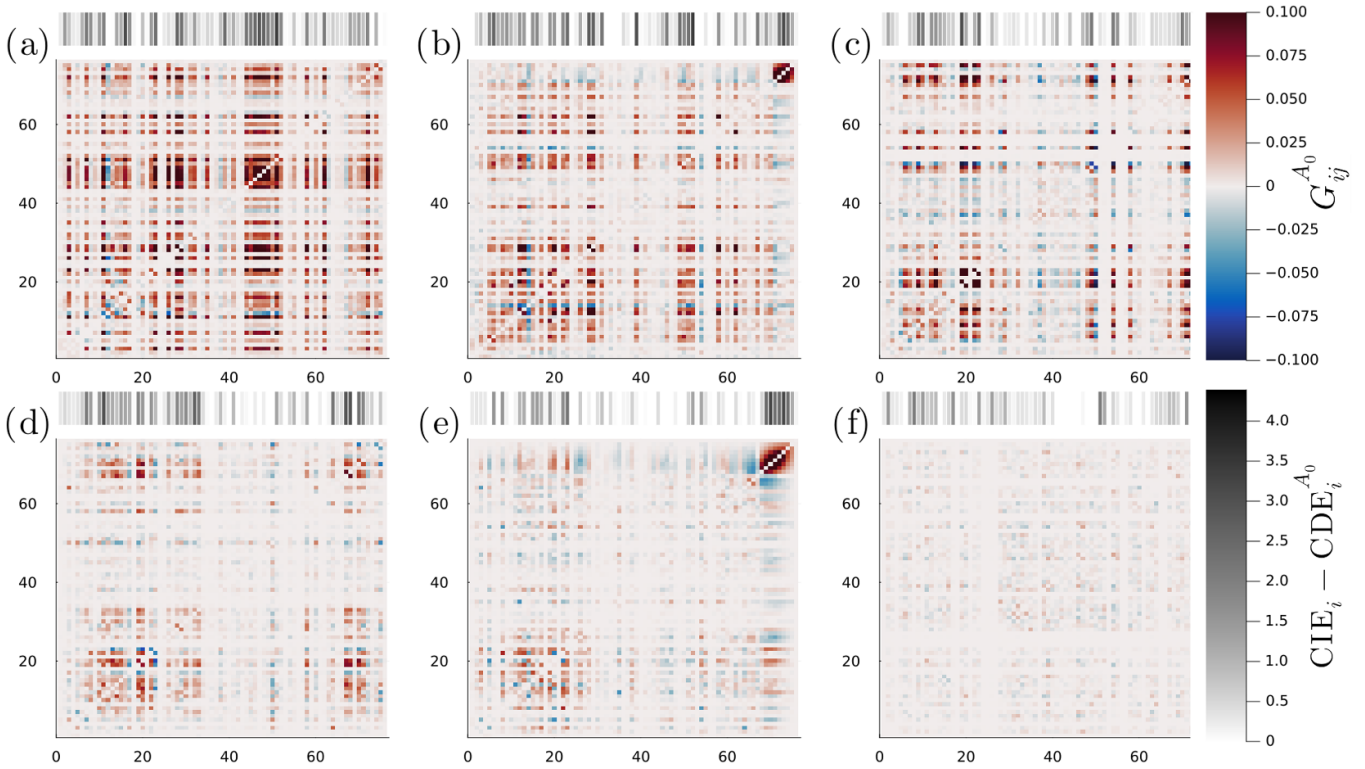
Figure 6(a) corresponds to an ancestor (red in figure 2) for which a large peak and a longer time to equilibrate are observed in the evolution of  $\chi_{\text{dyn}}^{A_0}$ . The large value of  $\chi_{\text{dyn}}^{A_0}$ , which is just the sum of all matrix elements, implies a large number of positively correlated sites. Figure 6(e) corresponds to an ancestor (green in figure 2) for which  $\chi_{\text{dyn}}^{A_0}$  has a smaller peak before saturating. We notice in this case that  $G_{ij}^{A_0}(t^*)$  still displays some strong correlations, but the majority of the sites are uncorrelated. Finally, figure 6(f) corresponds to an ancestor (blue in figure 2) for which  $\chi_{\text{dyn}}^{A_0}$  saturates to a small value. In this case correlations are almost absent at  $t^*$ , which means that essentially all sites mutate independently. It is thus evident how different initial sequences give rise to different patterns in the  $G_{ij}^{A_0}$  matrix. At large times the chains reach equilibrium and are almost statistically indistinguishable from natural sequences. This means that  $G_{ij}^{A_0}(t \rightarrow \infty) \approx G_{ij, \text{nat}}^{A_0}$ , with

$$G_{ij, \text{nat}}^{A_0} = \langle \delta_{a_i^{\text{nat}}, a_i^0} \delta_{a_j^{\text{nat}}, a_j^0} \rangle_{\text{nat}} - \langle \delta_{a_i^{\text{nat}}, a_i^0} \rangle_{\text{nat}} \langle \delta_{a_j^{\text{nat}}, a_j^0} \rangle_{\text{nat}} ,$$

where  $\langle \dots \rangle_{\text{nat}}$  is the average computed over the natural sequences.



