

# Skoltech

Skolkovo Institute of Science and Technology

## GENOMIC PATTERNS OF EPISTASIS AT MACRO- AND MICROEVOLUTIONARY SCALES

*Doctoral thesis*

*by*

Anastasia Stolyarova

Doctoral Program in Life Sciences

Supervisor

Professor Georgii Bazykin

Moscow 2021

I hereby declare that the work presented in this thesis was carried out by myself at Skolkovo Institute of Science and Technology, Moscow, except where due acknowledgment is made, and has not been submitted for any other degree.

Candidate (Anastasia Stolyarova)

Supervisor (Prof. Georgii Bazykin)

# Abstract

Many lines of evidence indicate that amino acid sites of protein-coding genes are involved in tight networks of epistatic interactions. In the course of evolution, amino acid preferences may change due to the coevolution in epistatically interacting sites or due to global landscape changes that are external to the genome. In this thesis, we use comparative genomics methods to study how the shape of the fitness landscape impacts the patterns of genetic differences observed on various evolutionary scales.

To address the contribution of epistatic selection to patterns of within-population genetic diversity, we study population genomics data of the fungus *Schizophyllum commune*, the most polymorphic species known. Throughout its genome, we observe the excess of short-range linkage disequilibrium between nonsynonymous polymorphisms, caused by attraction of rare alleles. This effect is especially pronounced for pairs of sites that are located within the same gene. Together with elevated LD between pairs of sites that encode physically interacting amino acids, and a substantial correlation between LDs between shared pairs of nonsynonymous polymorphisms in two *S. commune* populations, these patterns indicate that selection in *S. commune* involves positive epistasis due to coevolution between nonsynonymous alleles.

To find evidence of epistatic selection in the divergence of closely related species, we examine the phylogenies of the Baikal Lake amphipods and of primates, which contain very short internal edges. We detect six salient bursts of evolution of individual proteins during such short time periods, each involving between 6 and 38 amino acid substitutions and limited to phylogenetic edges  $< 0.001 dS$ . These bursts are extremely unlikely to have occurred neutrally and are apparently caused by correlated positive selection.

Using analytical modeling and simulations, we show that we can distinguish between different factors of fitness landscapes changes by looking at the dynamics of the fitness of the allele currently occupying the site: it tends to increase with the time since its origin due to coevolution between epistatically interacting sites (“entrenchment”) and

to decrease due to random environmental fluctuations (“senescence”). By comparing the genomes of diverged species (vertebrates and insects), we show that the amino acids originating at negatively selected sites experience strong entrenchment. By contrast, the amino acids originating at positively selected sites experience senescence.

These findings are indicative of the complex structure of selective constraints shaping the patterns of genetic variation within and between species.

# Acknowledgments

First of all, I want to thank my co-authors for their contribution to the research projects included in this thesis: Vasily Ptushenko for the theoretical framework on the decline of the current allele fitness driven by random changes of the fitness landscape (Chapter 5); Anna Fedotova for *S. commune* DNA library construction and sequencing (Chapter 3); Elena Nabieva, Alexey Neverov, Anfisa Popova and Alexander Favorov for discussing the projects and helping with methods and writing (Chapter 5). I would say a special thank you to Tatiana Neretina and Elena Zvyagina for maintaining the collection of *S. commune* samples and extracting genomic DNA, and to Timothy James, Anna Baikalo, and all members of the Bazykin-Kondrashov group who participated in *S. commune* sampling both in Russia and the US. I believe keeping this collection is particularly important and will prove useful in future studies. I would also like to thank Shamil Sunyaev and the members of his lab for useful comments on the draft of the *S. commune* paper. I am grateful to the reviewers for providing comments and suggestions that helped to improve the thesis text.

I am deeply grateful to my supervisors Georgii Bazykin and Alexey Kondrashov for their guidance and support during my Ph.D. and for sharing their expertise and perspective with me. Georgii encouraged me to try new things and provided me freedom of choosing research topics and methods. I highly appreciate our regular discussions with Alexey, which I find invaluable and insightful.

I want to thank my fellow students Aleksandra Bezmenova, Ksenia Safina, and Marina Kalinina, with whom we did this journey together. All the colleagues from the Bazykin-Kondrashov group and other labs from our bioinformatics community – with many of them also being my close friends – for the unique atmosphere of support and the opportunity to share and discuss ideas with you.

I want to thank Fyodor Kondrashov and the members of his lab for the opportunity to implement my idea of yeast fitness experiment as part of my internship at IST Austria and support during the complicated pandemic times, and Skoltech for making this internship possible.

The completion of my thesis wouldn't have been possible without the love and help of my family – especially my mom Elena Stolyarova, who has always believed in me and supported my choice of studying biology. Last but not least, I want to thank Ruslan Soldatov for always being here, for understanding and encouragement, and for having enough patience to listen about epistasis on a daily basis.

# Publications

## Papers

- AV Stolyarova, E Nabieva, VV Ptushenko, AV Favorov, AV Popova, AD Neverov, GA Bazykin. **Senescence and entrenchment in evolution of amino acid sites.**  
*Nature Communications* 2020; 14;11(1):4603.  
doi: 10.1038/s41467-020-18366-z
- AV Stolyarova, GA Bazykin, TV Neretina, AS Kondrashov. **Bursts of amino acid replacements in protein evolution.**  
*Royal Society Open Science* 2019; 6(3):181095. doi: 10.1098/rsos.181095

## Preprints

- AV Stolyarova, TV Neretina, EA Zvyagina, AV Fedotova, AS Kondrashov, GA Bazykin. **Complex fitness landscape shapes variation in a hyperpolymorphic species.**  
*bioRxiv* 2021. doi.org/10.1101/2021.10.10.463656

## Conference presentations

- Within-gene epistatic selection in genetically diverse populations.  
*MCCMB, Moscow 2021* (Oral presentation)
- Within-gene epistatic selection in genetically diverse populations.  
*EMBL: Predicting evolution 2021* (Poster presentation)
- Within-gene epistatic selection shapes polymorphism in natural populations of the world's most variable eukaryotic species.  
*EMBL: Molecular Mechanisms of Evolution & Ecology 2020* (Oral presentation)
- Prevalent epistatic interactions between amino acid sites in *S. commune*.  
*MCCMB, Moscow 2019* (Oral presentation)
- Prevalent epistatic interactions between amino acid sites in *S. commune*.  
*SMBE, Manchester UK, 2019* (Poster presentation)
- Bursts of nonsynonymous replacements in protein evolution.  
*SMBE, Yokohama 2018* (Poster presentation)
- Causes of single position fitness landscape changes.  
*SMBE, Austin 2017* (Poster presentation)

# Contents

|   |           |
|---|-----------|
| <b>Abstract</b>                             | <b>2</b>  |
| <b>Acknowledgments</b>                      | <b>4</b>  |
| <b>Publications</b>                         | <b>6</b>  |
| <b>Contents</b>                             | <b>8</b>  |
| <b>List of abbreviations</b>                | <b>11</b> |
| <b>Chapter 1: Introduction</b>              | <b>12</b> |
| <b>Chapter 2: Literature review</b>         | <b>14</b> |
| Fitness landscapes                          | 14        |
| Epistasis                                   | 16        |
| Models of fitness landscapes                | 20        |
| Evolution on fitness landscapes             | 26        |
| Populations on fitness landscapes           | 26        |
| Valley crossing                             | 28        |
| Dynamics of adaptation                      | 29        |
| Mutational robustness                       | 31        |
| Predictability of evolution                 | 32        |
| Epistasis and recombination                 | 37        |
| Dynamic fitness landscapes                  | 43        |
| Environmental fluctuations                  | 43        |
| Frequency-dependent selection               | 45        |
| Epistatic changes of allele's fitness       | 46        |
| Empirical fitness landscapes                | 49        |
| Landscapes of homologous sequences          | 49        |
| The complexity of the empirical landscapes  | 50        |
| Empirical inference of historical evolution | 53        |

|  |            |
|--|------------|
| Phylogenetic evidence of epistasis   | 56         |
| Compensated pathogenic deviations  | 56         |
| Patterns of divergent and convergent evolution   | 56         |
| Correlated evolution of interacting sites  | 58         |
| Phylogenetic clustering of interacting sites   | 59         |
| Epistasis in within-population variation   | 61         |
| Statistical epistasis  | 61         |
| Epistasis-driven linkage disequilibrium  | 61         |
| <b>Chapter 3: Complex fitness landscape shapes variation in a hyperpolymorphic species</b> | <b>64</b>  |
| Introduction   | 65         |
| Materials and methods  | 67         |
| <i>S. commune</i> sampling, sequencing and assembly  | 67         |
| Data on <i>H. sapiens</i> and <i>D. melanogaster</i> populations                           | 70         |
| Estimation of LD   | 71         |
| Haploblocks annotation   | 71         |
| Estimation of LD between physically interacting amino acid sites                           | 71         |
| Simulations of epistasis   | 72         |
| Results  | 74         |
| Epistatic selection is more efficient in genetically diverse populations                   | 74         |
| Elevated LD between nonsynonymous polymorphisms  | 76         |
| Physically interacting amino acid sites are under stronger LD                              | 86         |
| Excess of LD <sub>nonsyn</sub> is more pronounced in distinct regions of high LD           | 88         |
| Excess of LD <sub>nonsyn</sub> requires stable polymorphism                                | 92         |
| Correlated LDs between shared SNPs in two populations                                      | 95         |
| Discussion   | 99         |
| <b>Chapter 4: Correlated positive selection leads to bursts of amino acid replacements</b> | <b>101</b> |
| Introduction   | 102        |
| Materials and methods  | 104        |

|  |            |
|--|------------|
| Phylogenies of closely related species   | 104        |
| Inference of bursts of nonsynonymous substitutions   | 104        |
| Filtering of candidate bursts  | 105        |
| Results  | 107        |
| Discussion   | 114        |
| <b>Chapter 5: Changes of single-position fitness landscapes affect evolution of amino acid sites</b> | <b>119</b> |
| Introduction   | 120        |
| Materials and methods  | 122        |
| Multiple alignments of protein-coding sequences  | 122        |
| Simulations of amino acid evolution on dynamic landscapes  | 124        |
| Substitution subtrees  | 125        |
| Inference of senescence or entrenchment for groups of alleles  | 126        |
| Summary statistics   | 128        |
| ABC validation   | 128        |
| Results  | 131        |
| Environmental fluctuations decrease the fitness of the current allele                                | 131        |
| Senescence and entrenchment result in opposite substitution patterns                                 | 133        |
| Senescence and entrenchment at single-allele resolution  | 137        |
| Heterogeneity of alleles leads to an artifactual signal of entrenchment                              | 139        |
| Inferring senescence and entrenchment from phylogenetic distribution of substitutions                | 142        |
| Positively selected sites show strong senescence   | 142        |
| Discussion   | 146        |
| <b>Chapter 6: Conclusions</b>  | <b>148</b> |
| <b>References</b>  | <b>151</b> |
| <b>Appendix A</b>  | <b>170</b> |

## List of abbreviations

ABC — approximate Bayesian computations

AOD — associative overdominance

BGS — background selection

CPD — compensated pathogenic deviations

DCA — direct-coupling analysis

DMS — deep mutational scanning

FDR — false discovery rate

GWAS — genome-wide association studies

HoC — house-of-cards model

HRI — Hill-Robertson interference

LD — linkage disequilibrium

LoF — loss-of-function mutation

LRLD — long-range linkage disequilibrium

LTEE — long-term evolution experiment

NFDS — negative frequency-dependent selection

OR — odds-ratio

SNP — single-nucleotide polymorphism

SPFL — single-position fitness landscape

WT — wild-type

## Chapter 1: Introduction

The term “epistasis” was initially introduced by W. Bateson in 1909 to define the phenomenon of under-representation of some phenotypic classes in dihybrid crosses, which he explained by some mutations masking the impact of other ones. R. A. Fisher later used the related term “epistacy” to denote any deviation from the independence of effects of single mutations in different loci — the meaning generally used by this day (Ronald Aylmer Fisher 1918). The term “epistasis” comprehends genetic interactions of various nature: it can describe general laws of selection, features of adaptive paths, the relation between genotype and phenotype, functional interactions between specific genes, or susceptibility to complex diseases. From an evolutionary point of view, epistasis is not a phenomenon but a general characteristic of the fitness landscape underlying the process of adaptation.

Given the complexity and hierarchical structure of biological systems, the presence of ubiquitous epistasis by itself is not surprising. However, the question is whether the complex structure of fitness landscapes actually influences the evolution of natural populations — and what we can conclude on the evolutionary mechanisms underlying the observed patterns of evolutionary change. Describing evolution on epistatic fitness landscapes may help to explain various aspects observed in the evolution of natural populations, such as hybrid incompatibilities, evolutionary advantage of sex, predictability and repeatability of adaptive paths, compensation of pathogenic variants, historical contingency in evolution, *etc.*

Currently available genomic data make it possible to study the action of natural selection without knowing the fitness or any other phenotype directly by analyzing the accumulation of genetic differences on the microevolutionary scale (*i.e.* between individuals comprising the same population) or on the macroevolutionary scale (*i.e.* between individuals representing diverged species).

Since epistasis implies any kind of dependency between effects of mutations in different loci, the potential evolutionary consequences of it may be various. Population genomics

data provide a snapshot of the current population state, allowing the assessment and interpretation of the deviations of the joint distribution of alleles within a population. By comparing the genomes of different species and reconstructing their phylogenies, we can trace the changes of natural selection acting on the substituted alleles and infer the causes of these changes.

The goal of this thesis is to infer how epistasis affects the evolution of protein-coding sequences in natural populations on different evolutionary scales.

The main objectives of this thesis are:

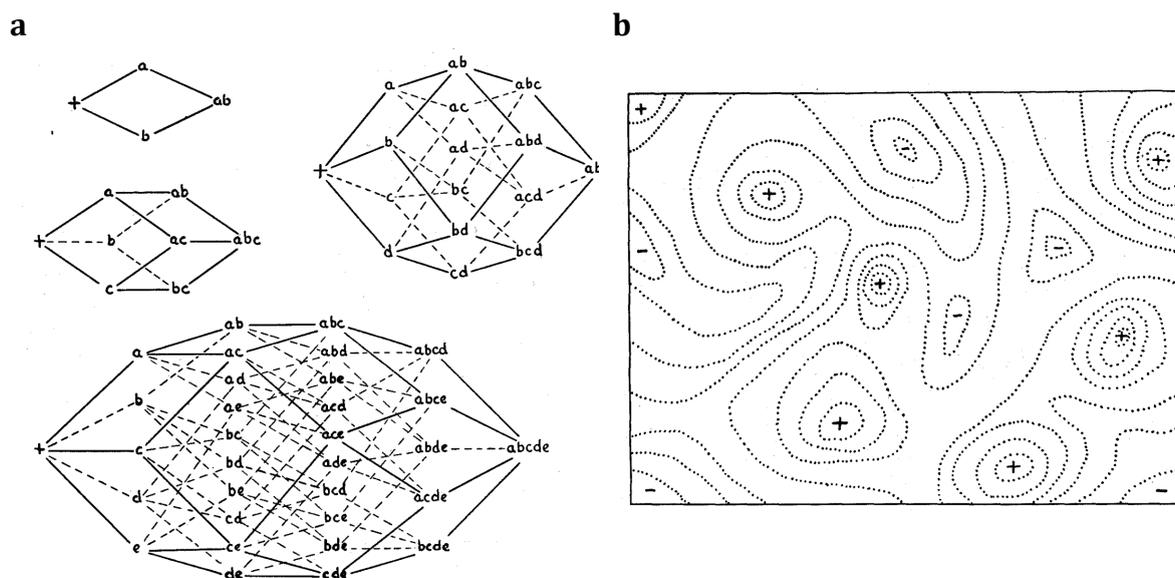
1. to assess the evidence of epistasis shaping the within-population variation based on patterns of linkage disequilibrium between natural polymorphisms;
2. to detect bursts of adaptive evolution between recently diverged species, presumably caused by correlated positive selection between multiple genomic sites;
3. to estimate how selection acting on the allele currently occupying a genomic site changes in the course of species divergence.

To achieve these objectives, we use comparative genomics methods to assess the distribution of polymorphisms within populations or the patterns of substitutions between diverged species. To examine patterns of standing genetic variation, we consider population genomics data on populations of varying levels of genetic diversity: from the widely studied, but less polymorphic populations of *Homo sapiens* and *Drosophila melanogaster* to the most genetically variable eukaryotic species known *Schizophyllum commune*. As for studying the adaptive evolution of closely related species, we make use of dense phylogenies of Lake Baikal amphipods and of primates. To infer changes of the current allele fitness on the macroevolutionary scale, we analyze multiple alignments of orthologous genes of vertebrates and insects and mitochondrial genes of Metazoa. In the first and second tasks, we focus on the short-range interactions (*e.g.* within genes or between sites located in neighboring genes), while in the third task we study genome-wide patterns of fitness landscape changes.

## Chapter 2: Literature review

### Fitness landscapes

**Fitness landscape**, or adaptive landscape, is a key concept of evolutionary biology. It is a function mapping the genotype space to **fitness**, *i.e.* the measure of the evolutionary success, of an individual carrying the corresponding genotype. The genotype space represents all possible combinations of alleles in a number of genetic loci and has the size of  $K^L$ , where  $K$  is the alphabet size (the number of permissible alleles) and  $L$  is the number of loci.



**Figure 2.1. Representation of fitness landscapes** (Wright 1932). **(a)** The full space of possible genotypes of  $L$  loci (with  $L$  from 2 to 5) with  $K=2$  possible alleles in each locus are  $L$ -dimensional hypercubes. **(b)** Geometric representation of the fitness landscape on the two-dimensional genotype space. Dotted lines connect genotypes with equal fitness, “+” mark adaptive peaks, and “-” mark regions of low fitness.

The full landscape of  $L$  biallelic sites is a hypercube in a  $L$ -dimensional space of genotypes (Figure 2.1a). The high dimensionality of the landscape makes it hard to conceive its structure and to predict how its shape conducts the evolution of the

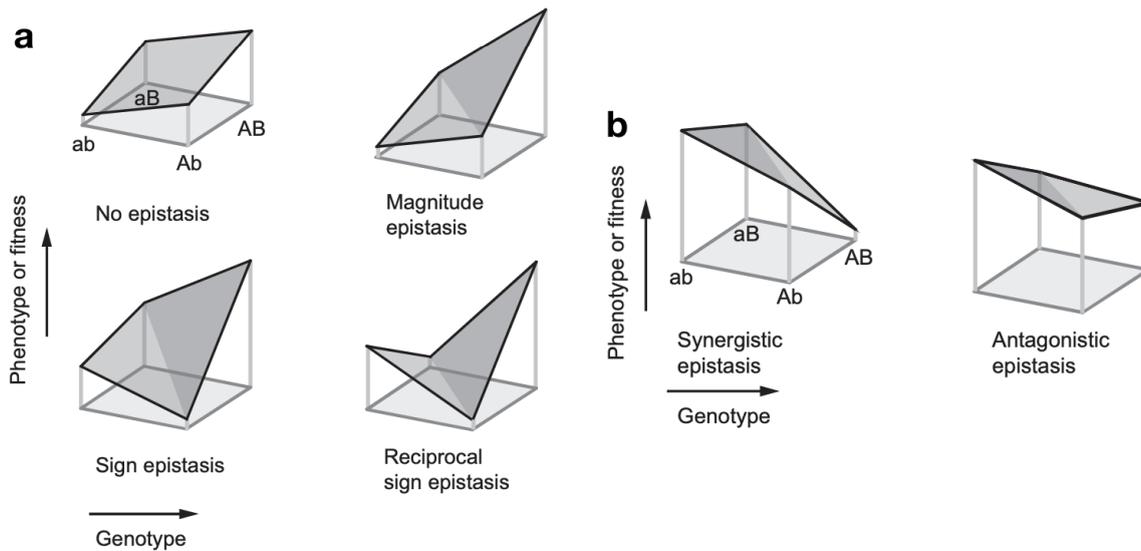
genomic sequences. In empirical fitness landscape studies, while using experimental methods to directly measure fitness of the genotypes, the dimensionality is usually reduced by constraining the set of examined genotypes. Indirect studies of fitness landscapes use comparative genomic methods to infer general features of the landscapes or to describe the low-dimensional projections of the full landscape (a classic example of a fitness landscape on the two-dimensional genotype space by S. Wright is shown in Figure 2.1b). In order to study evolutionary dynamics of a population on the fitness landscape, S. Wright also described it in terms of relation between population mean fitness and genotype frequencies. Similarly, the landscape can be defined on the space of a quantitative phenotype, allowing to estimate selection acting on the distribution of the phenotype values (Lande 1976). In this thesis, we'll define the landscape only as the function relating genotypes to fitness (Provine 1989).

## Epistasis

If fitness effects of mutations in different genomic loci are independent, the effect of a combination of mutations can be calculated as a sum of the effects of individual mutations (or, in the multiplicative model, log fitness is assumed to be additive). In this case, the fitness landscape is linear and has a single adaptive peak (Figure 2.2a, upper left). Any deviation from this additivity means non-independence of mutation effects, or **epistasis** between loci. Under epistasis, the effect of a mutation depends on the genetic background — the allelic state of interacting sites. Therefore, the fitness of a combination of mutations can't be predicted from the additive compounds without knowing the epistatic coefficients between the mutations.

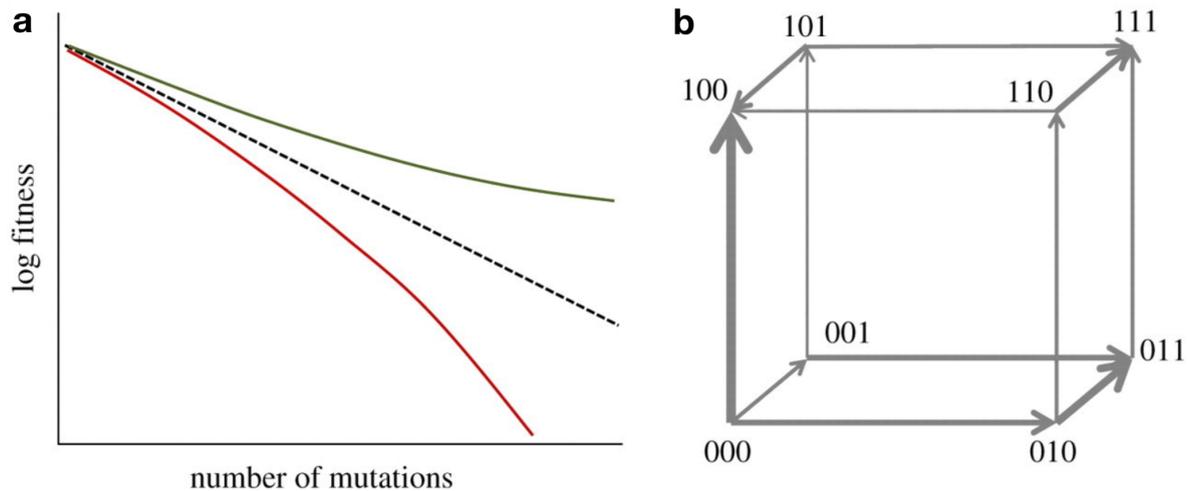
Since epistasis is any form of non-additivity of fitness effects of single mutations, it can be classified based on various criteria:

- Based on whether the sign of the fitness effect of a mutation depends on the presence of other mutations, epistasis can be **magnitude (monotonic)** and **sign**; in the case of monotonic epistasis, only the strength of a mutation can change but not its sign. Under reciprocal sign epistasis between two mutations, each of them can be both deleterious or beneficial depending on which allele occupies another site (Figure 2.2a) (Weinreich, Watson, and Chao 2005).
- Considering the direction of the shift of fitness value of the combination of mutation as compared to the additive expectation, epistasis can be **negative** (if the fitness of a genotype carrying multiple mutations is lower than expected) and **positive** (if it's higher than expected). If all mutations have the same sign, it can be redefined as **synergistic** (negative epistasis between deleterious mutations / positive between beneficial mutations) and **antagonistic** (positive epistasis between deleterious mutations / negative between beneficial) (Figure 2.2b). Positive/negative epistasis is polarized in regard to alleles occupying the interacting sites: for example, positive epistasis between mutations  $a \rightarrow A$  and  $b \rightarrow B$  corresponds to negative epistasis between mutations  $A \rightarrow a$  and  $b \rightarrow B$  and to positive epistasis between  $A \rightarrow a$  and  $B \rightarrow b$ .