**Figure 5.** Scatter plot of  $\text{CDE}^{A_0}$  against the maximum of  $\chi_{\text{dyn}}^{A_0}$  (left) and against  $\log_{10}(t_{90})$  (right) for each protein family we considered. For each family, the points correspond to 200 sequences  $A_0$  extracted from the natural MSA with weights (see SM). In (b) colored stars are used to highlight the sequences that correspond to the blue, green, and red curves in figure 2.



**Figure 6.** Correlations between sites at time  $t^*$  at which  $\chi_{\text{dyn}}^{A_0}(t)$  reaches its maximum (with a cutoff at  $t_{\text{th}} = 1000$  sweeps), for different ancestral sequences. The bar above each snapshot encodes the strength of epistatic constraints for the sites in the ancestral sequence.

In section 4.3, we showed that the presence of a peak in  $\chi_{\text{dyn}}^{A_0}$  is correlated to  $\text{CDE}^{A_0}$ , i.e. to the amount of epistatically constrained sites in the ancestral sequence. In this section, we show that these sites also display large dynamical correlations between themselves. The bar above each plot in figure 6 is shaded to represent  $\text{CIE}_i - \text{CDE}_i^A$ , capturing the information about site  $i$  encoded in the background sequence  $A_{\setminus i}$ . This coloring reflects the degree to which each site is constrained by epistatic interactions. The conserved and variable sites are indicated in white, while the epistatically constrained ones are shown in black. As expected, we see good agreement between these residues and the ones that give a large contribution to  $\chi_{\text{dyn}}^{A_0}$ . Hence, we conclude that epistatically constrained sites mutate cooperatively around the time scale corresponding to the peak of  $\chi_{\text{dyn}}^{A_0}$ .

#### 4.5. Response to environmental variation

In the literature on disordered physical systems, it has been established that the dynamical correlations discussed in sections 4.3 and 4.4 are related to the linear response of the dynamics to a change in temperature through a kind of fluctuation-dissipation relation [44].

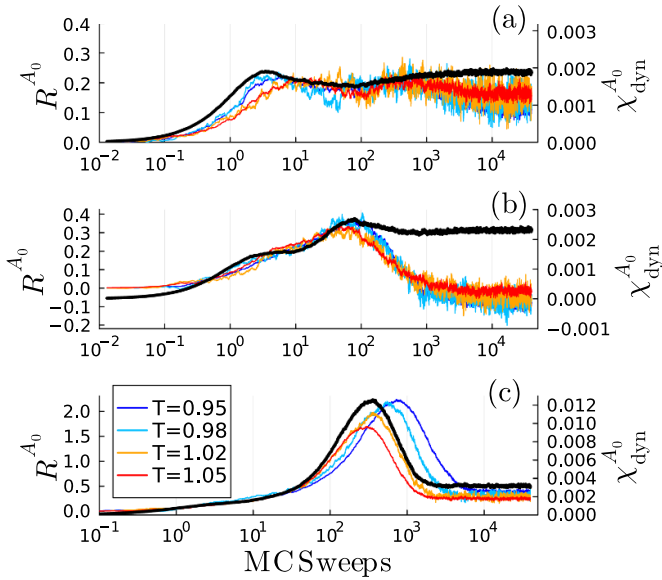
In the present context, the ‘temperature’ is a parameter that controls the probability with which mutations are accepted. High temperature corresponds to all mutations being accepted (very low selection), while low temperature corresponds to only the most beneficial mutations being accepted (very

strong selection, i.e. directed evolution). The value of  $T = 1$  corresponds to the conditions at which the DCA model is trained on natural sequences, hence  $T \sim 1$  corresponds to a neutral drift dynamics during which the neutral space of sequences that have comparable fitness to natural ones is explored [25]. Indeed, [25] has shown that two distinct *in vitro* evolution experiments, realized with selection pressures comparable to the natural one, can be described by fitting the temperature in a range slightly above  $T \sim 1$ .

Hence, in our modeling framework, a change of temperature corresponds to a change of selection strength. Following [44], we check whether the dynamical susceptibilities are related to the response of the average dynamics with respect to the temperature variation. To estimate such linear response, we consider two evolutionary dynamics starting from the same ancestor, one at temperature  $T \neq 1$  and the other at temperature  $T = 1$ , and we compare the average number of accepted mutations  $[H^{A_0}(t; T)]$  between the two evolutions. The dynamical linear response is given by

$$R^{A_0}(t) = \lim_{T \rightarrow 1} \frac{[H^{A_0}(t; T)] - [H^{A_0}(t; T = 1)]}{T - 1}, \quad (13)$$

and we emphasize that this quantity depends on the ancestor  $A_0$  and on time  $t$ . In figure 7, we compare the time dependence of  $\chi_{\text{dyn}}^{A_0}(t)$  with  $R^{A_0}(t)$  for the same three ancestors as in figure 2. Note that we cannot take the limit  $T \rightarrow 1$  due to statistical noise, and we estimate  $R^{A_0}(t)$  by observing that the curves for several values of  $T$  close to 1 are almost superimposed to



**Figure 7.** Evolution of the dynamical susceptibility  $\chi_{\text{dyn}}^{A_0}(t)$  compared with the linear response of the average Hamming distance to a change of temperature  $R^{A_0}(t)$ . The three panels correspond to the three ancestral sequences highlighted in figure 2, (a) blue, (b) green, (c) red.

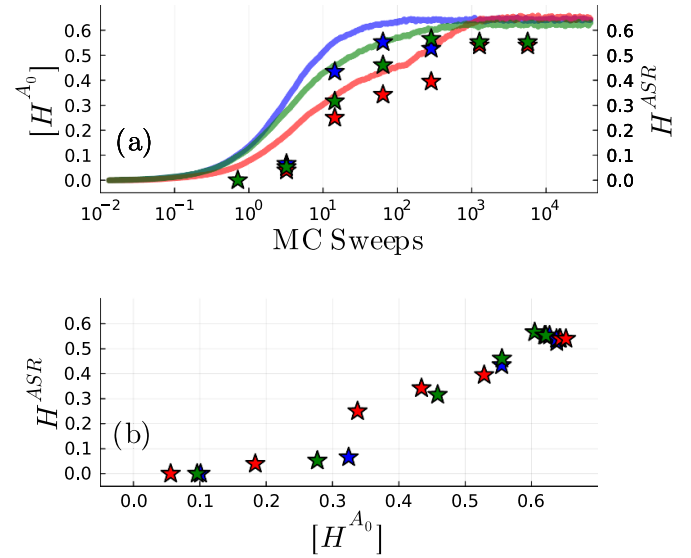
each other. Figure 7 shows that the fluctuation-dissipation relation seems to hold quite well, such that  $\chi_{\text{dyn}}^{A_0}(t) \propto R^{A_0}(t)$ , at least for times  $t$  smaller than the peak of both quantities.

The results in figure 7 suggest that more epistatically constrained ancestral sequences, which display a stronger peak in the dynamical susceptibility, will also display a stronger response of the evolutionary dynamics to a small change in selection strength. Such response is stronger around the time of the peak in  $\chi_{\text{dyn}}^{A_0}(t)$ , which is also the time at which epistatically constrained sites mutate cooperatively, see section 4.4. In other words, *sequences with more epistatically constrained sites are more sensitive to environmental changes*. In the SM, we give an additional analytical argument based on the short-time dynamics, that supports the same conclusion.

#### 4.6. Limits of ancestral sequence reconstruction (ASR)

In an evolutionary process, predictions can be made in two directions: either starting from the ancestral sequence and predicting its future evolution over a given time; or starting from a set of evolved sequences and attempting to infer the ancestral sequence from which the evolution began. An attempt at predicting future evolution has been made in [25] (see [51–55] for related studies in disordered systems), where however the prediction was limited to small evolutionary time scales. Because we have repeatedly hinted that our analysis can provide insight into the time scale over which reconstructing the ancestral sequence of an evolutionary process is possible, before concluding, we consider here the ASR problem more explicitly.

For each of the three sequences we analyzed in figure 2, we simulate 200 independent evolutionary trajectories, thus obtaining several MSAs of 200 evolved sequences at different



**Figure 8.** (a) The full lines represent the evolution of the average Hamming distance for the same three ancestors as in figure 2. The star symbols represent the Hamming distance  $H^{\text{ASR}}$  between the ground truth ancestor and the one reconstructed using the simulated MSAs at different times. (b) Hamming distance  $H^{\text{ASR}}$  between the reconstructed ancestor and the ground truth as a function of the average distance  $[H^{A_0}(t)]$  between the ancestor and the MSA used for the reconstruction. The colors refer to the same sequences of panel (a).

evolutionary times  $t$ . Then, we use the FastML online tool [80] to infer the ancestral sequence starting from those MSAs and using a star-shaped phylogenetic tree. We then check how close the reconstructed sequence is to the actual ancestor.