**Figure 2.2. Classification of epistasis between mutations in two loci** (Kogenaru, de Vos, and Tans 2009). **(a)** Without epistasis, both mutations  $a \rightarrow A$  and  $b \rightarrow B$  are beneficial. Magnitude epistasis changes the fitness of the double mutant, but both mutations keep their sign whether a mutation in another locus is present or not. Under sign epistasis, the sign of the effect of the mutation  $a \rightarrow A$  depends on the state of the second locus: it is beneficial in the presence of allele  $b$  and deleterious in the presence of allele  $B$ . In the case of reciprocal sign epistasis, both mutations change their effect sign depending on the genetic background. **(b)** Under synergistic epistasis, the absolute value of the fitness effect of the double mutation is larger than expected without epistasis; under antagonistic epistasis, it is less than expected without epistasis. Since both mutations  $a \rightarrow A$  and  $b \rightarrow B$  are deleterious, synergistic interactions between them correspond to negative epistasis, and antagonistic — to positive epistasis.

- If the fitness of a genotype carrying a combination of mutations can be calculated as a simple, usually monotonic, nonlinear function of the sum of the additive effects of these mutations, epistasis is considered **unidimensional** (global). The assumption for the unidimensional epistasis is that fitness is defined by nonlinear scaling of some hidden feature (fitness potential), which is by itself additive. On the opposite, if fitness cannot be approximated with such function of the fitness potential, epistasis is **multidimensional** (F. A. Kondrashov and Kondrashov 2001a; de Visser, Cooper, and Elena 2011; Sailer and Harms 2017a) (Figure 2.3).



**Figure 2.3. The dimensionality of epistasis** (de Visser, Cooper, and Elena 2011). **(a)** Unidimensional epistasis between deleterious mutations, with the number of mutations in a genotype acting as the fitness potential function. Dashed line — no epistasis, green — narrowing (antagonistic) epistasis, red — widening (synergistic) epistasis. **(b)** Multidimensional epistasis between three biallelic sites. Arrows point towards the more fit genotypes, the thickness of the arrows indicates the size of fitness gain.

- In the case of unidimensional epistasis, it can be distinguished as **narrowing** and **widening** epistasis based on whether it reduces or increases the variance of some quantitative trait (Figure 2.3a) (Shnol and Kondrashov 1993).
- Under multidimensional epistasis, the non-additivity of fitness effects might be fully explained by **pairwise** interactions between considered sites, assuming that the epistatic coefficient of a pair of mutations doesn't depend on the allelic state of other loci. If this is not the case, the **high-order** epistatic terms should be considered while describing the fitness landscape. The number of epistatic coefficients of order  $n$  drastically grows with  $n$  as  $\binom{L}{n}$ , so the high-order interactions are hard to estimate and are often prone to overfitting (Hinkley et al. 2011; Zhou and McCandlish 2020).

In many biological systems, epistasis can be in large part reduced to unidimensional or low-order interactions; however, high-order epistasis is

detectable in multiple experimental landscapes and may affect the accessibility of evolutionary paths (G. Yang et al. 2019; Sailer and Harms 2017b; Weinreich et al. 2013; Sailer and Harms 2017a).

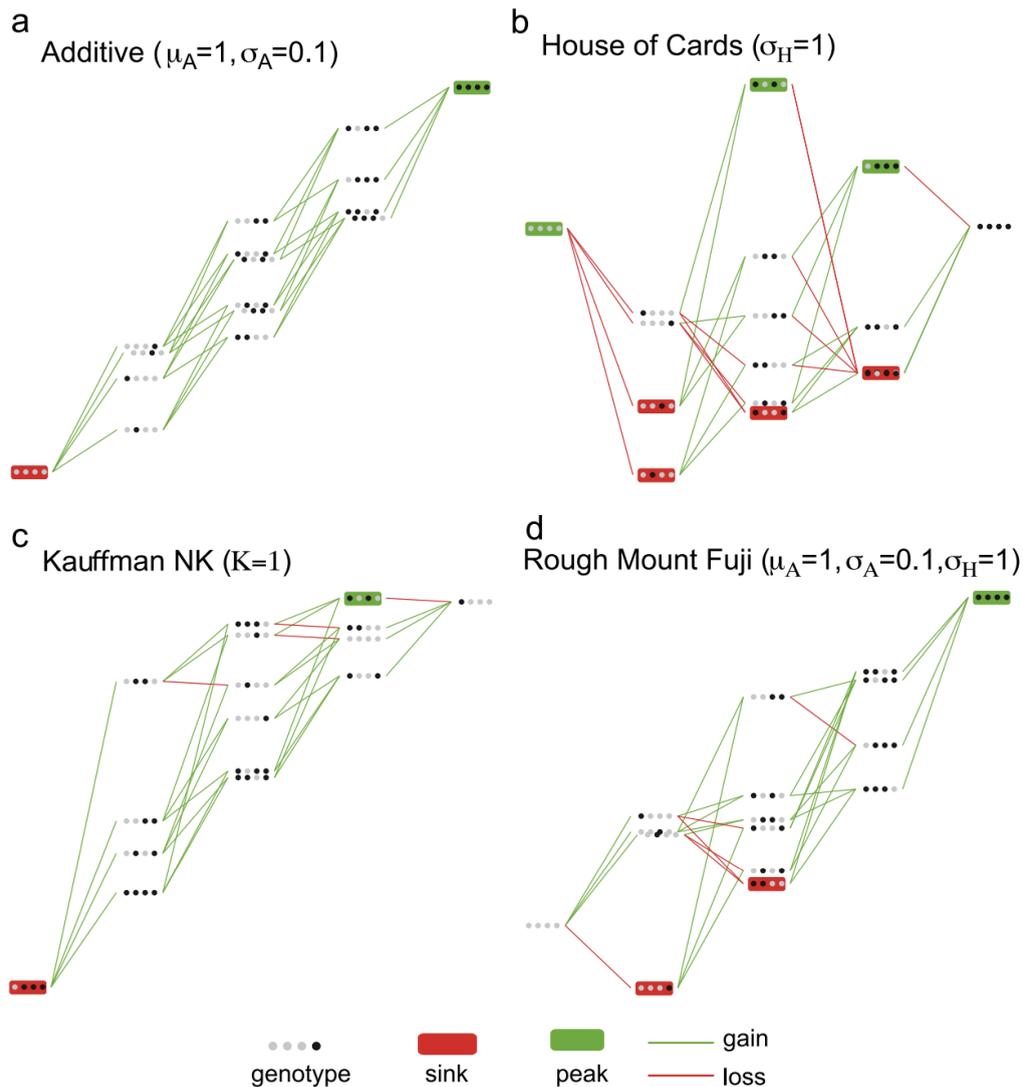
## Models of fitness landscapes

Since the full fitness landscape relates all possible genotypes to their fitness, it provides full information on epistasis acting between the considered loci (and *vice versa*, the landscape can be described by additive fitness effects of individual mutations and epistasis between them). For example, reciprocal sign epistasis is necessary to create rugged adaptive landscapes, *i.e.* the ones containing multiple local adaptive peaks (Poelwijk et al. 2011; Weinreich, Watson, and Chao 2005). Ruggedness is an important feature of a landscape, defining the accessibility of adaptive paths on the landscape and the predictability of evolution (Maynard Smith 1970; Weinreich et al. 2006; Poelwijk et al. 2007; Kvittek and Sherlock 2011; Ferretti et al. 2018).

Various modes and shapes of epistasis and selection can be implemented in different theoretical models of fitness landscapes. The landscape models are useful to interpret empirical datasets and to predict evolution under diverse regimes of selection (de Visser, Cooper, and Elena 2011; Kryazhimskiy, Tkacik, and Plotkin 2009; de Visser and Krug 2014; de Visser et al. 2018; Fragata et al. 2019; Bank et al. 2016). The low-dimensional landscape models imply fitness as a function of some features of the genotypes, which are assumed to be additive (e.g. FGM or power law landscapes), while multidimensional landscapes describe epistatic interactions between specific loci (e.g. HoC or NK landscapes) (Orr 2005). Here we will briefly describe some of them:

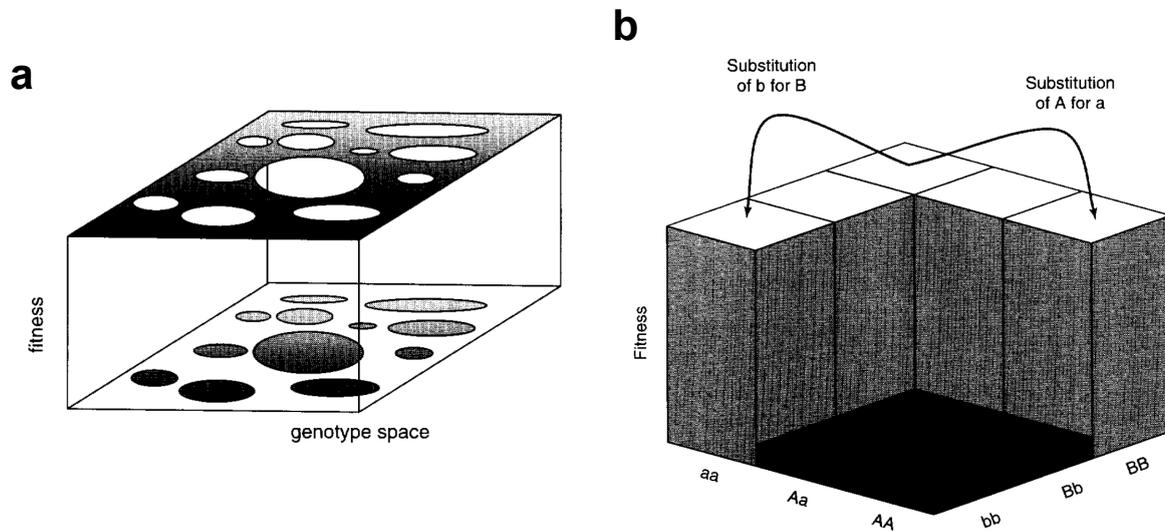
- In the house-of-cards (**HoC**) model, the fitness of each genotype is assumed to be randomly drawn from some distribution without regard to the fitness of the neighboring genotypes (Kingman 1978; S. Kauffman and Levin 1987). Therefore, it is impossible to predict the effect of a mutation on some genetic background knowing its effect on a different background. HoC model assumes extreme epistasis and produces rugged landscapes with multiple sharp peaks (Figure 2.4b) (Franke et al. 2011; Ferretti, Schmiegel, and Weinreich 2016). The rough Mount Fuji model (**RMF**) is the mixture of HoC and the additive landscape and is used to achieve the adjustable degree of epistasis (Figure 2.4d) (Aita et al. 2000; Neidhart, Szendro, and Krug 2014).

- **NK** model also allows adjusting the abundance of epistasis with a tunable parameter (Figure 2.4c) (S. A. Kauffman and Weinberger 1989; Sergey Gavrillets 1999). It is defined by the number of considered sites  $N$  and the parameter  $K$ , which implies the number of epistatic interactions per site. NK landscape with  $K = 0$  is additive; landscape with  $K = N - 1$  (the largest possible  $K$ ) is equal to the fully epistatic HoC landscape.



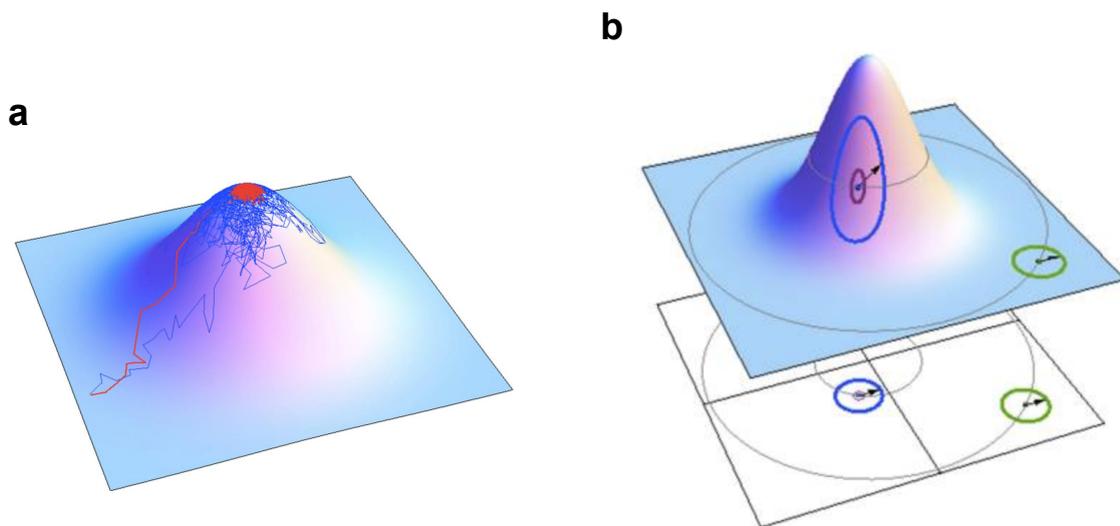
**Figure 2.4. Rugged models of fitness landscapes** (Ferretti, Schmiegelt, and Weinreich 2016). (a) Non-epistatic (additive) landscape ensures the convergence of all paths to the global maximum. (b) HoC model produces extremely rugged landscapes with multiple local peaks. (c, d) NK and RMF models with a tunable degree of epistasis can result in less or more rugged landscapes. Genotypes consisting of four biallelic loci are shown with dots; genotypes differing by only one mutation are connected.

- In contrast to the rugged landscape models described above, the **holey landscape** doesn't contain adaptive peaks and therefore doesn't imply adaptive evolution, or fitness gain. It represents the connected network of equally fit genotypes with holes, formed by lethal genotypes (Figure 2.5a) (S. Gavrillets 1997). Although the holey landscape is strongly epistatic (a mutation can be either neutral or lethal dependent on the initial genotype), the evolutionary process on such landscape is neutral. As a result of such evolution, diverging populations can get separated by the holes, creating hybrid incompatibilities without crossing adaptive valleys (T. Dobzhansky 1936; Orr 1995; A. S. Kondrashov, Sunyaev, and Kondrashov 2002). An example of a holey landscape is the Bateson-Dobzhansky-Muller (**BDM**) landscape model, that also assumes the existence of ridges connecting equally fit genotypes and is often considered in terms of incompatibilities and reproductive isolation (Figure 2.5b) (Theodosius Dobzhansky 1937; Bateson 1909; H. Muller 1942).



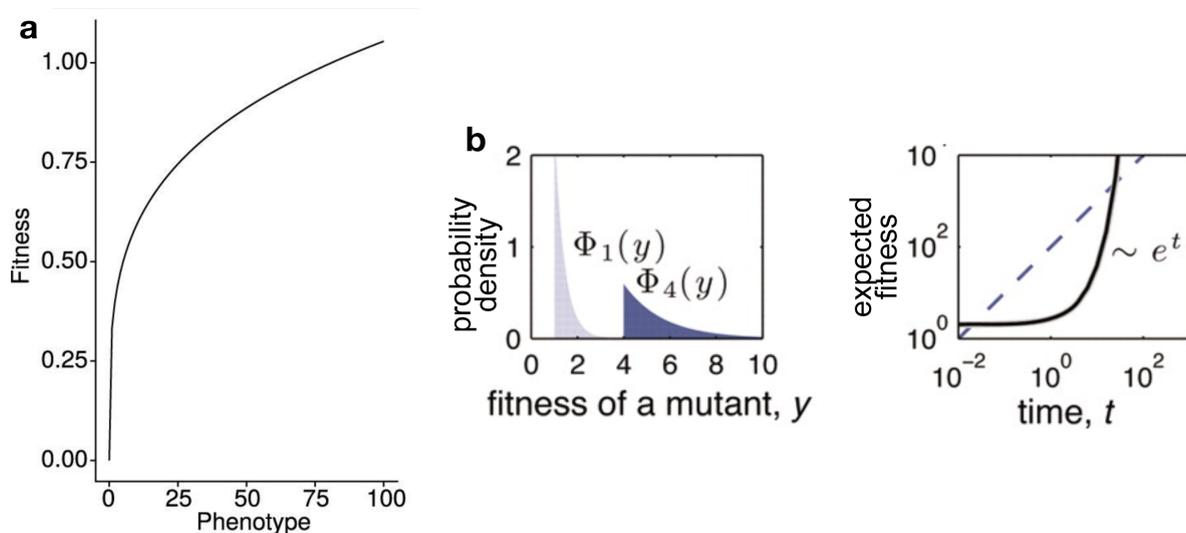
**Figure 2.5. Models of fitness landscapes containing ridges of equally fit genotypes** (S. Gavrillets 1997). (a) Holey landscape (b) BDM landscape on the space of two biallelic sites in a diploid population.

- Fisher's geometric model of the adaptive landscape (**FGM**) implies fitness as a smooth function on a space of several quantitative phenotypic traits (Ronald A. Fisher 1930; O. Tenaillon 2014). FGM assumes there is a single adaptive peak, so that selection promotes evolution of the considered traits towards the optimum (directional selection) and then restrains them to it (stabilizing selection) (Figure 2.6a). Under FGM, the effect of a mutation and both direction and strength of epistasis depend on the shape of the peak and the distance to it (Figure 2.6b, (G. Martin, Elena, and Lenormand 2007; O. Tenaillon 2014)). Due to the simplicity of the fitness function, FGM is widely used to describe adaptive evolution, wherein a population is forced to move from some point far from the global maximum towards it (Orr 1998; Bank et al. 2014; Gros, Le Nagard, and Tenaillon 2009; Harmand et al. 2017).



**Figure 2.6. Fisher's geometric model of the fitness landscape** (O. Tenaillon 2014). **(a)** Under selection, populations are driven towards the adaptive peak. Evolutionary paths of a population with large (red) and small (blue) effective population sizes are shown. **(b)** The fitness effect of a mutation (shown as arrows) depends on the location of the background genotype in regard to the adaptive peak. In the example, it's nearly neutral when far from the peak slope (green) or beneficial when close to the fitness maximum (blue).

- Power law landscapes** represent the class of landscapes without a fitness maximum. Such landscapes imply infinite adaptation — the evolving population can never approach the fitness peak so that new beneficial mutations are always getting fixed (Wiser, Ribeck, and Lenski 2013; Passagem-Santos, Zacarias, and Perfeito 2018). Such landscapes are unidimensional, and epistasis is defined by a non-linear monotonically increasing function of fitness over some additive trait (*i.e.* fitness potential). Usually, functions with a negative second derivative are used, conducting negative (antagonistic) epistasis between beneficial mutations (for example, power law function, Figure 2.7a). However, there is a class of “stairway to heaven” landscapes, which assume positive (synergistic) epistasis between beneficial mutations and therefore increasing rates of evolution (Figure 2.7b) (Kryazhimskiy, Tkacik, and Plotkin 2009).



**Figure 2.7. Models of fitness landscapes with infinite adaptation.** (a) Power law landscape with antagonistic epistasis between beneficial mutations (Passagem-Santos, Zacarias, and Perfeito 2018) (b) Stairway to heaven landscape with synergistic epistasis between beneficial mutations (Kryazhimskiy, Tkacik, and Plotkin 2009).  $\Phi_1(y)$  – distribution of fitness obtained by single mutations in the context of the initial genotype of low fitness (here, 1);  $\Phi_4(y)$  – of high fitness (here, 4).

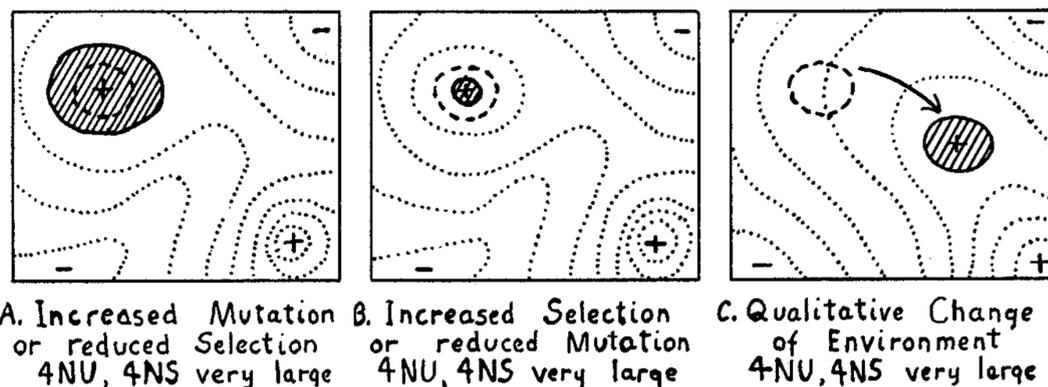
- There are multiple landscape models which are based on the assumptions on the features of the considered biological system (Fragata et al. 2019). In such models, epistasis is not defined explicitly but appears from these assumptions. For example, **biophysical** models define fitness of a protein or RNA sequence as a function of the corresponding structure, calculated expressed in terms of folding energy, stability, or other biophysical properties of the structure (Bershtein, Serohijos, and Shakhnovich 2017; Bertram and Masel 2020; Bershtein et al. 2006; Olson, Wu, and Sun 2014). Another class of models describes fitness on the level of **networks of interacting genes** (Reddy and Desai 2021; Friedlander et al. 2017; Yubero, Manrubia, and Aguirre 2017). Such models can be used to reconstruct interactions between genes within metabolic networks or interactions between amino acids within a single protein or between proteins.