The results are shown in figure 8(a), where we plot the evolution of the average Hamming distance for the same sequences as in figure 2, and we compare it with the Hamming distance (i.e. the percentage of errors) between the ground truth ancestor and the reconstructed one. At the beginning of the dynamics the FastML tool is able to reconstruct perfectly the original sequence in all of the three cases. This is not surprising, as at short times few mutations appear. The situation is similar at large times, where the MSAs generated from the three starting sequences are statistically indistinguishable and the FastML tool performs equally bad in all of them [81, 82]. Thus, it is more insightful to look at intermediate times, where a significant difference between the evolving sequences is present. After  $\sim 10$  sweeps the variable sites begin to mutate and the ASR tool cannot reproduce the initial sequence perfectly. The leap in the error percentage with respect to the previous point is however larger for the blue sequence, which is the one with fewer epistatic sites, while the red one is still recovered reasonably well. This difference is present for the subsequent points as well, until the error saturates. In figure 8(b) we show that the Hamming distance (percentage of errors) of the sequence reconstruction results grows roughly proportionally to the average sequence divergence of the evolutionary trajectories for all of the three studied sequences, but more so for the more epistatic red sequence. We conclude that, at least with the FastML tool, the possibility of reconstructing

ancestral sequences is mostly determined by the Hamming distance of the ancestor from the evolving sequences, i.e.  $[H^{A_0}(t)]$ , with more errors accumulating when using datasets of more diverged sequences.

Our analysis shows that, for a given evolutionary time, the amount of diversity  $[H^{A_0}(t)]$  depends strongly on the ancestor, due to highly non-trivial epistatic dynamical correlations. More epistatically constrained ancestors give rise to less diversity, thus allowing for reconstruction over longer evolutionary times (figure 8(a)). Yet, at comparable amount of diversity, more epistatic ancestors are more difficult to reconstruct (figure 8(b)), at least using the FastML algorithm that neglects correlations between sites [83]. These results confirm once more that the evolutionary time limit for predictability is linked to the presence of epistatic sites that carry information about the starting sequence only up to a certain sequence divergence. We further observe that more advanced tools (yet to be fully developed) for ASR might exploit the epistatic correlations to achieve a better reconstruction [84], which is not done in FastML and similar algorithms. We expect that such tools would perform better on more epistatic ancestors such as the red one in figure 8.

## 5. Conclusions

In this study, we explore the role of disorder and stochasticity in protein evolution, focusing on the interplay between epistasis, site correlations, and sequence-dependent fluctuations. We simulated evolution *in silico* using DCA to model the fitness landscape, and an MCMC algorithm to mimic the process of mutation and selection, a methodology that has been validated in previous work [24, 25, 29]. By employing tools from statistical physics [40–49], we quantified how the ancestral amino acid sequence significantly influences early evolutionary dynamics, as measured by the dynamical susceptibility. Our main results are the following.

- Our analysis shows that, during the initial phase of evolution, there is very significant heterogeneity in the dynamics, depending on the choice of ancestral sequence. Some ancestral sequences lead to a smooth evolution where each site mutates almost independently of the others, while others lead to a more complex dynamics characterized by intermediate plateaus and significant fluctuations (section 4.1).
- We have shown that, for a variety of protein families, the noise arising from the starting sequence dominates over stochastic evolutionary fluctuations, implying that one can, in principle, trace the evolutionary trajectory back to its origin (section 4.2). This effect is more or less pronounced, depending on the family, and our tools allow one to quantify it. However, as time progresses and epistatically constrained sites evolve, this traceability is lost due to the growing influence of the mutational stochasticity.
- The heterogeneity of the initial evolutionary dynamics can be traced back to the amount of epistatic constraints in the ancestral sequence. We introduced a quantity, the average over sites of the context-dependent entropy, and we have

shown that this quantity is strongly correlated with the time at which dynamical heterogeneity reaches its peak and with the strength of the fluctuations at the peak (section 4.3). This method for assessing sequence evolvability could guide experimentalists in selecting an appropriate starting sequence for experiments, based on the desired functional behavior.

- The amplitude of the global fluctuations can be expressed in terms of a sum of pairwise dynamical correlations between pairs of residues. More epistatically constrained ancestral sequences show groups of residues that mutate collectively, and those patterns can be identified from the correlation matrix  $G_{ij}^{A_0}(t)$  we introduced in section 4.4. The observed correlations between sites thus reflect a complex epistatic landscape, where certain residues evolve in a highly context-dependent manner. These correlations could also be leveraged in experimental settings; for instance, one may imagine to control the evolvability of a specific residue by targeting sites that exhibit strong correlations with it.
- We demonstrated that these epistatic sites, which evolve over long time scales, significantly affect the response of the dynamics to a change in environmental conditions. More epistatically constrained sequences lead to a larger response to a change in environment, suggesting potential implications on sequence evolvability imposed by environmental pressures, such as antibiotic concentration (section 4.5).
- Finally, we presented a preliminary study of the performance of an ASR algorithm (here, FastML), in light of our previous findings. We found that more epistatically constrained ancestors lead to less diversity at comparable time scales, which facilitates their reconstruction. Yet, at comparable diversity, they are more difficult to reconstruct. We believe that this analysis will be instrumental in improving the efficiency of ASR algorithms, which could in principle exploit the correlations identified in this work.