## Evolution on fitness landscapes

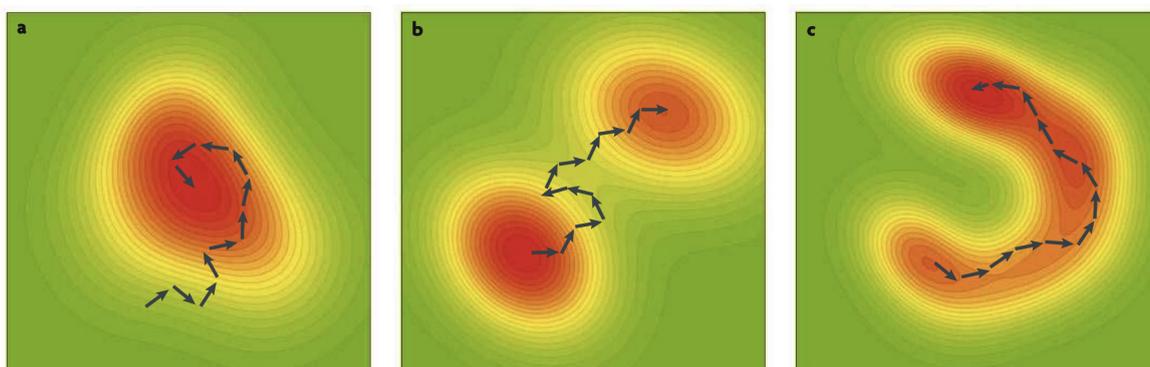
### Populations on fitness landscapes

Although a full fitness landscape relates each possible variant of the considered genomic sequence to its fitness, only a subset of these variants is actually present in a real evolving population at any time point. Since any population is not homogeneous, *i.e.* individuals within it differ to some extent, it occupies a certain area of the adaptive landscape — and most of the genotypes falling outside this area are never seen in the evolution of this population (Figure 2.8). In the course of microevolution, the density distribution of a population over the occupied region may vary as the frequency of present alleles changes due to natural selection. Under selection only, the frequency of highly fit genotypes will increase, and the frequency of low fit genotypes will decrease. Under the simple models of selection, this results in the increase of the mean population fitness and the population climbing onto a fitness peak (Figure 2.9a) (Ronald A. Fisher 1930; S. Kauffman and Levin 1987; Sergey Gavrilets 2004; Woodcock and Higgs 1996). However, there are selection models (*e.g.* multiple forms of balancing selection) preventing the population from converging to the fitness peak (R. A. Fisher 1941; Deborah Charlesworth 2006).

Also, new genotype variants occur due to mutation, introducing new alleles. In sexual populations, recombination also shuffles the present genotypes, composing new combinations of alleles at different recombining loci. In mutation-selection equilibrium, the size of the fitness landscape region occupied by a population depends on the strength of these forces: the higher is the mutation rate or the weaker is the selection, the larger will be the polymorphism level within the population (Figure 2.8) (Wright 1932; James Franklin Crow and Kimura 1970).



**Figure 2.8. Populations on fitness landscapes** (Wright 1932). (a) Under mutation-selection equilibrium, a highly polymorphic population occupies a large area of the adaptive peak. (b) A less polymorphic population is limited to the small area on the top of the peak. (c) After a sudden change of the landscape, selection drives adaptive evolution, pushing the population to the novel adaptive peak. Dashed regions show the subspace of the genotypes present in the population.



**Figure 2.9. Evolutionary paths on epistatic landscapes** (P. C. Phillips 2008). (a) In the course of adaptation on the FGM landscape, the evolutionary path consists of uphill steps of increasing fitness. (b) The path from one adaptive peak to another one requires valley crossing, *i.e.* transition through the disadvantageous intermediate genotypes. (c) Genetic drift allows random walks on the regions of nearly equal fitness. On the rugged landscapes, such regions may look like ridges creating a bypass between highly fit genotypes separated by a fitness valley.

## Valley crossing

On a rugged landscape, nearby adaptive peaks are separated by regions of low fitness, or adaptive valleys. For a population to transit from one peak to another one, downhill steps are needed to make the way through the intermediate states of low fitness — the process called **valley crossing** or tunneling (Figure 2.9b) (Sergey Gavrilets 2004; P. C. Phillips 2008). Generally, with natural selection preventing the fixation of disadvantageous mutations, the population gets constrained to local maxima and can't move to the nearby peak even if it's higher so that achieving it would be evolutionary beneficial.

Shifting balance theory, introduced by S. Wright in 1931, describes the mechanism of how a population can cross fitness valleys and occupy higher fitness peaks (Wright 1931; Coyne, Barton, and Turelli 1997). According to it, genetic drift allows the population to spread and makes it possible to occasionally move to the slope of the nearby fitness peak. Next, selection drives it to the top of the peak, leading to the division of the population into sub-populations, followed by the sub-population occupying the highest peak outcompeting another one. Genetic drift is stronger in small populations; at the same time, in large populations, the high level of standing genetic variation may lead to the appearance of genotypes with multiple mutations (Ochs and Desai 2015; Weissman et al. 2009; Nelson and Grishin 2019; Weinreich and Chao 2005; Meer et al. 2010). Complex mutation events or recombination (in sexual populations) also can give the opportunity to leap over the valley at once, escaping the negative selection acting against the survival of the intermediate states (Weissman, Feldman, and Fisher 2010; Belinky et al. 2019).

Fitness landscapes can contain plateaus: connected regions with nearly equal fitness, where selection is not effective and evolution proceeds mainly under genetic drift. On the rugged epistatic landscapes, particularly multidimensional, such regions may look like ridges and saddles, connecting points of high fitness which look like separate adaptive peaks in some low-dimensional projections and making it possible to tunnel between them (Figure 2.9c) (P. C. Phillips 2008; A. S. Kondrashov, Sunyaev, and Kondrashov 2002; Bakhtin et al. 2021; Gokhale et al. 2009; Katsnelson, Wolf, and Koonin 2019). Another possibility to make the path to the nearby adaptive peak accessible are

drastic changes of the landscapes, caused by environmental fluctuations which may eliminate the valley of low fitness between the peaks (Dodson and Hallam 1977; Steinberg and Ostermeier 2016).

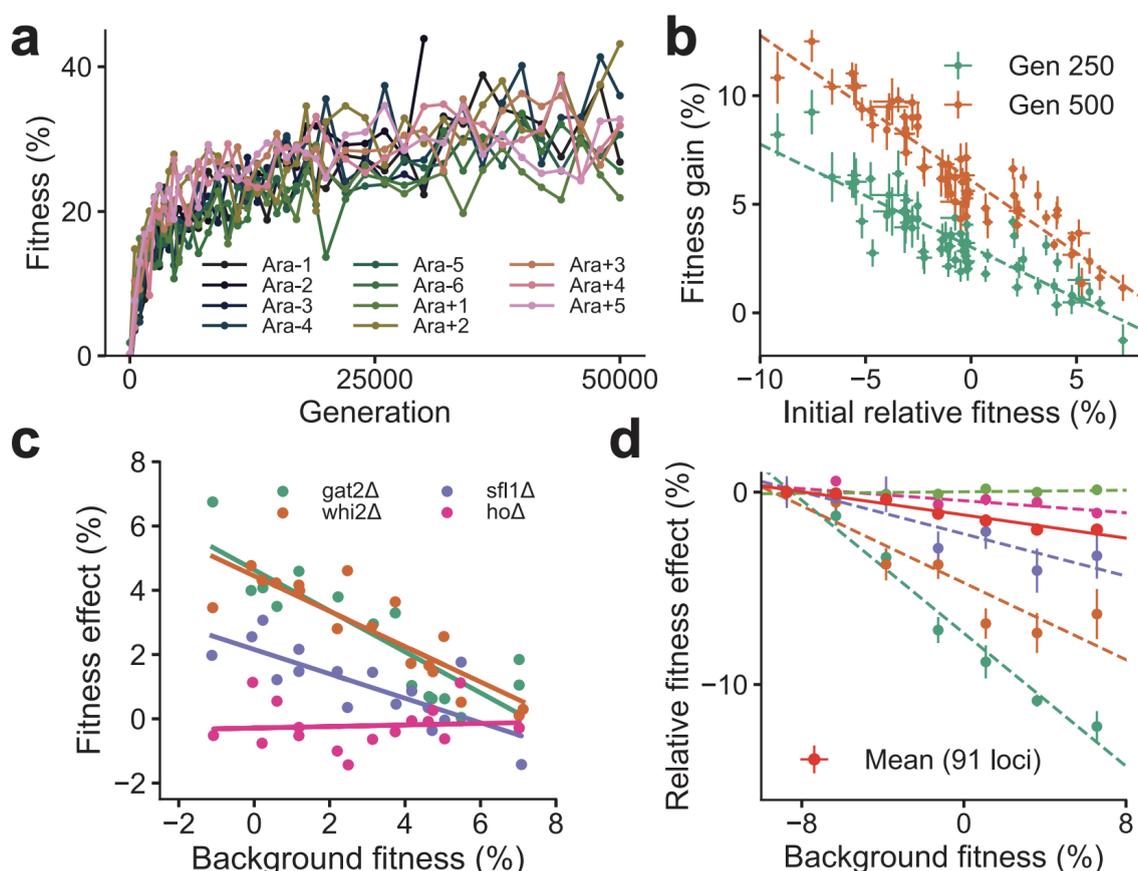
## Dynamics of adaptation

On various models of fitness landscapes, the dynamics of adaptation can be quite diverse. The most direct way to trace the adaptive evolution of a real population is evolution experiments. In the course of such experiments, a population is evolving in laboratory conditions for a long time, accumulating mutations that are beneficial in the new environment; moreover, the fitness of the combinations of obtained mutations may be measured afterwards (de Visser, Cooper, and Elena 2011; Barrick and Lenski 2013). Although evolutionary experiments can shed light on some aspects of adaptation or mutation accumulation, the specificity of evolutionary conditions (*e.g.* starting from genetically homogeneous populations) makes them not generally representative of natural evolution.

The pioneering evolutionary experiment on the T4 bacteriophage demonstrated that recombination increases the rate of adaptation, while epistasis is more efficient in clonal populations (Malmberg 1977). The longest evolutionary experiment, providing unique data on long-term adaptation, is Lenski's long-term experiment on *E. coli* evolution (LTEE). In this experiment, 12 lineages of *E. coli* (six of which acquired mutator phenotype at some point) have evolved for over 80,000 generations (~32 years) (Good et al. 2017; Richard E. Lenski et al. 1991).

Lenski's experiment, as well as other long-term experiments on bacteria or eukaryotes, reveal that the dynamics of fitness gain in the course of adaptation is non-linear: early substitutions give a large increase of fitness, while later ones have a smaller effect (Figure 2.10a) (Schoustra et al. 2016; Johnson et al. 2021; Good and Desai 2015; Good et al. 2017; Wünsche et al. 2017; Wisser, Ribbeck, and Lenski 2013). The decreasing rate of adaptability is at least to some extent explained by diminishing returns epistasis (antagonistic epistasis between beneficial mutations), confirmed by experiments on measuring the fitness of a number of mutations in multiple genetic contexts: a mutation that is advantageous on a low-fit background is less beneficial on a high-fit background (Figure 2.10b-d) (Wei and Zhang 2019; Kryazhimskiy et al. 2014).

Diminishing returns epistasis between beneficial mutations is predicted to emerge under multiple fitness landscape models. On unidimensional landscapes, it can be defined by the smooth shape of the adaptive peak, with smaller fitness effects near the fitness maximum (Schoustra et al. 2016). In multidimensional models, diminishing returns epistasis results from the appearance of multidimensional adaptive ridges, forming curved but smooth epistatic paths of increasing fitness (Lyons et al. 2020). Diminishing returns effects can also come from the modular structure of gene networks due to fitness “saturation” within modules (Wei and Zhang 2019; Reddy and Desai 2021).

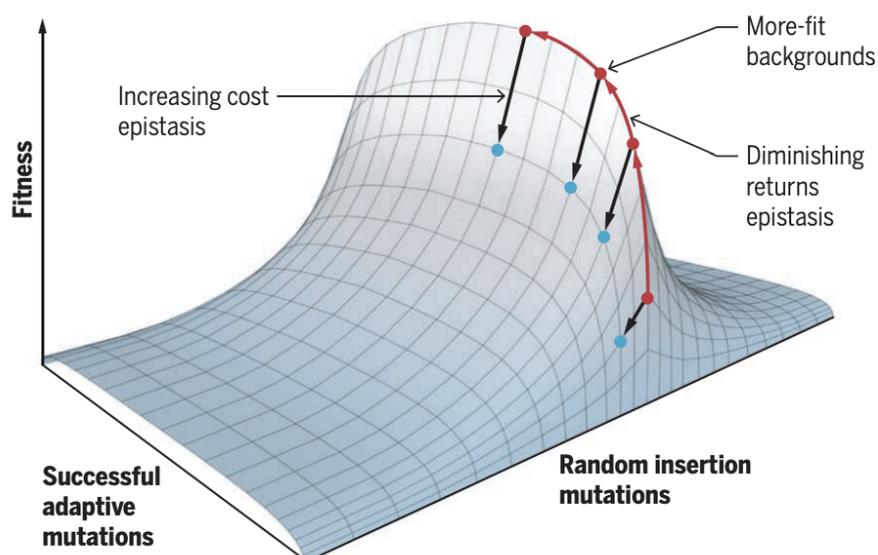


**Figure 2.10. Diminishing returns and increasing cost epistasis affecting adaptability in evolutionary experiments** (Reddy and Desai 2021). **(a)** During the *E. coli* LTEE experiment, the rate of the fitness gain declines with time (Wiser, Ribeck, and Lenski 2013). **(b)** Similarly, the less fit strains of *S. cerevisiae* increase their fitness faster than more fit strains (Kryazhimskiy et al. 2014). **(c)** Direct fitness measurements in *S. cerevisiae* show that the decline of fitness gain is explained by

diminishing returns epistasis (Kryazhimskiy et al. 2014). **(d)** Deleterious mutations demonstrate the opposite trend of increasing cost epistasis (Johnson et al. 2019).

## Mutational robustness

The opposite point of view on the adaptation process is how robust the adapted genotype is to the emergence of deleterious mutations and whether the robustness depends on the genotype's fitness. Multiple studies consider the relationship between robustness and adaptability in diverse epistasis models, presenting controversial results (Draghi et al. 2010; de Visser et al. 2003; Masel and Trotter 2010; Wagner 2008). An important phenomenon recently presented by mutational experiments in *S. cerevisiae* is increasing cost epistasis — deleterious mutations have a larger effect on highly fit genotypes and a smaller effect on less fit genotypes (Johnson et al. 2019) (Figure 2.10d). This implies that genotypes become less robust to deleterious mutations while increasing their fitness in the course of adaptation (Figure 2.11). The effects of diminishing returns and increasing cost epistasis are at first sight inconsistent and require specific landscape shape — so that adaptation makes the genotype both less adaptable and less robust to deleterious changes (Johnson et al. 2019). However, both evolutionary patterns are shown to emerge in multidimensional epistasis models (Lyons et al. 2020; Reddy and Desai 2021).



**Figure 2.11. The fitness landscape model with both diminishing returns and**

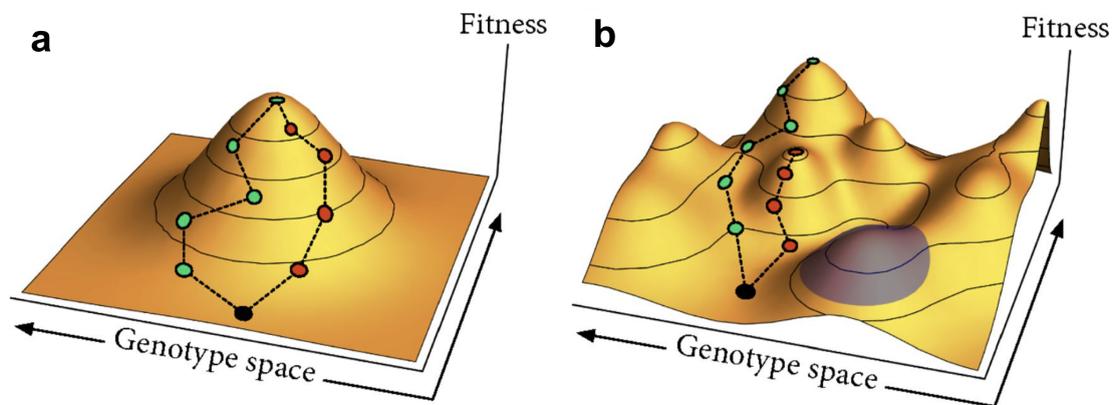
**increasing cost epistasis** (Miller 2019). Red arrows show the adaptive path consisting of beneficial mutations under diminishing returns epistasis (the fitness gain is smaller at highly fit genotypes). Black arrows show the effect of deleterious mutations, which are under increasing cost epistasis (the fitness loss is larger at highly fit genotypes).

## Predictability of evolution

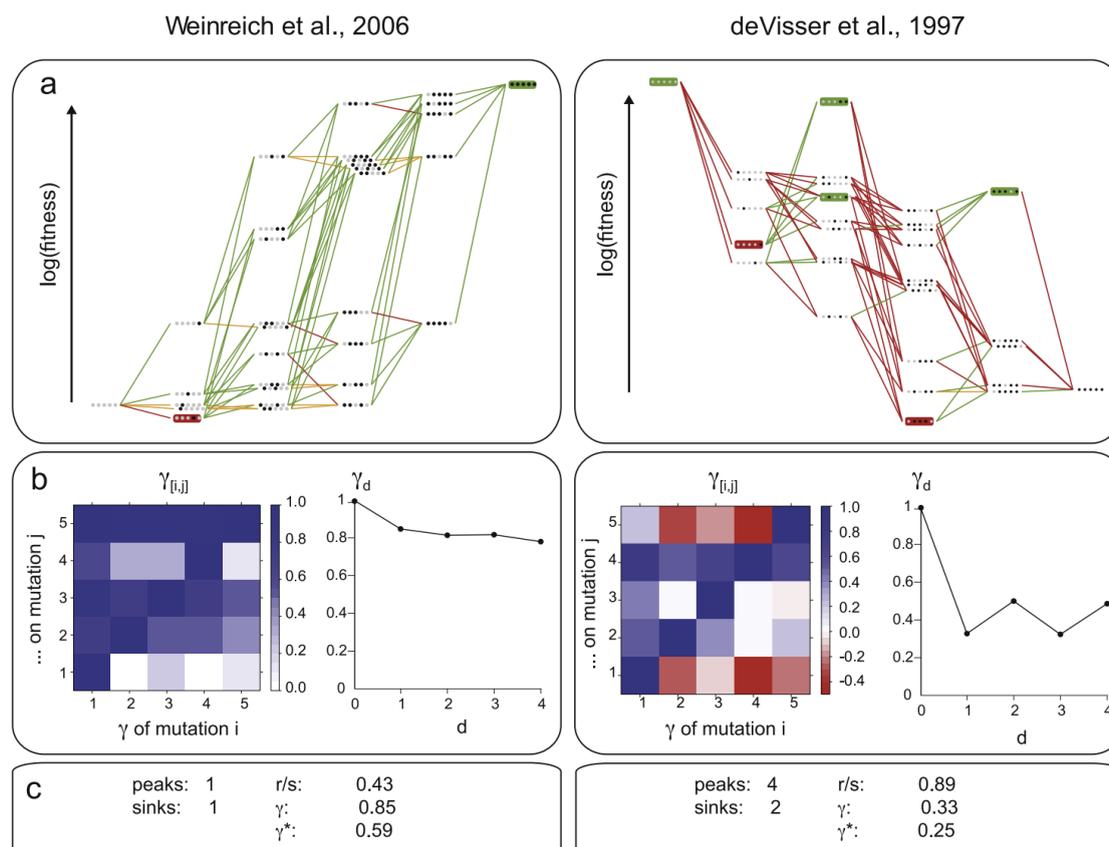
How unambiguously the path of adaptation can be defined by selection varies among different landscape models. The predictability of adaptation is an important feature of the fitness landscape, describing whether it is possible to predict the adaptive path along the landscape knowing the fitness effects of individual mutations on a certain background or the preceding evolutionary path (de Visser et al. 2018; de Visser and Krug 2014).

Predictability strongly depends on the ruggedness of the landscape. The ruggedness emerges under sign epistasis, which makes the direction of selection acting on a mutation dependent on the genetic background, forming complex multidimensional adaptive paths (Weinreich, Watson, and Chao 2005). Ruggedness constraints evolution, limiting the number of available adaptive paths (de Visser et al. 2018; Ferretti, Schmiegel, and Weinreich 2016; D. A. Kondrashov and Kondrashov 2015). If there is only one adaptive peak, the outcome of the adaptive evolution of a sequence on such a landscape is predictable: under selection, it will achieve the global maximum (Figure 2.12a). If the landscape is rugged, *i.e.* there are multiple local maxima, population may get stranded on the nearest adaptive peak and never achieve the global optimum (Figure 2.12b) (Fragata et al. 2019; D. A. Kondrashov and Kondrashov 2015; Poelwijk et al. 2007; de Visser and Krug 2014; Van Cleve and Weissman 2015). In an extreme case of a highly rugged HoC landscape, where there is no correlation between fitness effects of mutations in different genetic contexts, every evolutionary step changes the local landscape completely and makes it impossible to predict the subsequent beneficial mutations. Real landscapes are shown to be somewhere in-between, combining additive and rugged components (Figure 2.13) (de Visser and Krug 2014; Sarkisyan et al. 2016; Weinreich et al. 2006; Bank et al. 2016; Poelwijk et al. 2007; Ferretti, Schmiegel, and Weinreich 2016).

In the terms of ruggedness and predictability, landscapes can be functionally characterized by several measurable parameters, such as the number of local adaptive peaks and valleys (sinks), the number of accessible adaptive paths, the amount of sign epistasis or roughness/slope ratio (Figure 2.13c) (de Visser and Krug 2014; Ferretti et al. 2018; Ferretti, Schmiegelt, and Weinreich 2016; Van Cleve and Weissman 2015; Szendro, Schenk, et al. 2013).



**Figure 2.12. Predictability of evolution on landscapes of different complexity** (Van Cleve and Weissman 2015). **(a)** Evolution of a genotype on a smooth landscape with a single adaptive peak converges to the fitness maximum. **(b)** On the rugged landscape with multiple peaks, evolutionary trajectories might end on different local maxima. Circles show two evolutionary paths (green and red). The blue area shows a small part of the rugged landscape, which contains a single adaptive peak and appears smooth.



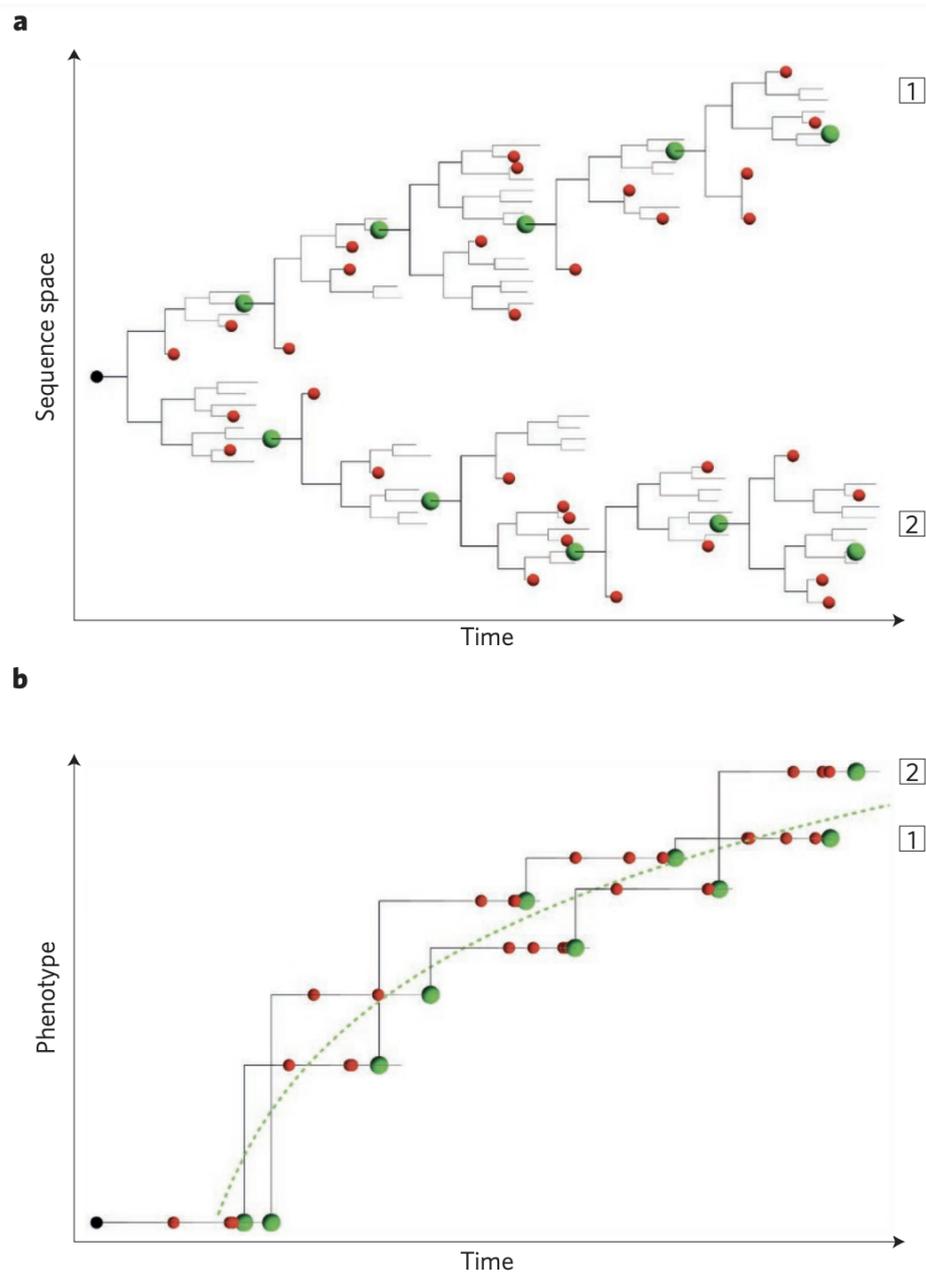
**Figure 2.13. The ruggedness of empirical fitness landscapes** (Ferretti, Schmiegelt, and Weinreich 2016). **(a)** The less rugged landscape of five biallelic sites in *E. coli*  $\beta$ -lactamase (Weinreich et al. 2006) and the more rugged landscape of five deleterious mutations in *Aspergillus niger* (de Visser, Hoekstra, and van den Ende 1997). Genotypes corresponding to adaptive peaks are highlighted in green, genotypes corresponding to sinks of fitness (or adaptive valleys) are highlighted in red. **(b)** Epistasis between mutations in these fitness landscapes: blue — no epistasis, white — non-sign epistasis, red — sign epistasis. **(c)** The measures of the ruggedness of the landscapes ( $\gamma$  — correlation of fitness effects of mutations in the neighboring genotypes,  $\gamma^*$  — correlation of the signs of fitness effects of mutations in the neighboring genotypes,  $r/s$  — roughness/slope ratio).

Evolution experiments make it possible to directly quantify the reproducibility and contingency in evolution (de Visser and Krug 2014; Poelwijk et al. 2007). On the smooth landscapes, all uphill adaptive paths converge to the same adaptive peak, independently of the starting point. On the rugged landscapes, diverging populations may follow distinct complex multidimensional trajectories, occupying different peaks and appearing in reproductive isolation because of Dobzhansky-Muller incompatibilities (Orr 1995; D. A. Kondrashov and Kondrashov 2015).

The convergence of evolution can be addressed on the level of phenotypes, including fitness, or genotypes. Paths of adaptation to similar environments often show convergence on the phenotypic level or on the level of genes, but less commonly by specific mutations (Figure 2.14) (Lässig, Mustonen, and Walczak 2017; Kryazhimskiy et al. 2014; Lang, Botstein, and Desai 2011; Lang and Desai 2014).

The effect of epistasis on the predictability of evolution is different in terms of convergence and repeatability of adaptive trajectories. On rugged landscapes, the outcome of the adaptation strongly depends on the starting genotype and initial evolutionary steps: the probability of several independently evolving populations achieving the same adaptive peak is low (Salverda et al. 2011; D. A. Kondrashov and Kondrashov 2015; Fragata et al. 2019). At the same time, epistasis reduces the set of available evolutionary pathways and therefore specifies the accessible order of mutations (de Visser et al. 2018; Weinreich et al. 2006; Franke et al. 2011; D. A. Kondrashov and Kondrashov 2015). Therefore, epistasis constrains adaptation, forcing the evolving populations to follow the same paths and increasing the repeatability of adaptation. The reproducibility of specific evolutionary steps in the course of adaptation is shown in multiple evolution experiments (usually not on the level of specific nucleotide mutations, but single genes) (Olivier Tenaillon et al. 2012; Graves et al. 2017; Woods et al. 2006). However, standing genetic variation may also facilitate access to the distinct fitness peaks, decreasing the predictability of evolution (Zheng, Payne, and Wagner 2019).

Clearly, only a small subset of theoretically accessible adaptive paths are actualized in the evolution of real populations, both in nature or experimental evolution. Moreover, there are additional factors affecting the predictability of evolution, such as weak or strong mutation regimes, population size, and clonal interference in asexual populations (Lässig, Mustonen, and Walczak 2017; de Visser and Krug 2014; Szendro, Franke, et al. 2013; Jain and Krug 2007). Therefore, empirical observations don't make it possible to reconstruct the full landscape — however, they can be used to describe the main principles of its shape and structure (D. A. Kondrashov and Kondrashov 2015). Moreover, the general ruggedness of the full landscape doesn't necessarily affect the sequence evolution on small scales, if it's already constrained to one local adaptive peak (Figure 2.12b) (Van Cleve and Weissman 2015).



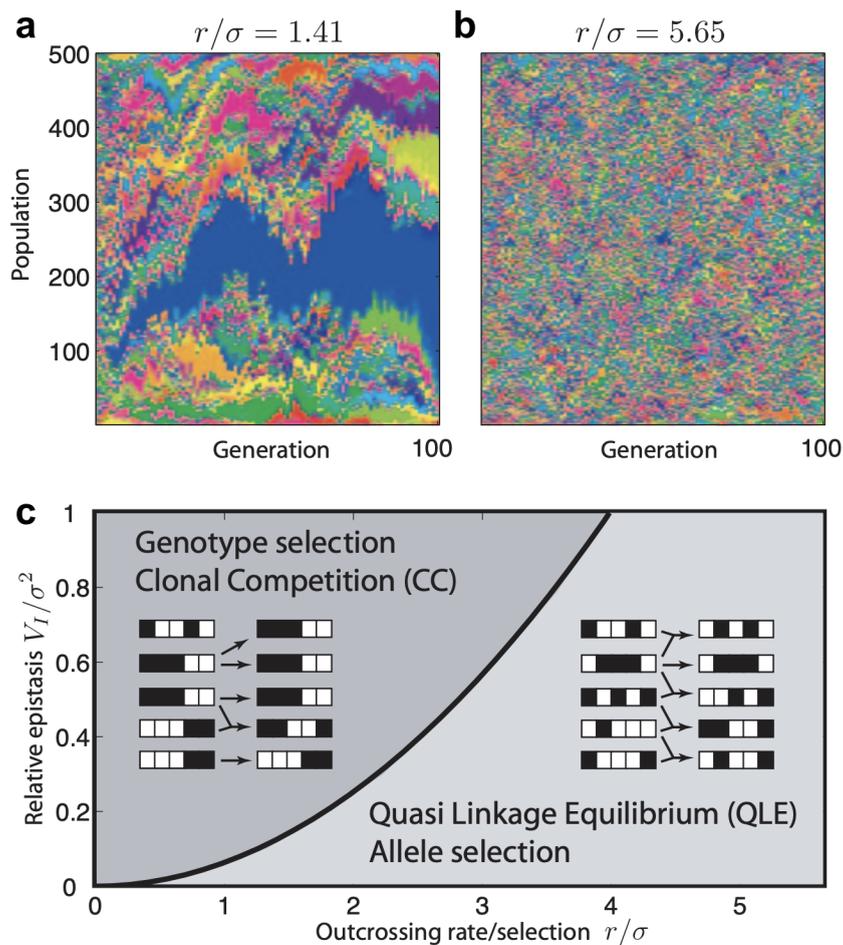
**Figure 2.14. Predictability of evolution on the level of genotypes and phenotypes** (Lässig, Mustonen, and Walczak 2017). **(a)** Two populations starting from the same initial genotype follow different adaptive paths, occupying distinct adaptive peaks. **(b)** Despite the sequence divergence, the dynamics of a phenotype trait in these populations is similar and converges to the same high fitness value. Green — positively selected genotypes, red — negatively selected genotypes.

## **Epistasis and recombination**

In asexual populations, the only mechanism of accumulating genetic variability is the emergence of new mutations. Under sexual reproduction (amphimixis), recombination is a way to reshuffle the parental genotypes, resulting in new combinations of alleles.

The selective effect of a specific recombination event depends on the complexity of the fitness landscape. On a fully additive landscape, reshuffling of the genotypes doesn't change fitness effects of the individual alleles constituting these genotypes. On an epistatic landscape, the consequences of a crossing-over event are harder to predict: selection coefficients of the present alleles may change in a newly established combination, decreasing the heritability of fitness (Falconer and Falconer 1989; Zuk et al. 2012). In the extreme case of HoC landscape, with no additive component of fitness, reshuffling of the parental genotypes makes fitness of the offspring independent of the parents' fitness.

Without recombination, selection operates on entire linked genotypes, leading to clonal competition (Figure 2.15a) (Ronald A. Fisher 1930; H. J. Muller 1932; Franklin and Lewontin 1970; R. A. Neher and Shraiman 2009). In this case, fitness effect of a particular mutation is defined only in the given genetic context. Recombination can break combinations of alleles, leading to quasi linkage equilibrium between loci so that selection is able to promote or eliminate specific alleles (Figure 2.15b) (Ronald A. Fisher 1930; M. Kimura 1965; Franklin and Lewontin 1970; N. H. Barton 1995). In the allele selection regime, the dynamics of the allele frequency is less dependent on the stochasticity of its appearance in a specific context and therefore is more predictable. Whether selection acts on genotypes or individual alleles depends both on the recombination rate and the strength and abundance of epistasis within the considered genotypes (Figure 2.15c) (Franklin and Lewontin 1970; R. A. Neher and Shraiman 2009).

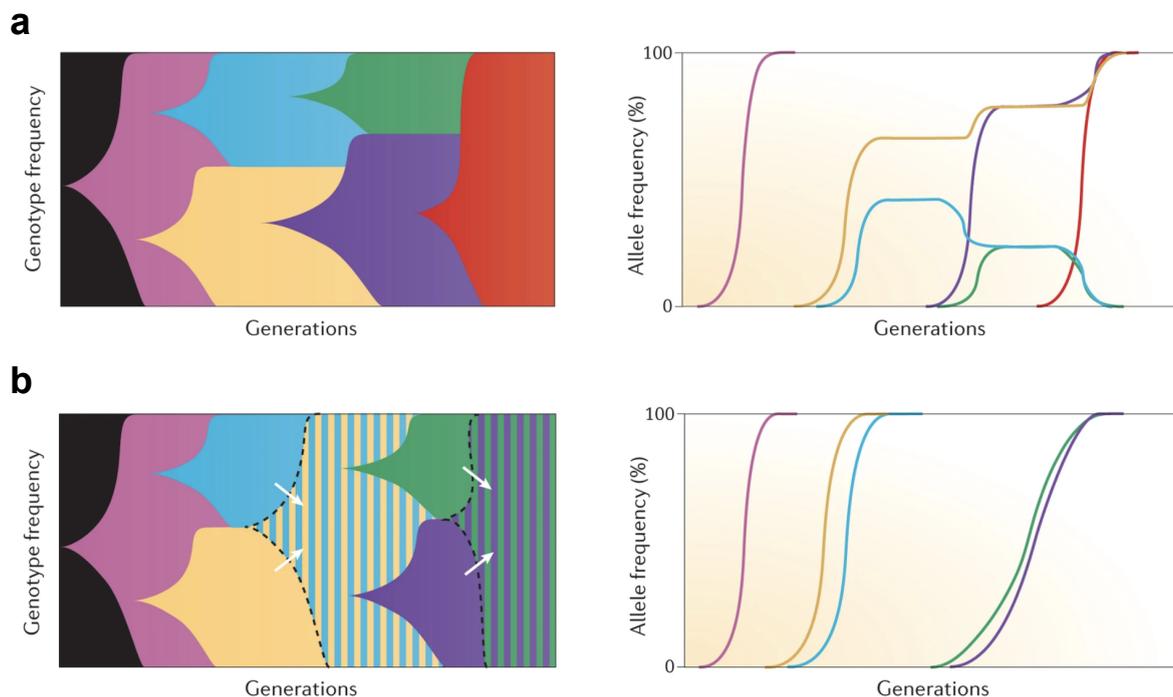


**Figure 2.15. Epistatic selection with or without recombination** (R. A. Neher and Shraiman 2009). (a) Under low recombination rate ( $r$ ) and strong epistasis, selection acts on the whole genotypes. (b) Under high recombination and weak epistasis, the units of selection are individual alleles, which can be reshuffled by recombination. In this model, the fitness variance of fitness ( $\sigma^2$ ) is fully explained by pairwise epistasis. Colors correspond to different genotypes. (c) The transition point between genotype and allele selection regimes depends on the recombination rate and the strength of additive and epistatic selection.  $V_I/\sigma^2$  — the proportion of fitness variance attributed to epistasis.

While considering the evolution of populations, the effects of recombination may be complex even in a non-epistatic case. One of the outcomes of linkage in adapting asexual populations is clonal interference. In a large clonal population, multiple beneficial mutations may arise and segregate simultaneously, competing with each other and mimicking the effect of stochastic drift (Figure 2.16a) (Ronald A. Fisher 1930; H. J. Muller 1932; Kim and Orr 2005; S.-C. Park and Krug 2007; Gerrish and Lenski 1998). In sexual populations, recombination increases the probability for beneficial mutations to

fix, because it unlinks beneficial mutations from their genetic background, allowing them to fix faster and therefore facilitating adaptation (Fisher-Müller hypothesis, Figure 2.16b) (Ronald A. Fisher 1930; H. J. Muller 1932; Kim and Orr 2005; N. H. Barton 1995). The dynamics of allelic frequency of deleterious variants is also prone to clonal selection: without the possibility to discriminate disadvantageous alleles from their background, their elimination is less efficient (Müller's ratchet) (H. J. Muller 1964; N. H. Barton 2010).

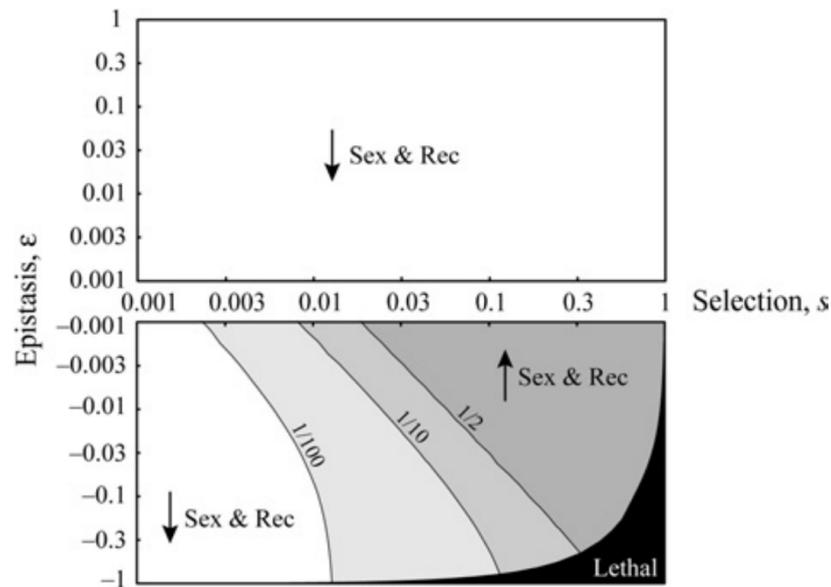
Selection acting in a specific genomic locus can also affect the patterns of the genetic variation in the entire linked region. For example, fast fixation of a positively selected allele results in the associated fixation of the neutral variants present in the linked genotype (genetic hitchhiking), leading to local reduction of genetic variation (selective sweeps) (Maynard and Haigh 1974; Paquin and Adams 1983; N. H. Barton 1998). The consequences of linked selection under mutation-selection equilibrium can be diverse and depend on the sign and strength of selection. If some loci are under strong purifying selection, linked neutral alleles will be also eliminated due to background selection (BGS), reducing genetic diversity (B. Charlesworth, Morgan, and Charlesworth 1993; N. H. Barton 2010). At the same time, the presence of recessive deleterious alleles is shown to lead to the associative overdominance (AOD) in the linked neutral loci (Ohta 1971; Zhao and Charlesworth 2016; Gilbert et al. 2020). This effect is akin to the selective advantage of heterozygotes and is able to, on the contrary, increase the genetic variability in the linked loci. Likewise, local genetic variation may be maintained due to linkage to a locus under balancing selection, caused, for example, by frequency-dependent selection or overdominance (Deborah Charlesworth 2006; James F. Crow 1987; Ayala and Campbell 1974). Another outcome of linkage demonstrated in a variety of evolutionary systems is Hill-Robertson interference: simultaneous segregation of linked weakly selected alleles can impede the efficiency of selection acting on them (W. G. Hill and Robertson 1966; Brian Charlesworth 2012; Comeron, Williford, and Kliman 2008; Roze and Barton 2006).



**Figure 2.16. Adaptation in asexual and sexual populations** (Barrick and Lenski 2013). **(a)** In asexual populations, the dynamics of genotypes' frequencies is clonal — multiple positively selected mutations interfere, retarding adaptation. **(b)** In sexual populations, beneficial mutations may be combined within a single genotype by recombination, allowing them to fix simultaneously. Different colors represent genotypes carrying different beneficial mutations.

The evolutionary consequences of recombination are even more complex in the context of epistasis. Generally, random shuffling of the genotypes by recombination can reduce the mean fitness of a population while at the same time increasing its variance and therefore evolvability (Brian Charlesworth 1990; N. H. Barton and Charlesworth 1998). Under multidimensional or unidimensional antagonistic epistasis, recombination breaks the coadapted combinations of positively interacting alleles, having a deleterious effect on fitness (Brian Charlesworth 1990; F. A. Kondrashov and Kondrashov 2001b; N. H. Barton 2010). Under synergistic epistasis between derived alleles, recombination may be advantageous by preserving their repulsion (Brian Charlesworth 1990; N. H. Barton 1995). However, these effects strongly depend on the regime of selection and on the linkage between considered loci (S. P. Otto and Feldman 1997; S. P. Otto and Gerstein 2006; Desai, Weissman, and Feldman 2007; Kouyos, Otto, and Bonhoeffer 2006). Given

these, the evolutionary advantage of sex and its abundance in natural populations remains controversial (Figure 2.17) (A. S. Kondrashov 2018; N. H. Barton and Charlesworth 1998; S. P. Otto and Gerstein 2006).



**Figure 2.17. Evolutionary advantage of sexual reproduction under different epistasis modes** (S. P. Otto and Gerstein 2006). Recombination between two linked loci is disadvantageous if mutations in these loci are positively interacting or are under weak negative selection (white); it can be favorable in the case of weak negative epistasis between strongly deleterious mutations (gray).

Linked selection is shown to strongly affect genetic diversity in natural populations. A striking example is decrease of diversity in the genomic region linked to a beneficial mutation fixed in the course of selective sweep. Naturally, linkage effects on diversity are more pronounced and involve larger segments of the genome if recombination rate is low (*e.g.* sex chromosomes, centromeres and telomeres) (Nordborg, Charlesworth, and Charlesworth 1996; B. Charlesworth 1996; N. H. Barton 2010; Ellegren and Galtier 2016). Positive correlation between recombination rate and diversity level along the genome was observed in multiple species, indicative of the action of background selection (Campos et al. 2014; Corbett-Detig, Hartl, and Sackton 2015; Sella et al. 2009; Hough et al. 2017). The decay of diversity caused by linkage effects competes with the mutagenic effect of recombination, although the impact of the latter on genetic diversity is not so high (Spencer et al. 2006; Hellmann et al. 2003; Arbeithuber et al. 2015).

Tightly linked loci, associated with complex phenotypes, are shown to form supergenes, maintained within populations by balancing selection, *e.g.* negative frequency-dependent selection, overdominance or associative overdominance (Mather 1950; Joron et al. 2011; Thompson and Jiggins 2014; D. Charlesworth and Charlesworth 1975; Gutiérrez-Valencia et al. 2021; Ohta 1971). Overdominance may arise due to epistatic interactions, resulting in the persistence of coadapted gene complexes (Brian Charlesworth and Charlesworth 1973; Faria et al. 2019). The first example of such coadapted gene complexes were inversions within *D. melanogaster* populations — large inversions suppress recombination, so that they are evolving as a single unit (T. Dobzhansky and Sturtevant 1938; Mather 1950).