More generally, our findings extend the analogy between protein evolution and disordered physical systems, reinforcing the idea that protein dynamics exhibit characteristics of complex and strongly correlated systems. However, when comparing these results with statistical physics models exhibiting glassy dynamics, we also found qualitative differences, emphasizing the unique constraints imposed by natural evolution on proteins. Future work should aim to further elucidate the connection between epistasis, evolvability, and environmental selection, which may offer insights into evolutionary strategies across diverse biological systems. Moreover, more realistic evolutionary dynamics could be considered. In this paper we focused on independent Monte Carlo chains, which corresponds to evolution on a star tree. Introducing a more complex tree structure could affect the results presented here, as the varying distance between the leaves of the tree and its root and the correlations coming from a common ancestor may affect the computation of  $\chi_{\text{dyn}}^{A_0}$ . Furthermore, we believe that most of the ideas presented in this work will soon be amenable to experimental testing, thanks to the increased power of *in vitro* evolution platforms.

## Data availability statement

The data cannot be made publicly available upon publication because they are not available in a format that is sufficiently accessible or reusable by other researchers. The data that support the findings of this study are available upon reasonable request from the authors.

## Acknowledgments

We thank Maria Chiara Angelini, Ludovic Berthier, Pierre Barrat-Charlaix, Simona Cocco, Dongkyu Lee, Arvind Murugan, Clément Nizak, Misaki Ozawa, Andrea Pagnani, Olivier Rivoire, Joe Thornton, Nobuhiko Tokuriki, and Alya Zeinaty for fruitful discussions. This research has been supported by first FIS (Italian Science Fund) 2021 funding scheme (FIS783 - SMAc—Statistical Mechanics and Complexity) from MUR, Italian Ministry of University and Research and from the PRIN funding scheme (2022LMHTET—Complexity, disorder and fluctuations: spin glass physics and beyond) from MUR, Italian Ministry of University and Research.