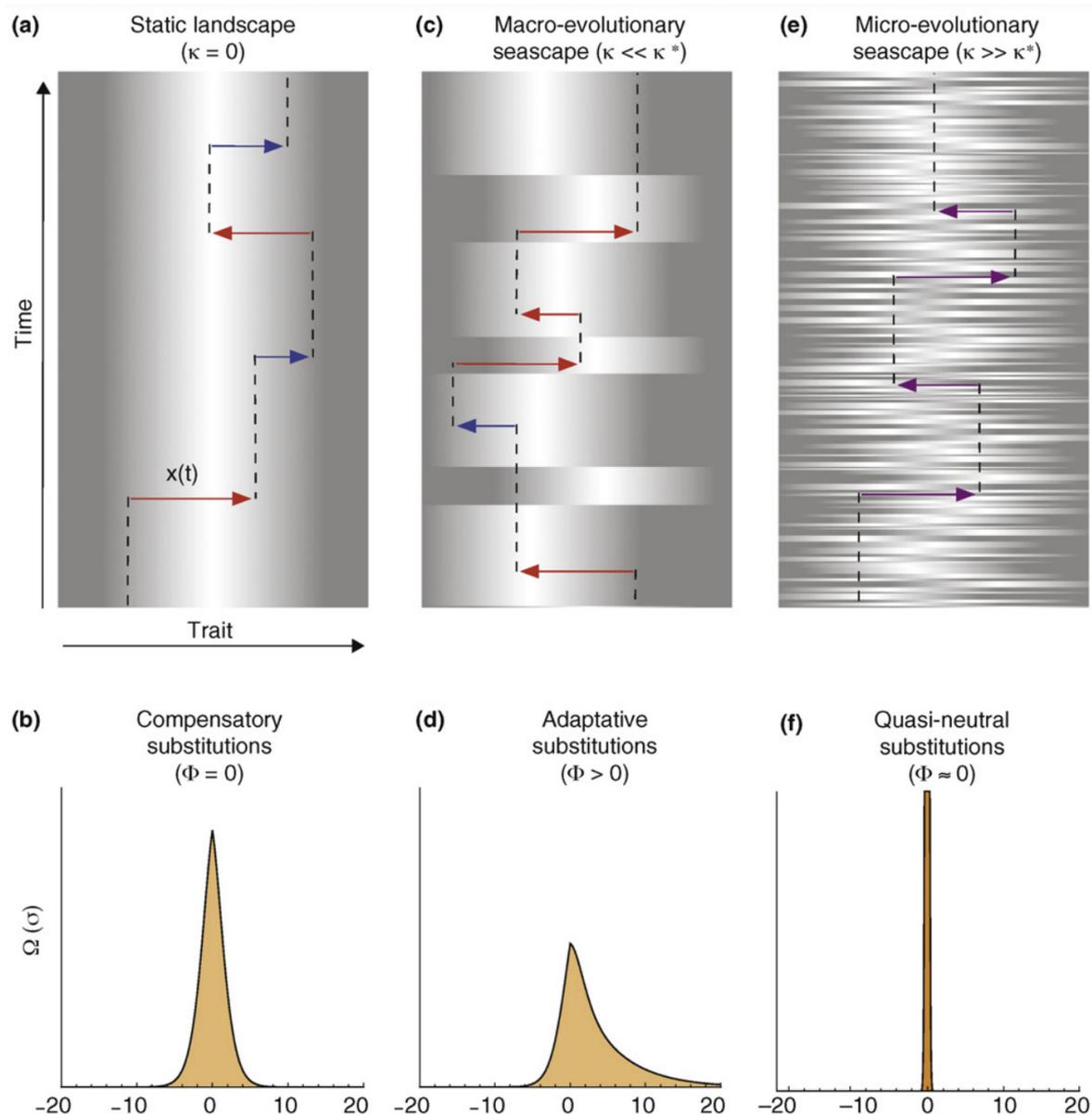
## Dynamic fitness landscapes

Fitness landscapes are usually described as static, or time-independent. However, fitness of a genotype can change due to external or internal factors. The landscape changes can alter the strength or even the sign of selection, promoting or restraining adaptive evolution.

### **Environmental fluctuations**

Fitness of a certain genotype is defined in the context of environmental and ecological conditions. If they change, a genotype highly fit under previous conditions may appear less adapted to the new environment. On a static landscape, the population eventually comes to mutation-selection equilibrium, when the mean fitness value of a population is constant. Sudden random change of the landscape can lead to disappearance or displacement of the currently occupied fitness peak, disrupting the equilibrium and driving adaptation to the newly established optimum (Figure 2.8c) (Wright 1932; V. Mustonen and Lässig 2007; Ronald A. Fisher 1930; John H. Gillespie 1991).

The studies on deformability of landscapes in changing environments also focus on the questions on pleiotropy and on the predictability of evolution in a new environment: what proportion of mutations are susceptible to environmental changes, what are the general patterns of the genotype-environment interactions, and whether a mutation neutral or deleterious under some conditions can be maintained cause it is beneficial in an alternative environment (C. Li and Zhang 2018; Fragata et al. 2018; Bajić et al. 2018; Ho and Zhang 2018; Hermsen, Deris, and Hwa 2012; Masel 2006; de Vos et al. 2013; Bergland et al. 2014; Hietpas et al. 2013; S. Wang and Dai 2019). Understanding patterns of adaptation under fluctuating conditions is necessary to predict the evolution of pathogens, such as the development of antibiotic resistance in bacteria or dynamics of virus-host co-evolution (Schrag, Perrot, and Levin 1997; R. E. Lenski 1998; Hegreness et al. 2008; Bhatt, Holmes, and Pybus 2011; R. A. Neher and Leitner 2010; Pennings, Kryazhimskiy, and Wakeley 2014).



**Figure 2.18. Macro- and microevolutionary fitness seascapes** (Ville Mustonen and Lässig 2009). **(a, b)** On the static landscape, the population reaches an equilibrium state. **(c, d)** Rare changes of the landscape lead to non-equilibrium population dynamics and sustained positive selection. **(e, f)** Frequent landscape changes /smooth/ the action of selection, resulting in a quasi-neutral regime of evolution.  $\Omega(\sigma)$  – distribution of selection coefficients of mutations,  $\Phi$  – fitness flux (a measure of fitness gain during adaptation, defined in (Ville Mustonen and Lässig 2009), with  $\Phi > 0$  corresponding to the action of positive selection).

If the rate of environmental fluctuations is comparable with the evolution rate, they cause the non-equilibrium dynamics of evolution, driving adaptation. However, if the landscape changes too frequently, *i.e.* fitness landscape changes faster than the allelic

composition of a population can adjust to the newly established selection regime, the efficacy of selection reduces (M. Lynch 1987; Ville Mustonen and Lässig 2009; V. Mustonen and Lässig 2007, 2010; Trubenová et al. 2019). The population fails to adapt to new conditions before they change, so that evolution becomes “quasi-neutral” (Figure 2.18). Such a dynamic fitness landscape is called “seascape”, reflecting its variability (Ville Mustonen and Lässig 2009).

### **Frequency-dependent selection**

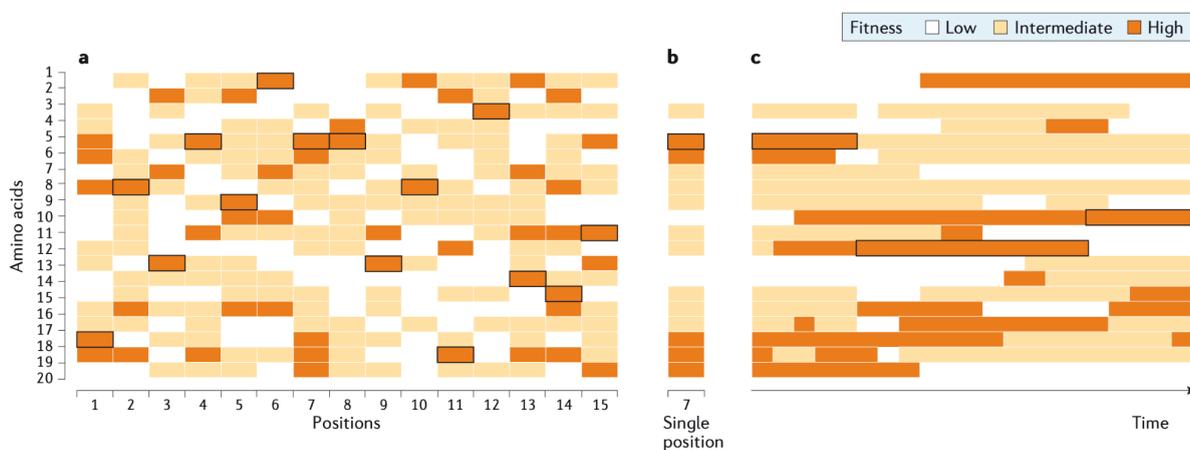
The shape of the fitness landscape can be dependent on the current state of the population composition. Under frequency-dependent selection, the fitness of an allele is defined by the fraction or by the absolute number of individuals carrying this allele (Ronald A. Fisher 1930; Ayala and Campbell 1974). In the case of negative frequency-dependent selection (NFDS), fitness and frequency are negatively correlated, so that rare alleles are under positive selection and common alleles are under negative selection. As a consequence, NFDS maintains persistence of multiple alleles within a locus, being the form of balancing selection (Takahata and Nei 1990). The long-term effect of NFDS is elevated genetic diversity in the linked genomic region (Deborah Charlesworth 2006).

NFDS can arise due to multiple aspects of the species' biology, for example, because of specific patterns of sexual reproduction (*e.g.* assortative mating or self-incompatibility) (Gigord, Macnair, and Smithson 2001; Sarah P. Otto, Servedio, and Nuismer 2008; Conover and Van Voorhees 1990; Delph and Kelly 2014). Frequency-dependent selection can be also caused by ecological interactions, *e.g.* co-evolution with another species or pathogens (Borghans, Beltman, and De Boer 2004; Barrett et al. 1988; Tellier and Brown 2011; Carius, Little, and Ebert 2001). However, it differs drastically from random environment-driven changes of the landscape, which aren't associated with the allelic composition of the population.

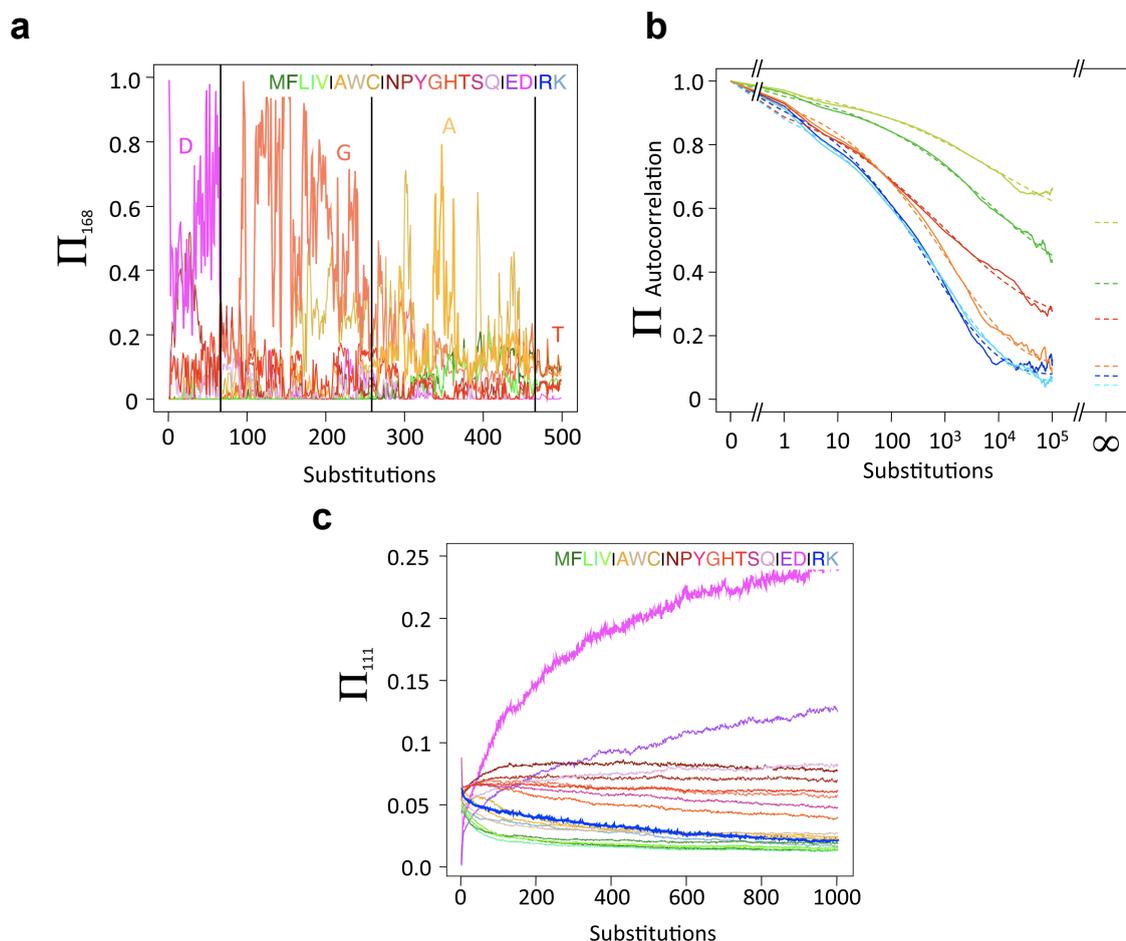
Another form of selection maintaining genetic diversity in the locus is overdominance, *i.e.* selection advantage of heterozygotes (Ronald Aylmer Fisher 1922; T. Dobzhansky 1950; Hedrick 2012; James F. Crow 1987; Takahata and Nei 1990). A similar effect can be produced by segregation of linked recessive alleles under weak negative selection, due to drift effects in finite populations (Ohta 1971; Zhao and Charlesworth 2016).

## Epistatic changes of allele's fitness

Full fitness landscapes are extremely high-dimensional. To characterize general principles of evolutionary dynamics on such landscapes, it can be to some extent reduced to the changes of properties of individual genomic loci. The fitness landscape at one locus (single-position fitness landscape, or SPFL) is a one-dimensional cross-section of the full landscape and can be represented with a vector of size  $K$ , where  $K$  is the number of possible alleles (e.g.  $K = 20$  if we consider single amino acid site) (Figure 2.19) (Bazykin 2015). The location of this cross-section is defined by the current state of the genetic background. Without epistasis, the shape of SPFL and therefore fitness effect of a mutation in the considered site doesn't depend on the genetic context. On epistatic landscapes, it may change with time even if the full fitness landscape remains static, due to substitutions in epistatically interacting sites (Figure 2.20a) (Bazykin 2015; Starr and Thornton 2016; David D. Pollock, Thiltgen, and Goldstein 2012; D. A. Kondrashov and Kondrashov 2015; Van Cleve and Weissman 2015). In this case, replacements in the background part of the genome are “external”, meaning that we can detect them only implicitly by their epistatic effect on the fitness of alleles in the considered locus, which can be inferred by changes of allelic frequencies or pattern of substitutions. Understanding the patterns of such changes may be used to infer epistasis in sequence divergence between species (Povolotskaya and Kondrashov 2010; A. S. Kondrashov et al. 2010; A. S. Kondrashov, Sunyaev, and Kondrashov 2002; Goldstein and Pollock 2017).



**Figure 2.19. Single-position fitness landscape** (Bazykin 2015; Storz 2016). **(a)** Single mutation fitness landscape of a protein sequence of length 15. Colors indicate the fitness of the allele, the currently present variants are outlined. **(b)** The fitness vector of all possible alleles for a specific position represents a single position fitness landscape. **(c)** Allelic preferences in a chosen position can change with time, resulting in the replacement of the currently present allele.



**Figure 2.20. Changes of allele fitness due to substitutions in epistatically interacting sites** (David D. Pollock, Thiltgen, and Goldstein 2012). **(a)** In simulations of evolution on a purple acid phosphatase structure, propensities of amino acid alleles at site 168 ( $\Pi_{168}$ ) change due to substitutions in epistatically interacting sites. Amino acid propensities are calculated as equilibrium frequencies in the current genomic background. Black lines show amino acid replacements in the considered site. **(b)** The correlation between initial and current propensities vector at a site declines with the accumulation of replacements in other genomic positions. Green and lime – amino acid sites buried in the protein structure; orange and red – partially exposed sites; blue and cyan – exposed sites. **(c)** The fitness of the allele currently occupying an amino acid site (here, aspartic acid D at site 111) increases with time due to epistasis, demonstrating entrenchment. Position 168 is the example of a site exposed in the protein structure, while position 111 is the example of spatially buried sites.

In multiple models of epistasis, the magnitude of fitness changes increases with the number of substitutions in the genetic background (Figure 2.20b). The characteristic lifespan of the correlation between the initial and the subsequent SPFLs is indicative of the strength and abundance of epistasis and of the predictability of evolution on the studied landscape (David D. Pollock, Thiltgen, and Goldstein 2012; Ferretti, Schmiegelt, and Weinreich 2016; Sarkisyan et al. 2016; Pokusaeva et al. 2019).

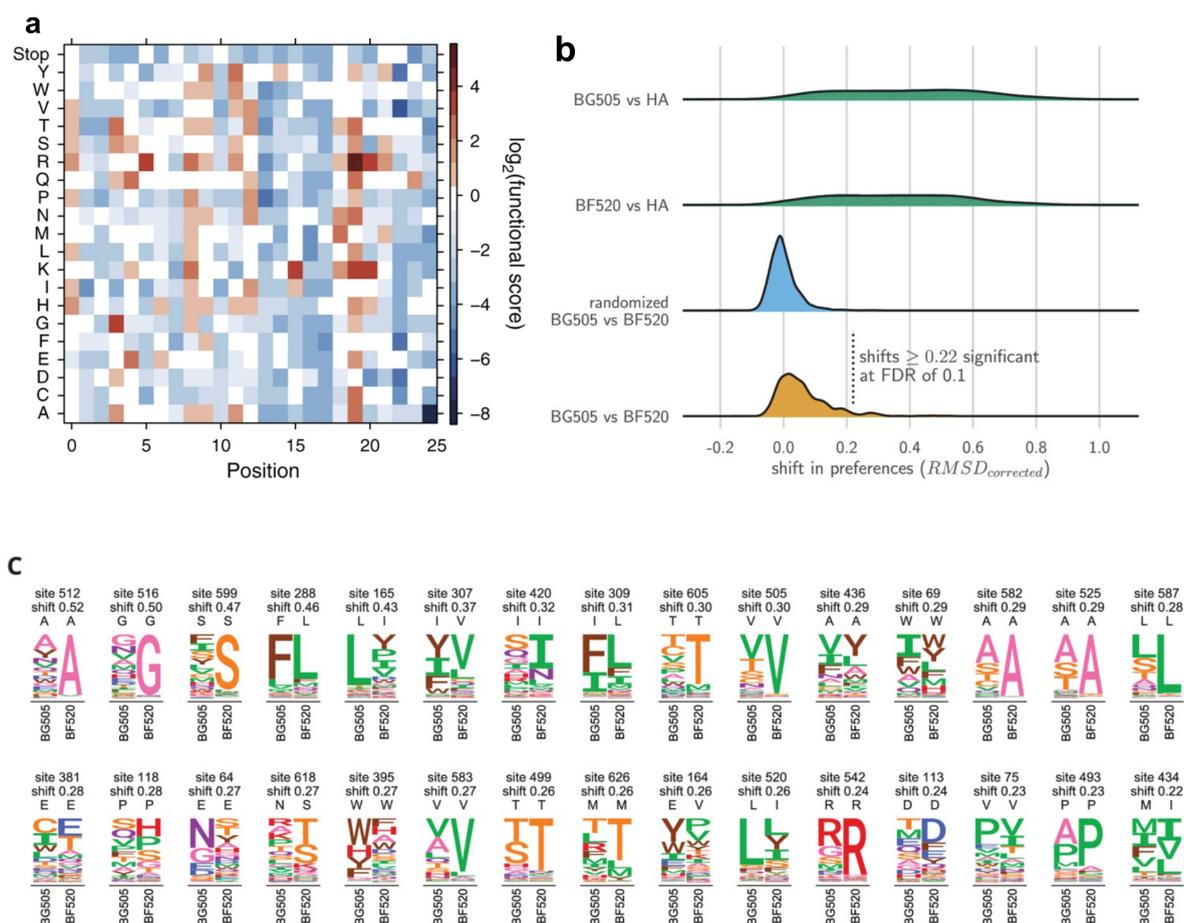
If the background substitutions are not accumulated randomly, but are fixed according to their selection coefficients (*i.e.* form an adaptive path), epistatic selection will advance the coexistence of combinations of positively interacting alleles. Due to such co-evolution, the fitness of the variant currently occupying a genomic site will on average increase, experiencing **entrenchment**, or evolutionary Stokes shift (Figure 2.20c) (David D. Pollock, Thiltgen, and Goldstein 2012; Flynn et al. 2017; Goldstein and Pollock 2017; Starr et al. 2018; Shah, McCandlish, and Plotkin 2015). Under entrenchment, climbing onto an adaptive peak will be coupled with the decline of robustness to deleterious mutations, making it related to increasing-costs epistasis (Lyons et al. 2020; Reddy and Desai 2021).

## Empirical fitness landscapes

### Landscapes of homologous sequences

Deep mutational scanning (DMS) is a high-throughput method of measuring fitness values of a set of genotypes (Fowler and Fields 2014). The DMS experiment consists of three steps: generation of the library of mutant genotypes, selection and fitness estimation. Estimation of fitness is performed by sequencing the genotype pool before and after selection and inferring changes of genotype frequencies in the course of selection. In DMS, it may be challenging or impossible to explicitly define the fitness of a mutated organism, so other quantitative traits are used as a proxy of fitness, such as protein stability, ligand binding affinity, fluorescence intensity, *etc.*

Usually, the mutant library for DMS contains all genotypes which can be obtained by a single mutation in a wild-type sequence (Figure 2.21a). Such data describe a small neighborhood of the fitness landscape and don't provide information on epistasis between these mutations or how the revealed fitness landscape shape affects the evolution of the sequence. One way to link single-mutation DMS data to evolutionary data is to compare mutational scans of two diverged homologous sequences (Figure 2.21bc) (Doud, Ashenberg, and Bloom 2015; Haddox et al. 2018; Chan et al. 2017; Lee et al. 2018). By comparing the SPFLs for the same site on different genetic backgrounds we can conclude to which extent the SPFLs are conserved among the homologs and whether the changes of alleles preferences between species are the response to environmental changes or is mediated by epistatic interactions (Chan et al. 2017; Lee et al. 2018). Such experiments show that the changes of the favorable allele generally resemble substitutions patterns (Doud, Ashenberg, and Bloom 2015), but not always; even strong SPFL shifts aren't necessarily followed by allelic replacements (Haddox et al. 2018).



**Figure 2.21. Example of using deep mutational scanning to detect evolutionary changes of allele preferences.** (a) Hypothetical DMS of a 25 amino acid long protein sequence: the colors show the value of some functional score for any possible genotype differing from the initial sequence by no more than one mutation (Fowler and Fields 2014). (b-c) Changes of SPFL inferred by DMS of the envelope protein of two HIV strains (BG505 and BF520) (Haddox et al. 2018). (b) The distribution of SPFL shifts between pairs of sequences, measured as RMSD between SPFLs of the same position in these two sequences corrected for the experimental noise. The distribution of  $RMSD_{corrected}$  in envelope protein of BG505 and BF520 (orange) is biased as compared to the randomized expectation (blue). The distribution of  $RMSD_{corrected}$  between the non-homologous sequences (here, influenza hemagglutinin protein HA) is shown in green. (c) The examples of sites with significantly shifted SPFLs. Logos show amino acid preference in two strains, black letters indicate the wild-type allele.

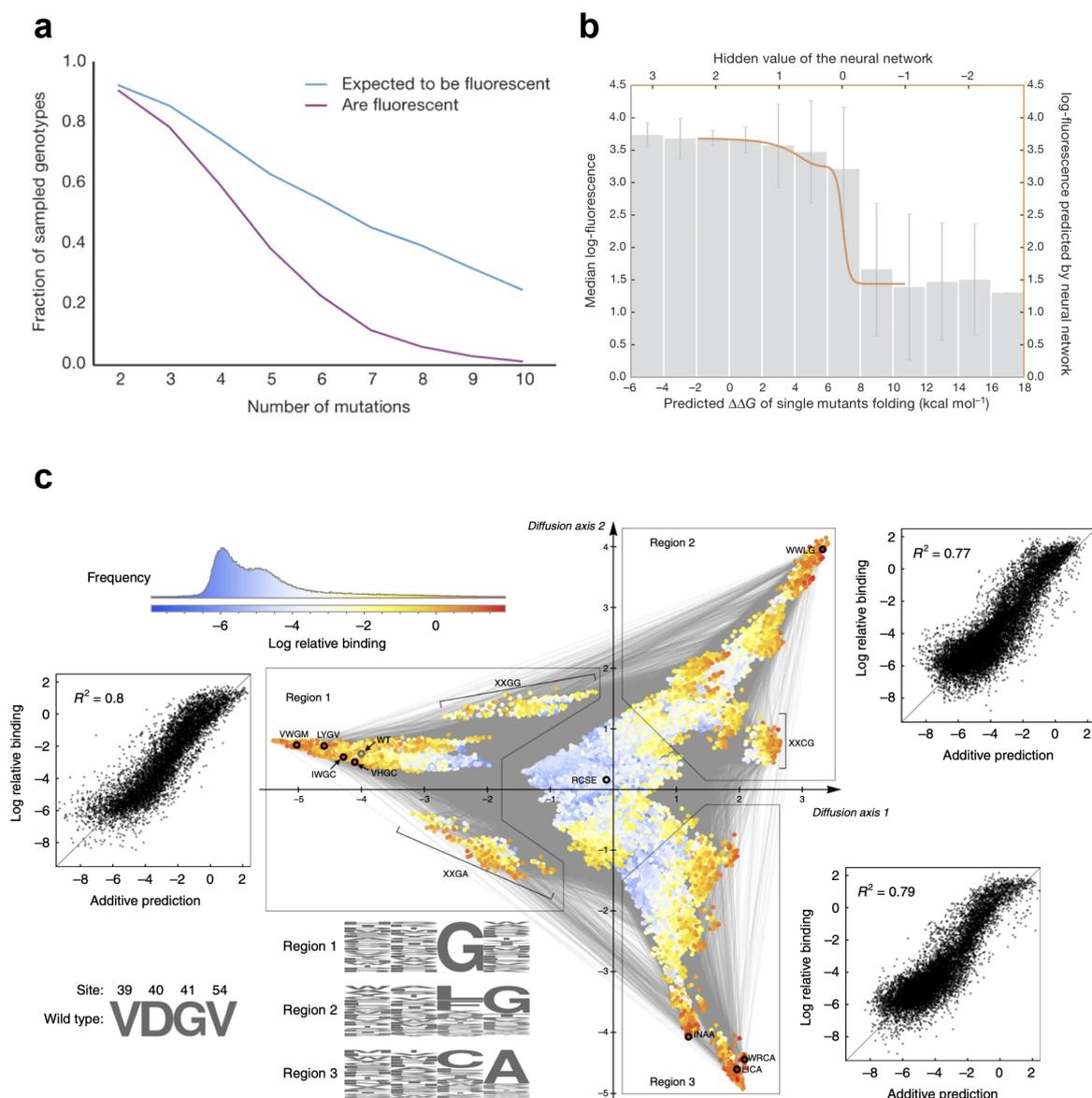
## The complexity of the empirical landscapes

DMS experiments may cover not only single-position mutants but also genotypes carrying combinations of two or more mutations as compared to the wild-type genotype. By limiting the number of considered sites and/or alleles in these sites, it's

possible to measure all possible combinations of the considered alleles — the full landscape on the space of these sites.

The empirical landscapes are highly epistatic. Multiple studies show that a large part of detected epistatic interactions can be attributed to unidimensional epistasis, with the additive fitness potential function representing the influence of mutations on some sequence trait (*e.g.* protein stability) (Figure 2.22a,b) (Sarkisyan et al. 2016; Jacquier et al. 2013; Kryazhimskiy et al. 2014; Otwinowski, McCandlish, and Plotkin 2018; Diss and Lehner 2018). The shape of the unidimensional epistasis informs on the evolvability of the protein: negative epistasis detected in some empirical landscapes prevents accumulation of combinations of deleterious mutations, reducing mutational load and constraining evolution as compared to the non-epistatic case (Figure 2.22a) (M. Kimura and Maruyama 1966; Otwinowski, McCandlish, and Plotkin 2018; Sarkisyan et al. 2016).

However, not all patterns of fitness variability can be reduced to unidimensional epistasis. Empirical landscapes reveal the presence of pairwise and high-order interactions that configure the rugged structure of protein landscapes and restrict the accessibility of evolutionary paths (Wu et al. 2016; Zhou and McCandlish 2020; Sailer and Harms 2017b; Lunzer, Golding, and Dean 2010). The pairwise epistasis is shown to originate from physical interaction between sites (Diss and Lehner 2018; Podgornaia 2014; Rollins et al. 2019; Stiffler et al. 2020). The landscapes can combine multi- and unidimensional effects: for example, reciprocal sign epistasis is shown to create distinct isolated fitness peaks, while unidimensional epistasis shapes the generally smooth surface of the peaks (Figure 2.22c) (Zhou and McCandlish 2020). Different regimes of epistasis may affect patterns of evolution on micro- and macroscale differently, decreasing the predictability of evolution.



**Figure 2.22. Uni- and multidimensional epistasis in empirical fitness**

**landscapes.** (a-b) Unidimensional epistasis between deleterious mutations shapes the fitness landscape of GFP (Sarkisyan et al. 2016). (a) The fraction of fit genotypes carrying multiple deleterious mutations (purple) is less than expected under non-epistatic expectations (blue). (b) Protein stability (measured as  $\Delta\Delta G$ ) is shown to be a good proxy of the fitness potential in GFP. (c) Reciprocal sign epistasis shapes the rugged landscape of four sites of GB1 protein (Zhou and McCandlish 2020; Wu et al. 2016). Gray lines connect genotypes differing by one mutation. The landscape is visualized with the dimensionality reduction method (McCandlish 2011). The corners of the triangle (regions 1-3) represent three adaptive peaks. Within peaks, epistasis is mostly unidimensional and can be approximated with a sigmoid function.

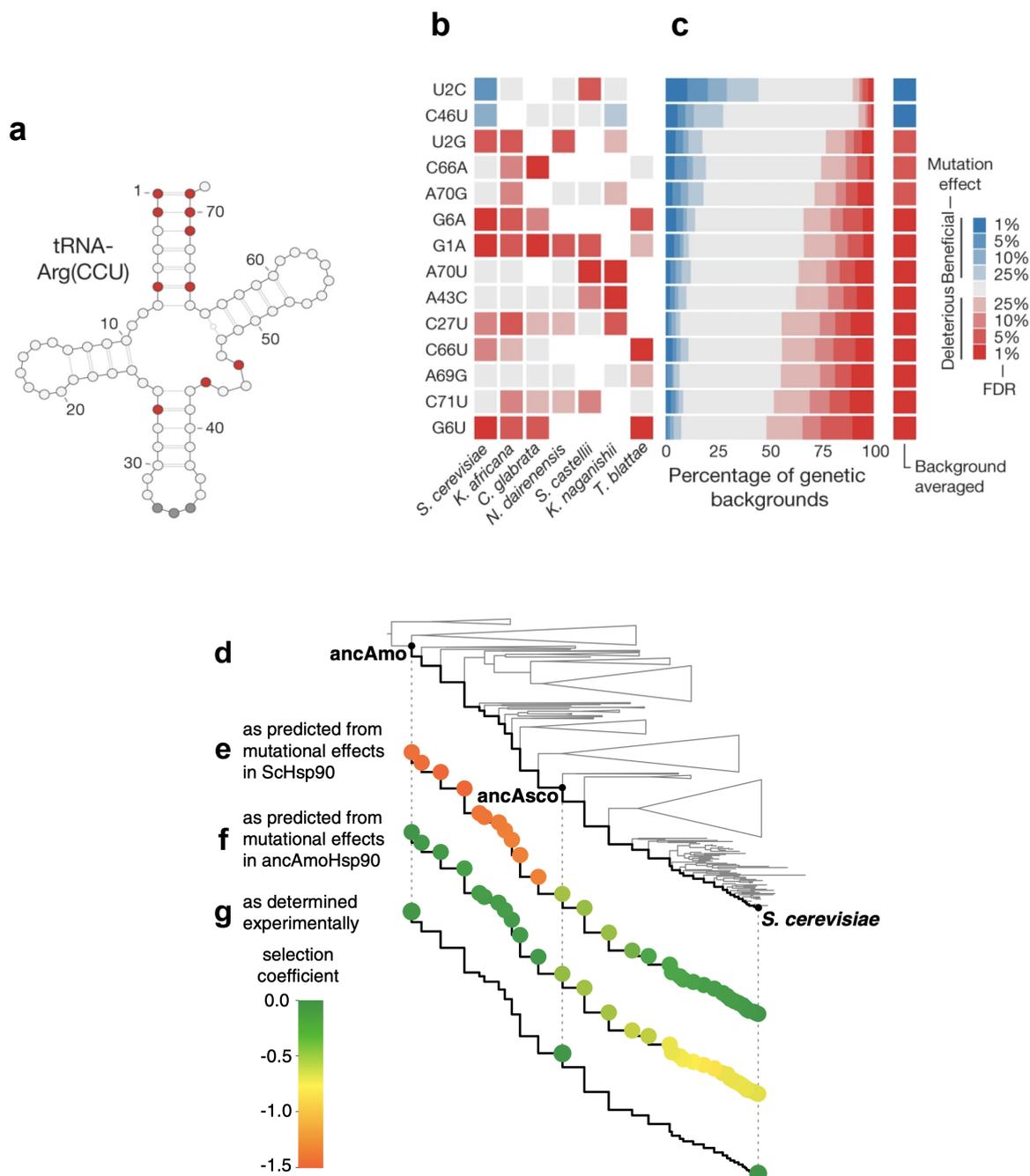
## **Empirical inference of historical evolution**

As was discussed above, only a small part of the full fitness landscape is actually engaged in the evolution of real populations. Although DMS experiments reveal abundant epistasis between sampled mutations, they don't provide insight on whether epistatic selection influences the fixation of mutations in the course of adaptation. The complexity of the historical evolutionary paths can be addressed by combining experimental approaches with comparative genomics methods to reconstruct the substitutions which occurred in the evolution of the examined sequence (Poelwijk et al. 2007; de Visser and Krug 2014; Bloom 2014).

One way to construct a set of potentially interesting genotypes is to focus on the alleles present in the orthologous sequences of several related species, and the combinations of such alleles. Experiments on measuring the fitness of the corresponding genotypes reveal epistasis-driven incompatibilities between the diverged genotypes: the evolutionary pathways between highly fit orthologous sequences are not neutral, but rather shaped by abundant pairwise and high-order epistasis (Figure 2.23a-c) (Pokusaeva et al. 2019; Domingo, Diss, and Lehner 2018; Poelwijk, Socolich, and Ranganathan, n.d.). Additionally, phylogenetic methods can be used to reconstruct the evolution of the sequence. In this case, it's possible to directly trace the adaptive paths by measuring the fitness of the ancestral and derived alleles in the context of present and reconstructed genotypes (Figure 2.23d) (Starr et al. 2018; Bridgham, Ortlund, and Thornton 2009; Gong, Suchard, and Bloom 2013; A. M. Phillips et al. 2021; Pillai et al. 2020).

Such studies show that the historical substitutions are not independent. The evolutionary paths are to large extent constrained by epistatic selection: many of the observed substitutions are not generally advantageous, but become such by preceding permissive substitutions at epistatically interacting sites (Figure 2.23e,g) (Starr et al. 2018; Bloom, Gong, and Baltimore 2010; Natarajan et al. 2016; Gong, Suchard, and Bloom 2013). The fixed allele then becomes entrenched due to co-evolution in the epistatically interacting sites. Entrenchment makes the currently present allele more fit, and the reversion to the ancestral variant more deleterious with time (Figure 2.23f,g) (Starr and Thornton 2016; Bridgham, Ortlund, and Thornton 2009; Wu et al. 2020).

Such retrospective analysis allows concluding that epistasis carves long and curved evolutionary pathways with strict constraints on the accessible order of substitutions. Permissive substitutions open new adaptive paths, previously hidden due to the ruggedness of the landscape, while the latter co-adaptation entrenches the newly fixed alleles, making the evolution non-reversible (Bridgham, Ortlund, and Thornton 2009; de Visser and Krug 2014).



**Figure 2.23. Epistatic constraints shaping the historical adaptive paths.** (a-c) epistatic interactions between mutations fixed in the evolution of the yeast arginine-CCU tRNA (Domingo, Diss, and Lehner 2018). (a) Secondary structure of the tRNA, with positions differing between the orthologous sequences shown in red. (b) Fitness effects of the corresponding mutations in the context of tRNA of different species. (c) Fitness effects of mutations across all measured genetic backgrounds. (d-g) Epistasis in the evolution of Hsp90 (Starr et al. 2018). (d) The phylogeny of the Hsp90; the reconstructed ancestral sequences of the common ancestor of Ascomycota and Amorphea are shown in black. (e-f) Fitness of the ancestral and intermediate genotypes as predicted from fitness effects of individual mutations on the background of the extant genotype (e) and the ancestral genotype (f). (g) Experimental measurements show that ancestral genotypes are as fit as the extant genotype, demonstrating the effects of permissive mutations and entrenchment.

## Phylogenetic evidence of epistasis

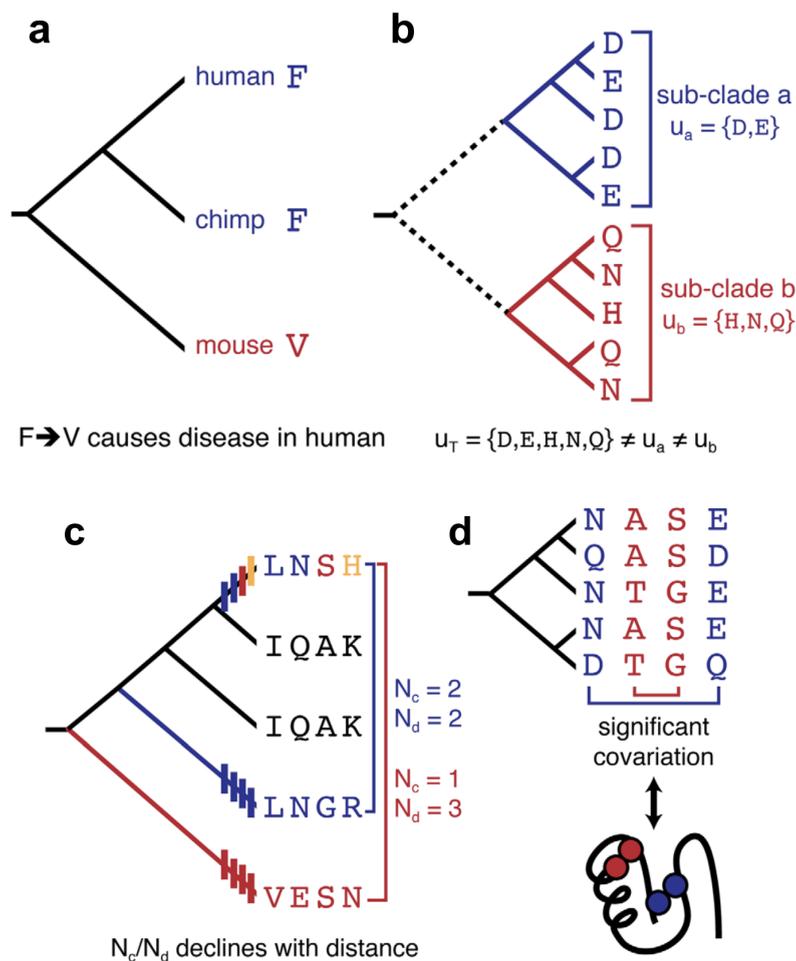
### Compensated pathogenic deviations

In the human population, some individuals carry strongly deleterious genomic variants shown to be causal on known diseases. Pathogenic variants may segregate in the human population, but can't achieve high frequency due to the associated fitness loss. A striking example of compensatory evolution is the fixation of disease-causing variants in the related species (compensated pathogenic deviations, or CPD) (Figure 2.24a) (A. S. Kondrashov, Sunyaev, and Kondrashov 2002; Kulathinal, Bettencourt, and Hartl 2004; Jordan et al. 2015). In this case, pathogenic allele becomes neutral due to substitutions in other genomic sites — usually, a single permissive substitution in the same gene is shown to be sufficient for the compensation (Jordan et al. 2015; Kern and Kondrashov 2004). In terms of fitness landscapes, compensation of deleterious alleles comprises evolutionary trajectories that go along the fitness ridges, resulting in the accumulation of genetic incompatibilities between the diverged genotypes (A. S. Kondrashov, Sunyaev, and Kondrashov 2002).

### Patterns of divergent and convergent evolution

While comparing homologous protein sequences from distant species, it's possible to infer changes of alleles' fitness between these species. If SPFL of a specific genomic site changed in the course of divergence of the lineages, distant sub-clades of the corresponding phylogenetic tree are expected to be enriched by different sets of alleles in this site (Figure 2.24b), as compared to the accumulation of amino acid changes in the neutral sites (Starr and Thornton 2016; Bazykin 2015). The dynamics of accumulation of genetic differences along the phylogenies is shown to be inconsistent with evolution under constant selection, but explainable by the rugged structure of the underlying fitness landscape (Povolotskaya and Kondrashov 2010; A. S. Kondrashov et al. 2010; Breen et al. 2012; McCandlish et al. 2013; Usmanova et al. 2015; Biswas et al. 2019). Complex networks of epistatic interactions, including compensatory and sign epistasis, elongate the adaptive evolution and slow the rate of fitness gain by constantly pushing the limits of the species divergence: epistasis constrains the number of accessible evolutionary paths at any given time point. Allelic substitutions in one genomic site

change the allelic preferences at other sites, opening new adaptive paths so that evolutionary trajectories go through curved fitness ridges (Povolotskaya and Kondrashov 2010; D. A. Kondrashov and Kondrashov 2015).



**Figure 2.24. Patterns of between-species variation evident of epistasis** (Starr and Thornton 2016). **(a)** Compensation of a disease-associated variant, **(b)** changes of the alleles usage between clades, **(c)** higher rate of convergent evolution in closely related species, **(d)** correlated evolution of physically interacting sites.

The divergence of alleles preferences along the phylogeny is also associated with specific dynamics of reversions (*i.e.* substitutions restoring the ancestral state) and convergent substitutions (*i.e.* recurrent substitution of the same allele in separate lineages) (Figure 2.24c). The dynamics of the rate of reversions to the ancestral allele

after substitution is indicative of the epistasis-driven changes of fitness of the alleles in the considered site (Starr and Thornton 2016; Storz 2016; Zou and Zhang 2015). Due to co-adaptation of the epistatically interacting sites, the newly established variant becomes entrenched, while the fitness of the initially beneficial allele that occupied the genomic site before the replacement declines with time (McCandlish, Shah, and Plotkin 2016). Shifts of the allelic preferences, including the loss of the epistatic “memory” of the ancestral allele, lead to the increase of the ratio of the rates of convergent and divergent evolution with time (Goldstein et al. 2015; Povolotskaya and Kondrashov 2010; Naumenko, Kondrashov, and Bazykin 2012). Shared allelic constraints increase the repeatability of evolution between closely related species, resulting in the negative correlation between the rate of convergent or parallel substitutions and evolutionary distance (Starr and Thornton 2016; Klink and Bazykin 2017; Soylemez and Kondrashov 2012).

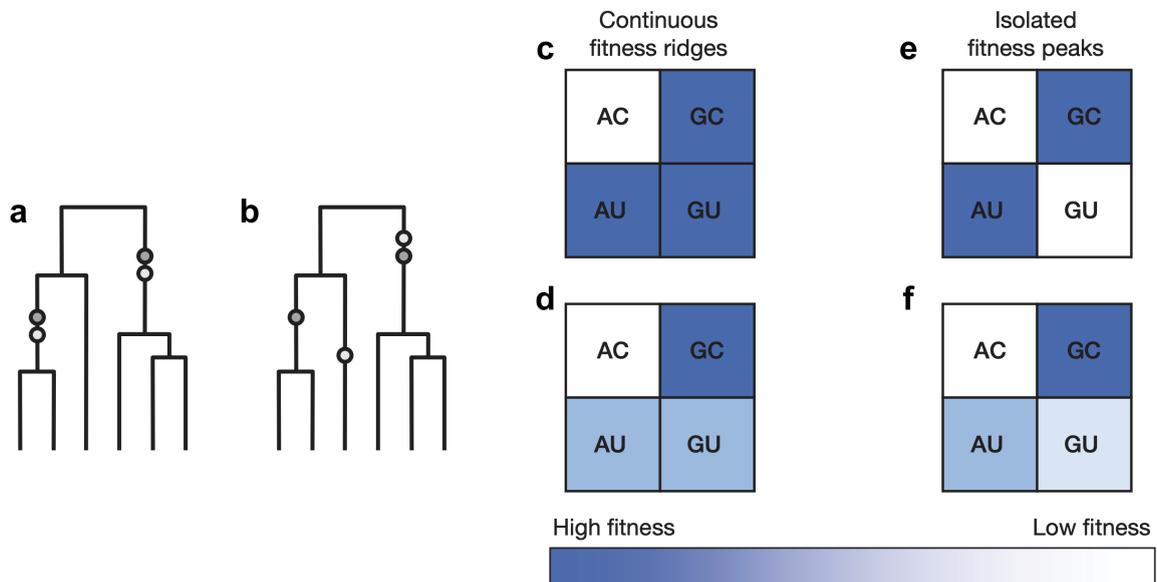
### **Correlated evolution of interacting sites**

Empirical studies of fitness landscapes reveal abundant pairwise epistasis between amino acid sites physically interacting in the protein structure (Rollins et al. 2019; Diss and Lehner 2018; Salinas and Ranganathan 2018; Olson, Wu, and Sun 2014). Comparative genomics studies on large phylogenies of distant species show that interacting sites tend to co-evolve, appearing as significant covariation between their evolution between species (Figure 2.24d) (Göbel et al. 1994; E. Neher 1994; Altschuh et al. 1987; Kamisetty and Ovchinnikov 2013). However, the inference of truly epistatic pairs of sites is impeded by abundant indirect correlations, which are hard to distinguish from the direct epistasis-driven correlations. This can be solved using direct-coupling analysis (DCA) — a group of statistical methods able to extract direct correlations and shown to be able to infer pairwise epistatic interactions based on thick between-species alignments or data on bacterial or viral divergence (Weigt et al. 2009; Morcos et al. 2011; J. P. Barton et al. 2016; Burger and van Nimwegen 2010; Puranen et al. 2018; Figliuzzi et al. 2016). Coupling analysis is successfully leveraged to reconstruct protein and RNA structures and protein-protein interactions based on epistatic constraints between physically interacting sites (Ovchinnikov, Kamisetty, and Baker

2014; Marks et al. 2011; Morcos et al. 2011; Sjodt et al. 2018; De Leonardis et al. 2015; Ovchinnikov et al. 2015, 2017).