## References

- [1] Blum M *et al* 2024 InterPro: the protein sequence classification resource in 2025 *Nucl. Acids Res.* **53** D444–56
- [2] Blum M *et al* 2024 InterPro: the protein sequence classification resource in 2025 *Nucl. Acids Res.* **53** D444–56
- [3] Consortium T U 2022 UniProt: the universal protein knowledgebase in 2023 *Nucl. Acids Res.* **51** D523–31
- [4] Burley S K *et al* 2022 RCSB protein data bank (RCSB.org): delivery of experimentally-determined PDB structures alongside one million computed structure models of proteins from artificial intelligence/machine learning *Nucl. Acids Res.* **51** D488–508
- [5] Kimura M 1983 *The Neutral Theory of Molecular Evolution* (Cambridge University Press)
- [6] Fowler D M and Fields S 2014 Deep mutational scanning: a new style of protein science *Nat. Methods* **11** 801–7
- [7] Sarkisyan K S *et al* 2016 Local fitness landscape of the green fluorescent protein *Nature* **533** 397–401
- [8] Notin P *et al* 2024 Proteingym: large-scale benchmarks for protein fitness prediction and design *Advances in Neural Information Processing Systems* vol 36
- [9] Hartl D L, Clark A G and Clark A G 1997 *Principles of Population Genetics* vol 116 (Sinauer Associates)
- [10] Harms M J and Thornton J W 2013 Evolutionary biochemistry: revealing the historical and physical causes of protein properties *Nat. Rev. Genet.* **14** 559–71
- [11] Weinreich D M, Lan Y, Wylie C S and Heckendorn R B 2013 Should evolutionary geneticists worry about higher-order epistasis? *Curr. Opin. Genet. Dev.* **23** 700–7
- [12] De Visser J A G M and Krug J 2014 Empirical fitness landscapes and the predictability of evolution *Nat. Rev. Genet.* **15** 480–90
- [13] Starr T N and Thornton J W 2016 Epistasis in protein evolution *Protein Sci.* **25** 1204–18
- [14] Poelwijk F J, Krishna V and Ranganathan R 2016 The context-dependence of mutations: a linkage of formalisms *PLoS Comput. Biol.* **12** e1004771
- [15] Cocco S, Feinauer C, Figliuzzi M, Monasson R and Weigt M 2018 Inverse statistical physics of protein sequences: a key issues review *Rep. Prog. Phys.* **81** 032601
- [16] Domingo J, Baeza-Centurion P and Lehner B 2019 The causes and consequences of genetic interactions (epistasis) *Annu. Rev. Genomics Hum. Genet.* **20** 433–60
- [17] Johnson M S, Reddy G and Desai M M 2023 Epistasis and evolution: recent advances and an outlook for prediction *BMC Biol.* **21** 120
- [18] Buda K, Miton C M and Tokuriki N 2023 Pervasive epistasis exposes intramolecular networks in adaptive enzyme evolution *Nat. Commun.* **14** 8508
- [19] Fantini M, Lisi S, De Los Rios P, Cattaneo A and Pastore A 2020 Protein structural information and evolutionary landscape by *in vitro* evolution *Mol. Biol. Evol.* **37** 1179–92
- [20] Stiffler M A, Poelwijk F J, Brock K P, Stein R R, Riesselman A, Teyra J, Sidhu S S, Marks D S, Gauthier N P and Sander C 2020 Protein structure from experimental evolution *Cell Syst.* **10** 15–24
- [21] Erdoğan A N, Dasmeh P, Socha R D, Chen J Z, Life B E, Jun R, Kiritchkov L, Kehila D, Serohijos A W R and Tokuriki N 2024 Neutral drift upon threshold-like selection promotes variation in antibiotic resistance phenotype *Nat. Commun.* **15** 10813
- [22] Park Y, Metzger B P H and Thornton J W 2022 Epistatic drift causes gradual decay of predictability in protein evolution *Science* **376** 823–30
- [23] Rix G, Williams R L, Hu V J, Spinner A, Pisera A (O), Marks D S and Liu C C 2024 Continuous evolution of user-defined genes at 1 million times the genomic mutation rate *Science* **386** 9073
- [24] de la Paz J A, Nartey C M, Yuvaraj M and Morcos F 2020 Epistatic contributions promote the unification of incompatible models of neutral molecular evolution *Proc. Natl Acad. Sci.* **117** 5873–82
- [25] Bisardi M, Rodriguez-Rivas J, Zamponi F and Weigt M 2022 Modeling sequence-space exploration and emergence of epistatic signals in protein evolution *Mol. Biol. Evol.* **39** msab321
- [26] Alvarez S, Nartey C, Mercado N and Morcos F 2022 Novel sequence space explored by functional proteins generated through computational evolution-based design *Biophys. J.* **121** 45a
- [27] Alvarez S, Nartey C M, Mercado N, de la Paz J A, Huseinbegovic T and Morcos F 2024 *In vivo* functional phenotypes from a computational epistatic model of evolution *Proc. Natl Acad. Sci.* **121** e2308895121
- [28] Biswas A, Choudhuri I, Arnold E, Lyumkis D, Haldane A and Levy R M 2024 Kinetic coevolutionary models predict the temporal emergence of hiv-1 resistance mutations under drug selection pressure *Proc. Natl Acad. Sci.* **121** e2316662121
- [29] Bari L Di, Bisardi M, Cotogno S, Weigt M and Zamponi F 2024 Emergent time scales of epistasis in protein evolution *Proc. Natl Acad. Sci.* **121** e2406807121
- [30] Morcos F, Pagnani A, Lunt B, Bertolino A, Marks D S, Sander C, Zecchina R, Onuchic J N, Hwa T and Weigt M 2011 Direct-coupling analysis of residue coevolution captures native contacts across many protein families *Proc. Natl Acad. Sci.* **108** E1293–301
- [31] Ferguson A L, Mann J K, Omarjee S, Ndung'u T, Walker B D and Chakraborty A K 2013 Translating HIV sequences into quantitative fitness landscapes predicts viral vulnerabilities for rational immunogen design *Immunity* **38** 606–17
- [32] Figliuzzi M, Jacquier H, Schug A, Tenaillon O and Weigt M 2016 Coevolutionary landscape inference and the context-dependence of mutations in beta-lactamase TEM-1 *Mol. Biol. Evol.* **33** 268–80