### **Phylogenetic clustering of interacting sites**

Co-evolution of epistatically interacting sites between species means that substitutions in these sites tend to occur simultaneously or within a short time. By reconstructing the evolutionary history of the sequence, it's possible to analyze joined phylogenetic distribution of substitutions in presumed interacting pairs of sites. Under epistasis, a substitution occurring in one site can drive co-adaptation in another site, causing the temporal clustering of allelic replacements along the phylogeny (Figure 2.25a,b) (Shapiro et al. 2006; Neverov et al. 2021, 2014; Bazykin 2015). The extent of clustering, which can be estimated by the lifespan of the intermediate states, characterizes the mode and strength of epistatic selection: whether the evolutionary path consisting of two subsequent substitutions goes along the adaptive ridge or crosses an adaptive valley (Figure 2.25c-f) (Meer et al. 2010; Gong, Suchard, and Bloom 2013; Nasrallah and Huelsenbeck 2013).



**Figure 2.25. Phylogenetic clustering of epistatically interacting substitutions.** (a-b) Distribution of compensatory (a) and independent (b) substitutions along the phylogeny (Shapiro et al. 2006). (c-f) Possible landscape models underlying the fast two-step transition from AU pair to GC in mitochondrial tRNAs (Meer et al. 2010): (c) flat fitness ridge, (d) ascending fitness ridge, (e) equal fitness peaks isolated by fitness valleys, (f) fitness peaks of different height isolated by intermediate states of different fitness.

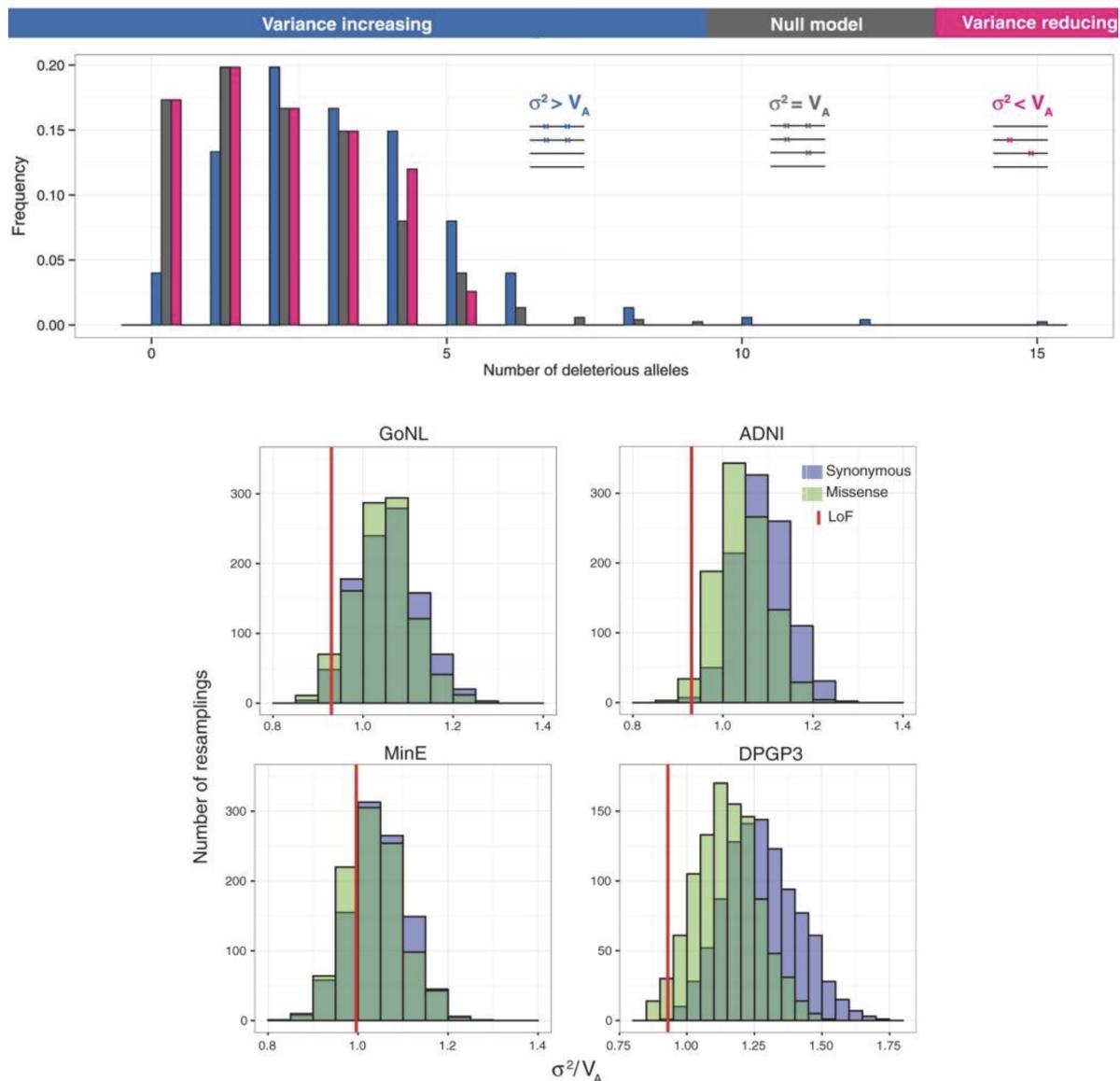
## Epistasis in within-population variation

### Statistical epistasis

Despite the abundant epistasis between substitutions observed on the macroevolutionary level, it's hard to detect between polymorphisms segregating within natural populations. Epistatic component of the variation of some complex trait can be the cause of missing heritability of the trait — however, there is no evidence of epistasis having large contribution to the trait variation, and to what extent missing heritability in natural populations is indeed explained by epistasis remains controversial (Hayman and Mather 1955; Falconer and Falconer 1989; Cheverud and Routman 1995; Hivert et al. 2021; Zuk et al. 2012; Visscher et al. 2007; Sackton and Hartl 2016). In genome-wide association studies (GWAS) phenotypic variation within populations is explained by genome content. The effects of polymorphisms associated with human polygenic phenotypes are shown to be generally additive: phenotypic variance attributed to epistasis between polymorphisms segregating within a population (statistical epistasis) is negligible (Cheverud and Routman 1995; James F. Crow 2010; Mackay and Moore 2014; William G. Hill, Goddard, and Visscher 2008). Complex linkage effects can also result in artificial signals of epistasis, impeding its inference in GWAS (Wood et al. 2014; Hemani et al. 2021).

### Epistasis-driven linkage disequilibrium

Epistasis can maintain favorable combinations of alleles at interacting sites, increasing linkage disequilibrium (LD) between them (Ronald A. Fisher 1930; Lewontin and Kojima 1960; N. H. Barton 2010; Takahasi and Tajima 2005; Kouyos, Silander, and Bonhoeffer 2007; Pedruzzi, Barlukova, and Rouzine 2018; Boyrie et al. 2021). In the absence of population structure, genomic admixtures or recent changes of population size, unlinked polymorphisms (*e.g.* located on different chromosomes or separated by large genomic distances within the same chromosome) are expected to segregate independently (Nei and Li 1973; Schaper et al. 2012; Lewontin and Kojima 1960; Rohlf, Swanson, and Weir 2010).



**Figure 2.26. Repulsion LD between loss-of-function polymorphisms in human and fruit fly populations (Sohail et al. 2017).** (a) The number of deleterious alleles per genotype in the absence of epistasis is expected to be Poisson-distributed so that its variance ( $\sigma^2$ ) is equal to the mean (additive variance  $V_A$ ) (gray). Under antagonistic epistasis, deleterious alleles are overdispersed (blue), while under synergistic epistasis they are underdispersed (red). (b) Underdispersion of LoF alleles (red) in human (Netherlands GoNL, European ancestry ADNI and Dutch MinE datasets) and fruit fly (Zambian DPGP3 dataset) populations.

Epistasis between distant polymorphisms may keep them long-range linkage disequilibrium (LRLD), which can be detected by analyzing population genomic datasets (Koch, Ristroph, and Kirkpatrick 2013; L. Park 2019). An example of unidimensional epistasis creating LRLD is synergistic epistasis between loss-of-function (LoF)

polymorphisms within populations of *H. sapiens* and *D. melanogaster* (Sohail et al. 2017). The decreased variance of the number of LoF alleles per genome as compared to the Poisson distribution indicates their repulsion, or negative, LD: the probability of a genotype to contain several LoF alleles is reduced to what is expected if they segregate independently (Figure 2.26).

In sexual populations, recombination competes with epistasis, disrupting coupling LD between distant interacting sites (Ronald A. Fisher 1930; H. J. Muller 1932; Franklin and Lewontin 1970; R. A. Neher and Shraiman 2009; Pedruzzi and Rouzine 2019). Nevertheless, within a single gene, physical proximity may suffice to limit recombination, so sets of coadapted variants may evolve (Lewontin and Kojima 1960; T. Dobzhansky 1950). However, population structure and complex effects of genetic drift and linkage impedes detection of short-range epistatic interactions (Ragsdale 2021; Good 2020). Despite these difficulties, recent studies describe repulsion between nonsynonymous and LoF polymorphisms within populations, which may be explained by negative epistasis acting on them (Garcia and Lohmueller 2021; Sandler, Wright, and Agrawal 2021). The opposite phenomenon of coupling LD between derived nonsynonymous variants was detected in bacterial populations (Arnold et al. 2020). The effect was largely restricted to regions of high genetic diversity, presumably generated by ongoing positive or balancing selection.

## Chapter 3: Complex fitness landscape shapes variation in a hyperpolymorphic species

It is natural to assume that patterns of genetic variation in hyperpolymorphic species can reveal large-scale properties of the fitness landscape that are hard to detect by studying species with ordinary levels of genetic variation. Here, we study such patterns in a fungus *Schizophyllum commune*, the most polymorphic species known. Throughout the genome, short-range linkage disequilibrium caused by attraction of rare alleles is higher between pairs of nonsynonymous than of synonymous sites. This effect is more pronounced if both sites are located within the same gene, especially if a large fraction of the gene is covered by haploblocks, genome segments where the gene pool consists of two highly divergent haplotypes, which is a signature of balancing selection.

Haploblocks are usually shorter than 1000 nucleotides, and collectively cover about 10% of the *S. commune* genome. LD tends to be substantially higher for pairs of nonsynonymous sites encoding amino acids that interact within the protein. There is a substantial correlation between LDs at the same pairs of nonsynonymous sites in the USA and the Russian populations. These patterns indicate that selection in *S. commune* involves positive epistasis due to compensatory interactions between nonsynonymous alleles. When less polymorphic species are studied, analogous patterns can be detected only through interspecific comparisons.

## Introduction

Alleles do not affect fitness and other phenotypic traits independently and, instead, often engage in epistatic interactions (Maynard Smith 1970; Wright 1932; Fenster, Galloway, and Chao 1997; John H. Gillespie 1994; Povolotskaya and Kondrashov 2010; de Visser, Cooper, and Elena 2011; McCandlish et al. 2013; de Visser and Krug 2014; Good and Desai 2015; Kryazhimskiy et al. 2011). Epistasis is pervasive at the scale of between-species differences, where it is saliently manifested by Dobzhansky-Muller incompatibilities and results in low fitness of interspecific hybrids (T. Dobzhansky 1936; Orr 1995; A. S. Kondrashov, Sunyaev, and Kondrashov 2002). By contrast, at the scale of within-population variation, the importance of epistasis remains controversial (Ronald A. Fisher 1930; H. J. Muller 1932; Franklin and Lewontin 1970; R. A. Neher and Shraiman 2009; Sackton and Hartl 2016; James F. Crow 2010; Mäki-Tanila and Hill 2014; William G. Hill, Goddard, and Visscher 2008; Hivert et al. 2021). This may look like a paradox, because such variation provides an opportunity to detect epistasis through linkage disequilibrium (LD), non-random associations between alleles at different loci (Ronald A. Fisher 1930; H. J. Muller 1932; Franklin and Lewontin 1970; N. H. Barton 2010; Takahasi and Tajima 2005; Kouyos, Silander, and Bonhoeffer 2007; Pedruzzi, Barlukova, and Rouzine 2018; Boyrie et al. 2021). Indeed, epistatic selection generates LD which can be detected (M.-C. Wang et al. 2012; Beissinger et al. 2016; Zan, Forsberg, and Carlborg 2018; Garcia and Lohmueller 2021; Boyrie et al. 2021). Perhaps, the fitness landscape is complex macroscopically (between diverged species) but is more smooth microscopically (within populations) or, in other words, epistasis is genuinely more pronounced at a macroscopic scale (Ochs and Desai 2015). If so, studying epistasis in hyperpolymorphic populations, where differences between genotypes can be as high as those between genomes of species from different genera or even families, holds a great promise because variation within such a population can cover multiple fitness peaks or a sizeable chunk of a curved ridge of high fitness (Theodosius Dobzhansky 1937; Bateson 1909; H. Muller 1942; S. Gavrillets 1997; A. S. Kondrashov, Sunyaev, and Kondrashov 2002).

The basidiomycete fungus *Schizophyllum commune* possesses the highest genetic diversity among the studied eukaryotic species, with up to 20% of neutral sites differing between any two individuals within a population. High genetic diversity is at least partially caused by high mutation rates in *S. commune*: although the per-generation mutation rate measured in the laboratory conditions is not extremely high ( $2 \times 10^{-8}$  mutations per nucleotide per generation), the fungus is shown to accumulate mutations while the mycelium grows, so that the per-generation mutation rate in nature can be substantially higher (Baranova et al. 2015; Bezmenova et al. 2020). Another cause of such a high level of polymorphism may be large effective population size (Baranova et al. 2015). *S. commune* possesses relatively small genome (38.5 Mb, 11 chromosomes), with approximately half of the genome carrying protein-coding sequences (Ohm et al. 2010). Its genes have introns (on average four per gene), although relatively short (typically < 100 nt). It has more than 20,000 mating types, encoded by two mating-type loci (Kothe 1999).

Recent study on the distribution of the crossing-over events in F1 hybrids of a pair of individuals sampled from USA and Russia showed that they are more frequent in genomic regions where parental genotypes are similar to each other, including exons, where relatively low genetic diversity is maintained by negative selection (Seplyarskiy et al. 2014). *S. commune* has a haploid life stage and can be cultivated in laboratory conditions, making it a promising object for population genetics studies.

Here, we study the LD patterns in 55 complete genomes of *S. commune* from North America and Europe.

## Materials and methods

### ***S. commune* sampling, sequencing and assembly**

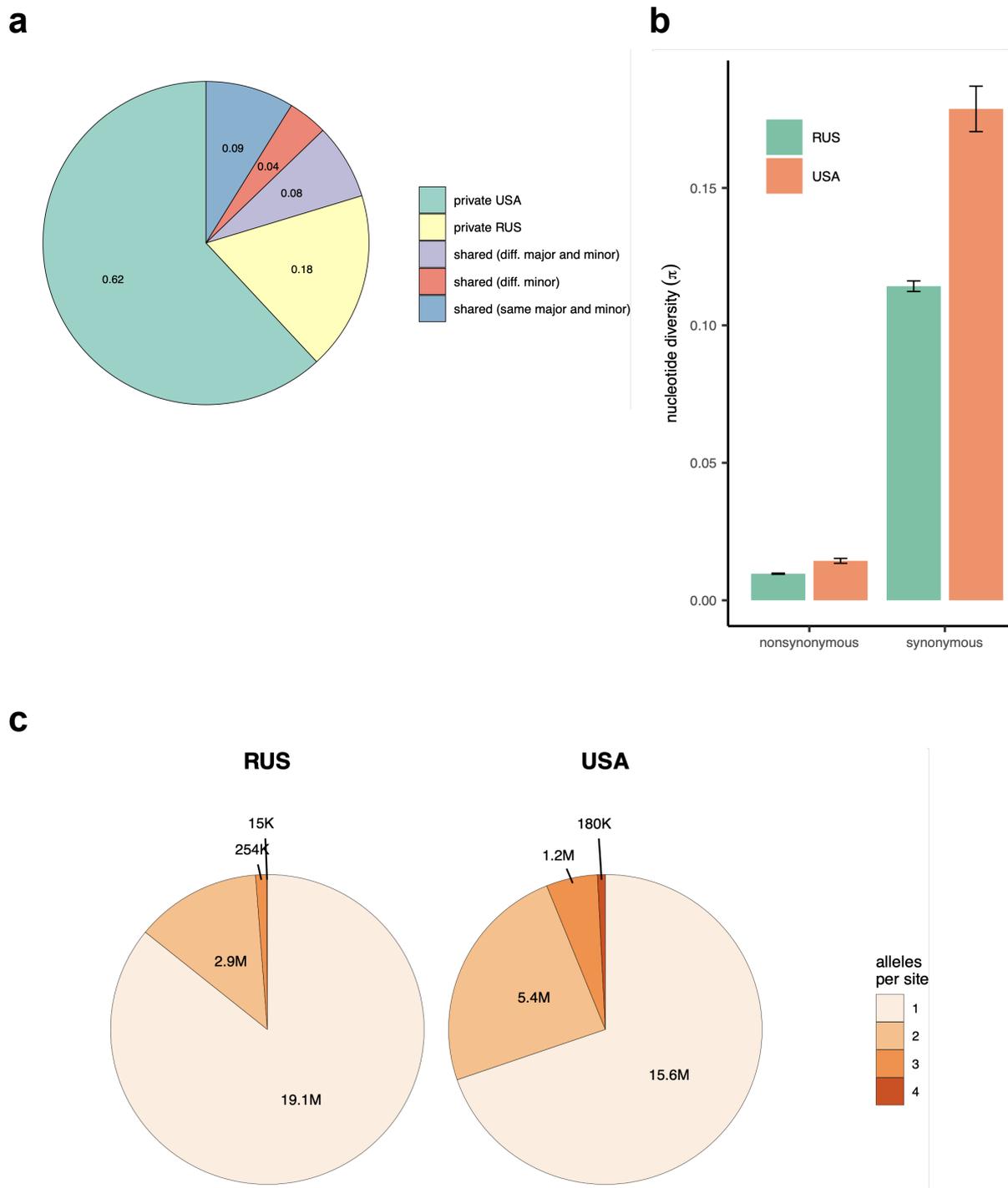
Haploid cultures of 24 isolates, each originated from a single haplospore, were obtained from fruit bodies collected in Ann Arbor, MI, USA by T. James and A. Kondrashov (17 samples) and in Moscow and Kostroma regions, Russia by A. Kondrashov, A. Baykalova and T. Neretina (7 samples) in 2009–2015. Specimen vouchers are stored in the White Sea Branch of Zoological Museum of Moscow State University (WS). To obtain isolates, wild fruit bodies were hung on the top lid of a 10 cm petri dish with agar medium. Petri dish was set at an angle of 60-70 degrees to the horizontal surface for 32 hours. A germinated spore was excised together with a square-shaped fragment (approximately 0.7x0.7 mm) of the medium from the maximally rarefied area of the obtained spore print under a stereomicroscope with 100x magnification. The obtained isolates were cultured in Petri dishes on 2% malt extract agar for a week. For storage, cultures were subcultured into 1.5 ml microcentrifuge tubes with 2% malt extract agar. To obtain sufficient biomass for DNA isolation, isolates were cultured in 20 ml 0.5% malt extract liquid medium in 50 ml microcentrifuge tubes in a horizontal position on a shaker at 100 rpm in daylight for 5 to 10 days. The tubes with the cultures were then centrifuged at 4000 rpm, and the supernatant was decanted. The resulting mycelium was lyophilized. DNA was extracted using Diamond DNA kit according to the manufacturer's recommendations.

DNA libraries were constructed using the NEBNext Ultra II DNA Library Prep Kit kit by New England Biolabs (NEB) and the NEBNext Multiplex Oligos for Illumina (Index Primers Set 1) by NEB following the manufacturer's protocol. The samples were amplified using 10 cycles of PCR. The constructed libraries were sequenced on Illumina NextSeq500 with paired-end read length of 151. The genomes were assembled *de novo* using SPAdes (v3.6.0) (Bankevich et al. 2012); possible contaminations were removed using *blobology* (Sujai Kumar et al. 2013). Average N50 was ~165kb for USA samples and ~70kb for Russian samples (assembly statistics shown in Table A1).

Together with the 30 samples sequenced previously (Baranova et al. 2015; Bezmenova et al. 2020), the obtained haploid genomes were aligned with TBA and *multiz* (Blanchette et al. 2004) and projected onto the reference scaffolds (Ohm et al. 2010). Ortholog sequences were extracted based on the reference genome annotation (Ohm et al. 2010) and realigned using *macse* codon-based aligner (Ranwez et al. 2011). Only the gap-free columns of the whole-genome alignment and the orthologs that were found in all 55 genomes were used for analysis. The total number of detected SNPs was 5.8 million for the USA population (82% of them biallelic) and 2.7 million for the Russian population (93% biallelic). 25% of the USA SNPs were shared with the Russian population (11% with the same major and minor alleles), and 53% of the Russian SNPs were shared with the USA population (23% with the same major and minor alleles, Figure 3.1).

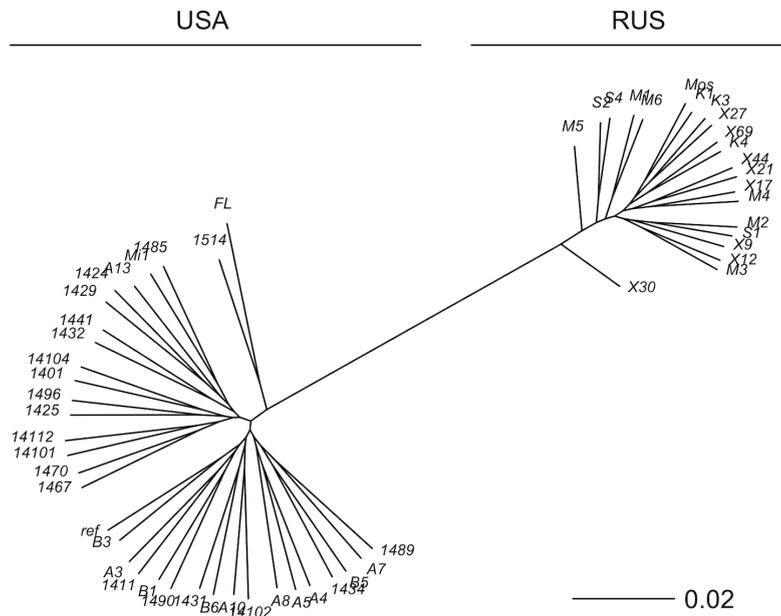
To represent the genomic distances between populations and the within-population structure, we reconstructed the whole-genome phylogeny of the sequenced genomes with RAxML (Stamatakis 2014) (Figure 3.2). Nucleotide diversity ( $\pi$ ) was estimated as the average frequency of pairwise nucleotide differences;  $\pi$  for different classes of sites is shown in Figure 3.1. Here, the number of synonymous and nonsynonymous sites was calculated using the method described in (Nei and Gojobori 1986). The value of  $\pi_{\text{syn}}$  calculated with this method was  $\sim 10\%$  lower than the value obtained using only the fourfold-degenerate sites.  $p_n/p_s$  for single genes was calculated as the ratio of the average number of synonymous and nonsynonymous differences between a pair of genomes divided by the number of corresponding sites, estimated the same way as  $\pi$ . Two samples from Florida (USA population; samples FL and s1514) represent a separate sub-population external to other USA samples (Figure 3.2,  $F_{\text{st}} = 0.11$ ), so they were excluded from the further analysis to minimize the possible effect of population structure.