- [33] Levy R M, Haldane A and Flynn W F 2017 Potts hamiltonian models of protein co-variation, free energy landscapes and evolutionary fitness *Curr. Opin. Struct. Biol.* **43** 55–62
- [34] Couce A, Caudwell L V, Feinauer C, Hindré T, Feugeas J-P, Weigt M, Lenski R E, Schneider D and Tenaillon O 2017 Mutator genomes decay, despite sustained fitness gains, in a long-term experiment with bacteria *Proc. Natl Acad. Sci.* **114** E9026–35
- [35] Vigué L and Tenaillon O 2023 Predicting the effect of mutations to investigate recent events of selection across 60,472 *Escherichia coli* strains *Proc. Natl Acad. Sci.* **120** e2304177120
- [36] Biswas A, Haldane A, Arnold E and Levy R M 2019 Epistasis and entrenchment of drug resistance in HIV-1 subtype B *eLife* **8** e50524
- [37] Lyons D M, Zou Z, Xu H and Zhang J 2020 Idiosyncratic epistasis creates universals in mutational effects and evolutionary trajectories *Nat. Ecol. Evol.* **4** 1685–93
- [38] Vigué L, Croce G, Petitjean M, Ruppé E, Tenaillon O and Weigt M 2022 Deciphering polymorphism in 61,157 *Escherichia coli* genomes via epistatic sequence landscapes *Nat. Commun.* **13** 4030
- [39] Chen J Z, Bisardi M, Lee D, Cotogno S, Zamponi F, Weigt M and Tokuriki N 2024 Understanding epistatic networks in the B1 beta-lactamases through coevolutionary statistical modeling and deep mutational scanning *Nat. Commun.* **15** 8441
- [40] Kirkpatrick T R and Thirumalai D 1988 Comparison between dynamical theories and metastable states in regular and glassy mean-field spin models with underlying first-order-like phase transitions *Phys. Rev. A* **37** 4439
- [41] Franz S, Donati C, Parisi G and Glotzer S C 1999 On dynamical correlations in supercooled liquids *Phil. Mag. B* **79** 1827–31
- [42] Franz S and Parisi G 2000 On non-linear susceptibility in supercooled liquids *J. Phys.: Condens. Matter* **12** 6335
- [43] Bouchaud J-P and Biroli G 2005 Nonlinear susceptibility in glassy systems: a probe for cooperative dynamical length scales *Phys. Rev. B* **72** 064204
- [44] Berthier L, Biroli G, Bouchaud J-P, Cipelletti L, Masri D E, L'Hôte D, Ladieu F and Pierno M 2005 Direct experimental evidence of a growing length scale accompanying the glass transition *Science* **310** 1797–800
- [45] Berthier L and Jack R L 2007 Structure and dynamics of glass formers: predictability at large length scales *Phys. Rev. E* **76** 041509
- [46] Franz S, Parisi G, Ricci-Tersenghi F and Rizzo T 2011 Field theory of fluctuations in glasses *Eur. Phys. J. E* **34** 1–17
- [47] Berthier L, Biroli G, Bouchaud J-P, Cipelletti L and van Saarloos W 2011 *Dynamical Heterogeneities in Glasses, Colloids and Granular Media* vol 150 (OUP Oxford)
- [48] Franz S, Jacquin H, Parisi G, Urbani P and Zamponi F 2013 Static replica approach to critical correlations in glassy systems *J. Chem. Phys.* **138** 12A540
- [49] Seoane B and Zamponi F 2018 Spin-glass-like aging in colloidal and granular glasses *Soft Matter* **14** 5222–34
- [50] Folea G, Biroli G, Charbonneau P, Hu Y and Zamponi F 2022 Equilibrium fluctuations in mean-field disordered models *Phys. Rev. E* **106** 024605
- [51] Widmer-Cooper A, Perry H, Harrowell P and Reichman D R 2008 Irreversible reorganization in a supercooled liquid originates from localized soft modes *Nat. Phys.* **4** 711–5
- [52] Schoenholz S S, Cubuk E D, Sussman D M, Kaxiras E and Liu A J 2016 A structural approach to relaxation in glassy liquids *Nat. Phys.* **12** 469–71
- [53] Bapst V *et al* 2020 Unveiling the predictive power of static structure in glassy systems *Nat. Phys.* **16** 448–54
- [54] Jung G, Biroli G and Berthier L 2024 Dynamic heterogeneity at the experimental glass transition predicted by transferable machine learning *Phys. Rev. B* **109** 064205
- [55] Jung G *et al* 2025 Roadmap on machine learning glassy dynamics *Nat. Rev. Phys.* **7** 1–14
- [56] Otwinowski J 2018 Biophysical inference of epistasis and the effects of mutations on protein stability and function *Mol. Biol. Evol.* **35** 2345–54
- [57] Otwinowski J, McCandlish D M and Plotkin J B 2018 Inferring the shape of global epistasis *Proc. Natl Acad. Sci.* **115** E7550–8
- [58] Reddy G and Desai M M 2021 Global epistasis emerges from a generic model of a complex trait *eLife* **10** e64740
- [59] Schulte A O, Alqatari S, Rossi S and Zamponi F 2025 Functional bottlenecks can emerge from non-epistatic underlying traits *bioRxiv Preprint* (<https://doi.org/10.1101/2025.05.20.655048>) (posted online 20 May 2025, accessed 20 May 2025)
- [60] Sailer Z R and Harms M J 2017 Detecting high-order epistasis in nonlinear genotype-phenotype maps *Genetics* **205** 1079–88
- [61] Sailer Z R and Harms M J 2017 High-order epistasis shapes evolutionary trajectories *PLoS Comput. Biol.* **13** e1005541
- [62] Domingo J, Diss G and Lehner B 2018 Pairwise and higher-order genetic interactions during the evolution of a tRNA *Nature* **558** 117–21
- [63] Poelwijk F J, Socolich M and Ranganathan R 2019 Learning the pattern of epistasis linking genotype and phenotype in a protein *Nat. Commun.* **10** 4213
- [64] Ballal A, Laurendon C, Salmon M, Vardakou M, Cheema J, Defernez M, O'Maille P E and Morozov A V 2020 Sparse epistatic patterns in the evolution of terpene synthases *Mol. Biol. Evol.* **37** 1907–24
- [65] Phillips A M, Lawrence K R, Moulana A, Dupic T, Chang J, Johnson M S, Cvijovic I, Mora T, Walczak A M and Desai M M 2021 Binding affinity landscapes constrain the evolution of broadly neutralizing anti-influenza antibodies *eLife* **10** e71393
- [66] Miton C M, Buda K and Tokuriki N 2021 Epistasis and intramolecular networks in protein evolution *Curr. Opin. Struct. Biol.* **69** 160–8
- [67] Lunzer M, Golding G B and Dean A M 2010 Pervasive cryptic epistasis in molecular evolution *PLoS Genet.* **6** e1001162
- [68] Miton C M and Tokuriki N 2016 How mutational epistasis impairs predictability in protein evolution and design *Protein Sci.* **25** 1260–72
- [69] Rivoire O, Reynolds K A and Ranganathan R 2016 Evolution-based functional decomposition of proteins *PLoS Comput. Biol.* **12** e1004817
- [70] Starr T N, Flynn J M, Mishra P, Bolon D N A and Thornton J W 2018 Pervasive contingency and entrenchment in a billion years of HSp90 evolution *Proc. Natl Acad. Sci.* **115** 4453–8
- [71] Bakerlee C W, Nguyen Ba A N, Shulgina Y, Rojas Echenique J I and Desai M M 2022 Idiosyncratic epistasis leads to global fitness–correlated trends *Science* **376** 630–5
- [72] Papkou A, Garcia-Pastor L, Escudero J A and Wagner A 2023 A rugged yet easily navigable fitness landscape *Science* **382** 901
- [73] Somermeyer L G, Fleiss A, Mishin A S, Bozhanova N G, Igolkina A A, Meiler J, Pujol M E A, Putintseva E V, Sarkisyan K S and Kondrashov F A 2022 Heterogeneity of the GFP fitness landscape and data-driven protein design *eLife* **11** e75842
- [74] Schulz S, Tan T J C, Wu N C and Wang S 2025 Epistatic hotspots organize antibody fitness landscape and boost evolvability *Proc. Natl Acad. Sci.* **122** e2413884122

- [75] Olson C A, Wu N C and Sun R 2014 A comprehensive biophysical description of pairwise epistasis throughout an entire protein domain *Curr. Biol.* **24** 2643–51
- [76] Romanowicz K J, Resnick C, Hinton S R and Plesa C 2025 Exploring antibiotic resistance in diverse homologs of the dihydrofolate reductase protein family through broad mutational scanning *bioRxiv Preprint* <https://doi.org/10.1101/2025.01.23.634126> (posted online 24 January 2025, accessed February 2025)
- [77] Figliuzzi M, Barrat-Charlaix P and Weigt M 2018 How pairwise coevolutionary models capture the collective residue variability in proteins? *Mol. Biol. Evol.* **35** 1018–27
- [78] Muntoni A P, Pagnani A, Weigt M and Zamponi F 2021 adabmDCA: adaptive Boltzmann machine learning for biological sequences *BMC Bioinform.* **22** 1–19
- [79] Russ W P *et al* 2020 An evolution-based model for designing chorisimate mutase enzymes *Science* **369** 440–5
- [80] Ashkenazy H, Penn O, Doron-Faigenboim A, Cohen O, Cannarozzi G, Zomer O and Pupko T 2012 FastML: a web server for probabilistic reconstruction of ancestral sequences *Nucl. Acids Res.* **40** W580–4
- [81] Gascuel O and Steel M 2010 Inferring ancestral sequences in taxon-rich phylogenies *Math. Biosci.* **227** 125–35
- [82] Evans W, Kenyon C, Peres Y and Schulman L J 2000 Broadcasting on trees and the ising model *Ann. Appl. Probab.* **10** 410–33
- [83] Felsenstein J 2003 *Inferring Phylogenies* (Oxford University Press)
- [84] De Leonardis M, Pagnani A and Barrat-Charlaix P 2025 Reconstruction of ancestral protein sequences using autoregressive generative models *Mol. Biol. Evol.* **42** msaf070