Genome sequence data are deposited at DDBJ/ENA/GenBank under accession numbers JAGVRL000000000-JAGVSI000000000, BioProject PRJNA720428. Sequencing data are deposited at SRA with accession numbers SRR14467839-SRR14467862.



**Figure 3.1. Patterns of nucleotide diversity in *S. commune*.** (a) The fraction of private and shared biallelic SNPs. (b) Within-population nucleotide diversity at different classes of sites. (c) The number of monomorphic and polymorphic sites in the

multiple whole-genome alignments of *S. commune* genomes. The USA dataset consists of 34 sequenced genomes, and the Russian dataset consists of 21 samples.



**Figure 3.2. The reconstructed phylogeny of *S. commune*.** USA and Russian populations of *S. commune* are highly divergent while having almost no within-population structure. Genetic distance is measured in nucleotide differences, the phylogeny is reconstructed based on the multiple whole-genome alignment.  $\pi$  between populations is approximately 0.34,  $F_{st} = 0.58$ .

### Data on *H. sapiens* and *D. melanogaster* populations

We used polymorphism data from 1,296 phased human genomes from African and European super-populations sequenced as part of the 1000 Genomes project (1000 Genomes Project Consortium et al. 2015). If several individuals from the same family were sequenced, we included only one of them. As a *D. melanogaster* dataset, we used 197 haploid genomes from the Zambia population (Lack et al. 2015). Only autosomes were analyzed in both datasets.

## Estimation of LD

As a measure of linkage disequilibrium between two biallelic sites, we used  $r^2$ , calculated as follows:

$$r^2 = \frac{(p(AB) - p(A)p(B))^2}{p(A)(1-p(A))p(B)(1-p(B))},$$

where  $p(A)$  and  $p(B)$  are the minor allele frequencies at these sites and  $p(AB)$  is the frequency of the genotype carrying both minor alleles.

Multiallelic sites (4.9% of polymorphic sites in the USA population and 0.9% in the Russian populations, and singletons (sites with minor allele present only in one genotype) were excluded from the analysis.

## Haploblocks annotation

In order to annotate the haploblocks, we calculated LD along the *S. commune* genome in a sliding window of 250 nucleotides with a step of 20 nucleotides (only non-singleton SNPs are analyzed; the windows with less than 10 SNPs were excluded). Any continuous sequence of overlapping windows with LD ( $r^2$ ) larger than the threshold value was merged together in a haploblock. The LD threshold value was defined independently for each *S. commune* population as the heavy tail of the within-window LD distribution, as compared with the lognormal distribution with the same mean and variance as in the data.

## Estimation of LD between physically interacting amino acid sites

Of 16,319 annotated protein-coding genes of *S. commune* (Ohm et al. 2010) 9,941 were found in all 55 aligned genomes. We blasted the protein sequences of these orthologous groups against the PDB database of protein structures. About 52% of them (5,188) had a match (e-value threshold =  $1e-5$ ) amongst the proteins with the known structure. We realigned the sequences of *S. commune* protein and the matching PDB protein with clustal and calculated within-population LD and physical distance ( $\text{\AA}$ ) for each pair of aligned positions in the corresponding three-dimensional structure. A pair of amino acid

sites was considered physically adjacent if they were located within 10 Å from each other.

To compare LD between pairs of physically close and distant sites, we used the controlled permutation test: for each pair of physically close amino acid sites (within 10 Å) we sampled a pair of physically distant amino acids on the same nucleotide distance (measured in aa). Pairs of sites closer than 5 aa were excluded from the analysis.

To examine LD patterns within individual protein structures, we calculated contingency tables of pairs of SNPs being located in codons encoding physically close amino acids and having high LD (no less than 90% quantile for a given gene). Pairs of amino acid sites located closer than 30 aa or more distant than 100 aa from each other were excluded; genes with less than 5 pairs of physically close sites under high or low LD were also excluded. From these contingency tables, we calculated the odds ratio (OR) and chi-square test p-value for each gene. p-values were adjusted using BH correction (Benjamini and Hochberg 1995).

### **Simulations of epistasis**

To simulate evolution of populations with or without epistasis and balancing selection, we used an individual-based model implemented by *SLiM* (Haller and Messer 2019). Simulations are performed with diploid population size  $N=1000$  and recombination rate 0. To achieve the level of genetic diversity  $\pi$  similar to *S. commune*, mutation rate  $\mu$  is scaled as  $\mu=\pi/2N=5e-5$ , recurrent/reverse mutations at the same genetic site are allowed. The length of the simulated sequence is 100 bp. Each simulation starts with a monomorphic population and proceeds for  $100N$  generations. For calculations of synonymous and nonsynonymous LD, random 100 haploid genotypes are sampled from the population. Only SNPs with minor allele frequency  $> 5\%$  in the sample are analyzed.

We model two types of sites, depending on whether mutations in them are neutral (with selection coefficient  $s_{\text{syn}} = 0$ ) or weakly deleterious ( $s_{\text{nonsyn}} \leq 0$ ), representing synonymous and nonsynonymous sites correspondingly. There are twice as many

nonsynonymous as synonymous sites. Under the non-epistatic model,  $s$  is independent of the genetic background. We assume  $s_{\text{nonsyn}} = -0.01$  with the dominance coefficient  $h$  of 0.5.

Under the pairwise positive epistasis model, we assume that a mutation at one nonsynonymous site can be partially or fully compensated by a mutation at another site. In this model, all nonsynonymous sites are split into pairs. Each mutation of a pair individually occurring within a genotype is assumed to be deleterious, with selection coefficient  $s_{\text{nonsyn}} = -0.01$ ; however, the fitness of the double mutant is larger than expected under the additive (non-epistatic) model. We use multiple epistasis models, which vary on the epistasis strength and landscape shape.

In the NFDS model of balancing selection, a single mutation at a random position is subjected to frequency-dependent selection (so that it is positively selected at frequencies below 0.5, and negatively selected at frequencies above 0.5). In the AOD model, mutations in 10 random positions are fully recessive ( $h=0$ ) and weakly deleterious ( $s=-0.0025$ ).

To simulate evolution of populations with different levels of genetic diversity under epistasis, we use *FFPopSim* (Zanini and Neher 2012) (the simulation results obtained with *FFPopSim* and *SLiM* were checked to be similar for  $\pi=0.2$  as in *S. commune*, but *FFPopSim* calculation time was substantially shorter; unfortunately, it doesn't allow to simulate evolution of diploid population and therefore couldn't be used to simulate overdominance). To achieve different levels of genetic diversity  $\pi$ , mutation rate  $\mu$  is scaled as  $\mu = \pi/2N$  (we checked that this approach gives the same results as scaling of  $N$  instead of  $\mu$ , as long as we scale  $s$  and recombination rate  $\rho$  to maintain  $Ns$  and  $\rho N$  constant). The calculations are performed the same way as in *SLiM*, but In this case, we use haploid population size  $N=2000$ , population-scaled recombination rate 0.01 and the simulated sequence length of 300 nucleotides.

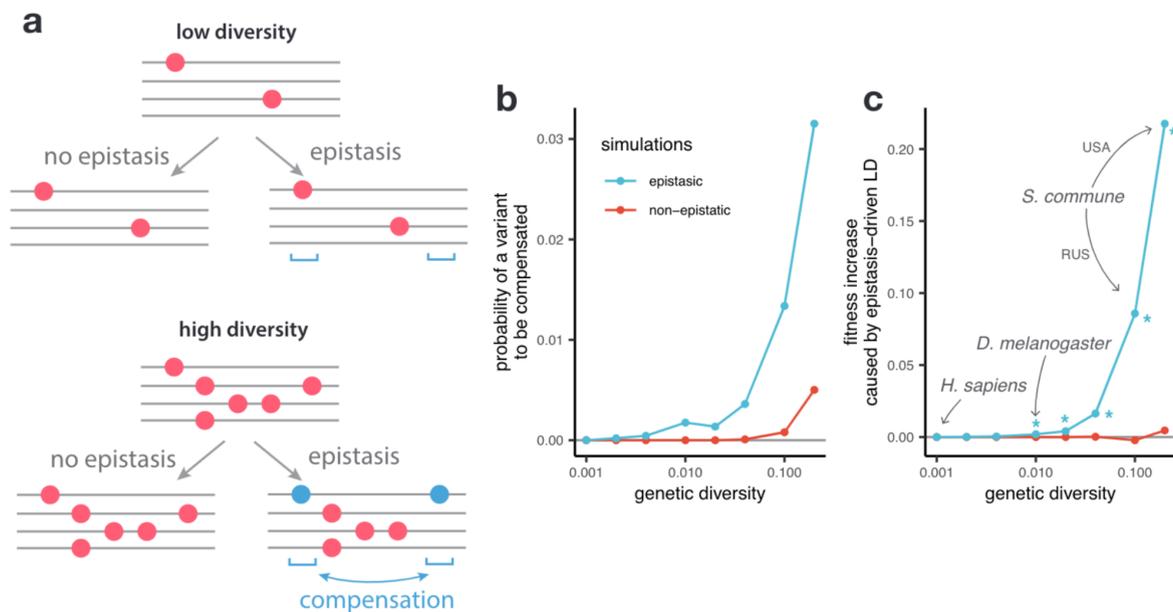
## Results

### **Epistatic selection is more efficient in genetically diverse populations**

Genetic interactions affect operation of selection only affect patterns of variation if sufficient variation is maintained. The potency of any kind of selection acting on the variants segregating within a population increases with the amount of variation. If the neutral level of genetic variation is low, the efficiency of selection will also be low since it won't be able to eliminate deleterious variants and promote beneficial variants if they are absent. For epistatic selection, however, this increase is expected to be faster than linear, because it depends on the number of possible allele combinations. In a highly polymorphic population, a particular allele is more likely to co-occur in the same haplotype with an interacting, *e.g.*, compensatory, allele, which should increase the impact of epistasis on linkage disequilibrium (Figure 3.3a).

To illustrate this point, we modelled the evolution of a genome region in the presence and in the absence of positive epistasis in a panmictic population. We assumed that all mutations at a set of sites are individually deleterious, and that all these sites are involved in pairwise positive (*i.e.*, antagonistic) sign epistasis; specifically, each deleterious mutation can be fully compensated by another mutation at exactly one site elsewhere in the genome, which is also deleterious when present alone. In the non-epistatic simulations, the effects of mutations were independent; however, at the end of the simulation we randomly assigned the “interacting” pairs of sites to account for the random coincidence of deleterious alleles. We found that in this model a higher polymorphism increases the probability that a deleterious mutation is compensated before being eliminated by selection (Figure 3.3b). This probability increases with genetic diversity even for the non-epistatic simulations, because increased diversity elevates the likelihood of randomly encountering a compensating allele in the same haplotype. For epistatic simulations, however, this increase is more radical, reflecting the effect of epistatic selection favoring compensated haplotypes. Despite the fact that a fraction of deleterious alleles are compensated in a fraction of genotypes, the average population fitness is higher in the non-epistatic simulations, consistent with positive epistasis increasing mutational load (Brian Charlesworth 1990).

After the mutation-selection equilibrium was reached, we measured the strength of epistatic selection between all segregating polymorphisms, asking to what extent the mutational load is reduced by epistasis maintaining combinations of compensatory mutations. As shown in Figure 3.3c, the ability of epistatic selection to reduce the mutation load (i. e., to increase the mean fitness) strongly depends on  $\pi$ . In less variable populations ( $\pi < 0.01$ ), epistasis is practically inefficient and doesn't affect LD (Wilcoxon test p-values  $> 0.33$ ); this is because the probability of occurrence of the favorable combination of alleles in the population for selection to act upon is low. In more diverse populations, however, such combinations may arise and be favored by epistatic selection, which increases LD between them (Wilcoxon test p-value  $< 0.01$  for  $\pi \geq 0.01$ ).

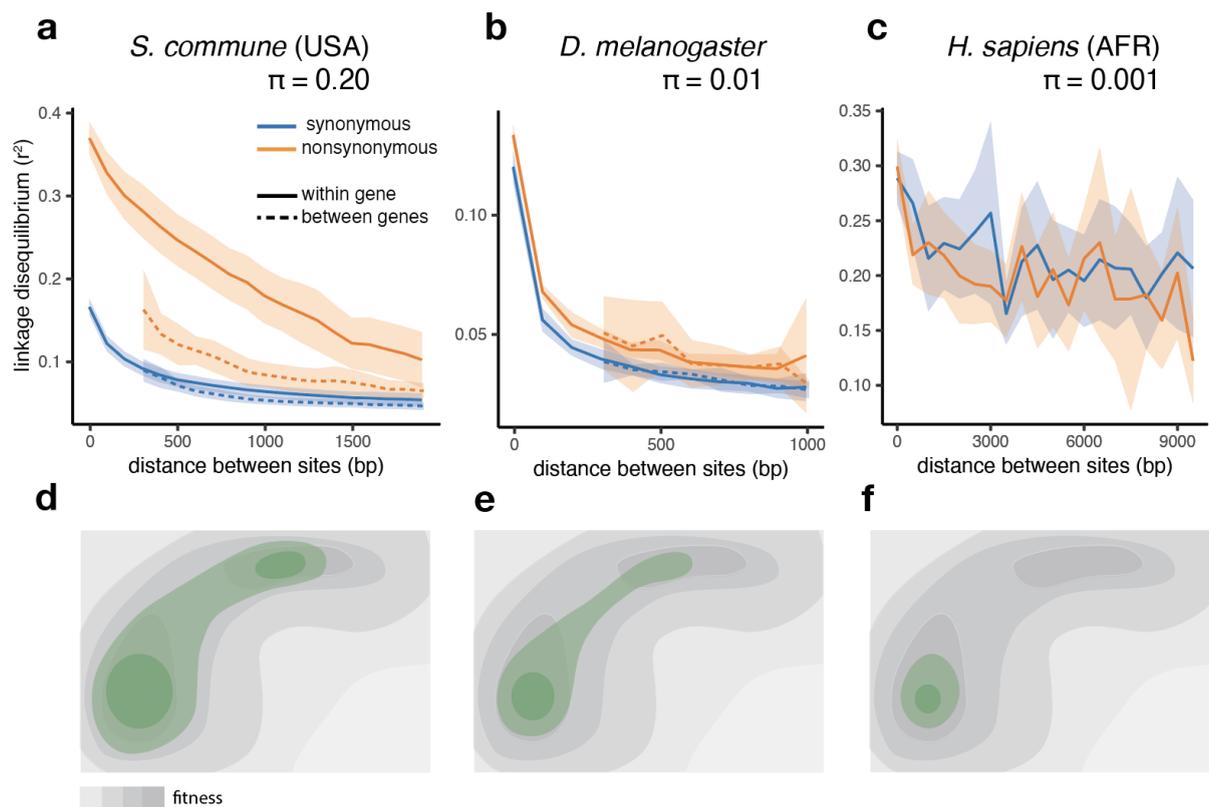


**Figure 3.3. The efficiency of epistasis in populations with different levels of genetic diversity.** (a) Under low genetic diversity, deleterious mutations (red dots) are unlikely to be compensated. If genetic diversity is high, epistatic selection maintains LD between SNPs in interacting sites (blue dots). (b) The probability that a deleterious variant is compensated by another variant within the same individual at the end of the simulation. (c) Increase in mean fitness of a population caused by epistatic selection maintaining LD between favorable allele combinations. The fitness is plotted relative to that of a population consisting of individuals with uncorrelated alleles at different sites, obtained by permuting alleles among individuals. The efficiency of epistatic selection in maintaining linkage is much higher in genetically variable populations. Asterisks in (c) indicate significant deviation from 0 (Wilcoxon paired test p-value  $< 0.01$ ). Each simulation was repeated for 100-10,000 times depending on genetic diversity.

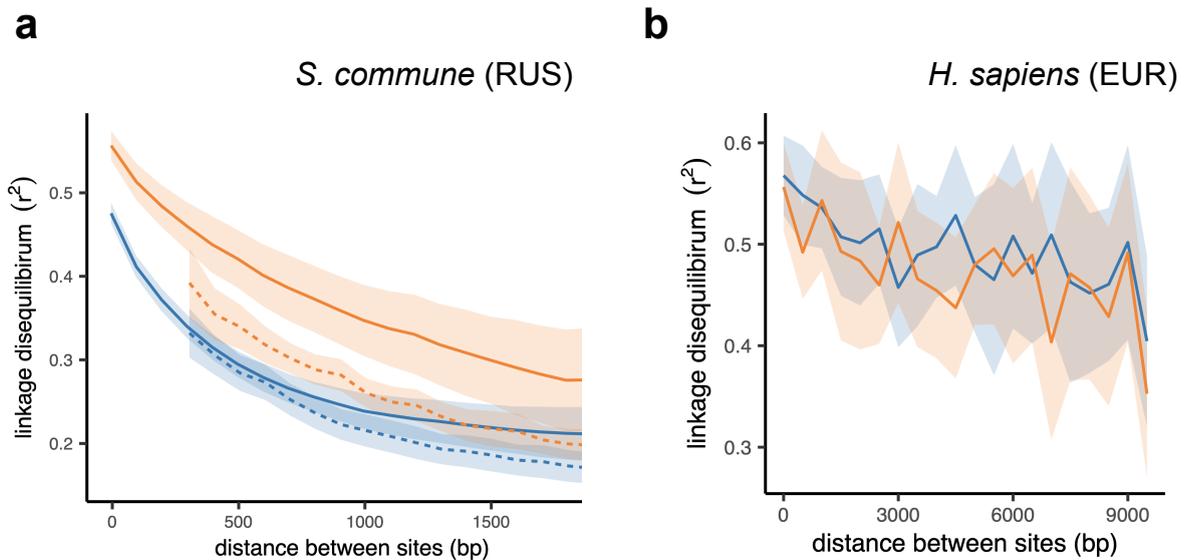
## Elevated LD between nonsynonymous polymorphisms

In a vast majority of species, nucleotide diversity  $\pi$ , the evolutionary distance between a pair of randomly chosen genotypes, is, at selectively neutral sites, of the order of 0.001 (as in *Homo sapiens*) or 0.01 (as in *Drosophila melanogaster*) (Leffler et al. 2012; Cutter, Jovelin, and Dey 2013). Still, a few hyperpolymorphic species with  $\pi > 0.1$  are known, of which the wood-decaying fungus *Schizophyllum commune* is the most extreme, where  $\pi = 0.20$  or 0.13 in the USA or the Russian populations, respectively (Baranova et al. 2015). We studied 34 haploid genotypes from the USA and 21 from Russia and compared the LD between nonsynonymous SNPs ( $LD_{\text{nonsyn}}$ ) to that between synonymous SNPs ( $LD_{\text{syn}}$ ).

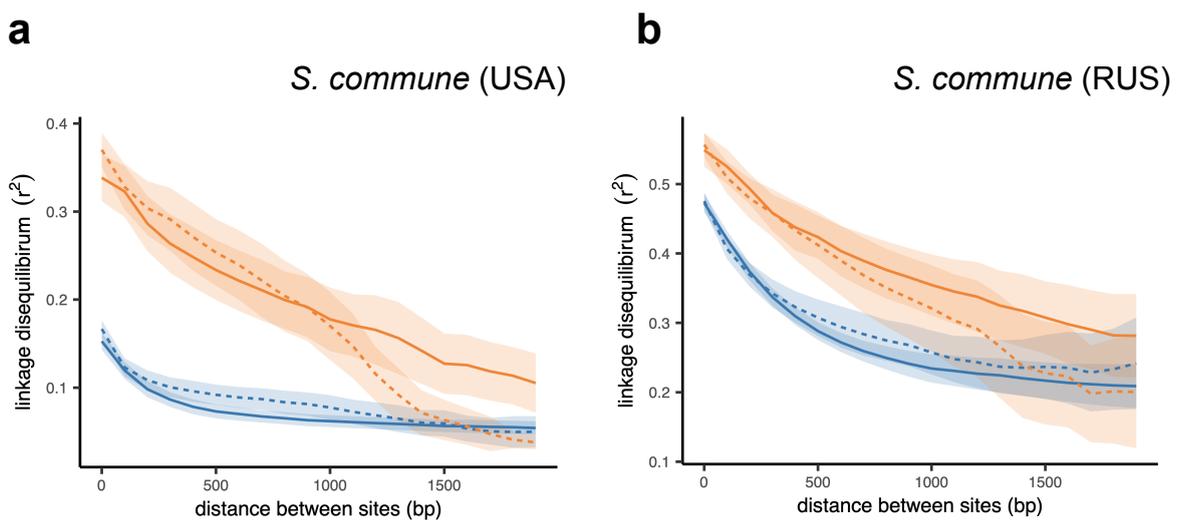
At sites with minor allele frequency (MAF)  $> 0.05$ , in both *S. commune* populations  $LD_{\text{nonsyn}}$  is much higher than  $LD_{\text{syn}}$  at the same nucleotide distance (Figure. 3.4a; for *S. commune* populations, the MAF  $> 0.05$  corresponds to excluding singletons only). This excess of  $LD_{\text{nonsyn}}$  is much stronger for pairs of SNPs located within the same gene, compared to pairs of SNPs from adjacent genes at the same distance. By contrast, the excess of  $LD_{\text{nonsyn}}$  is independent of whether the two SNPs are located within the same or in different exons of a gene (Figure. 3.6). In *S. commune*, the recombination rate is higher within exons (Seplyarskiy et al. 2014), which may affect the patterns of LD; however, this factor could only reduce within-gene LD, and in any case cannot explain the difference between  $LD_{\text{nonsyn}}$  and  $LD_{\text{syn}}$ . A much weaker excess of  $LD_{\text{nonsyn}}$  over  $LD_{\text{syn}}$  for MAF  $> 0.05$  is also observed in the less genetically diverse *D. melanogaster* population (Figure. 3.4b). In the still less polymorphic human populations,  $LD_{\text{nonsyn}}$  is indistinguishable from  $LD_{\text{syn}}$  at the same distances (Figure. 3.4c). The results are reproduced in the Russian population of *S. commune* and in the European ancestry super-population of human (Figure 3.5).



**Figure 3.4. The efficiency of epistatic selection in populations with different levels of genetic diversity.** (a-c) LD in natural populations for SNPs with MAF > 0.05. (a) USA population of *S. commune*, (b) Zambian population of *D. melanogaster*, (c) African superpopulation of *H. sapiens*. Filled areas in (a)-(c) indicate SE of LD calculated for each chromosome or scaffold separately. (d-f) A hyperpolymorphic population (d) may occupy a sizeable chunk of a complex fitness landscape, leading to pervasive positive epistasis, while variation within less polymorphic populations (e and f) is confined to smaller, and approximately linear, portions of the landscape, so that no strong epistasis and LD can emerge. The area of the landscape covered by the population is shown in green.

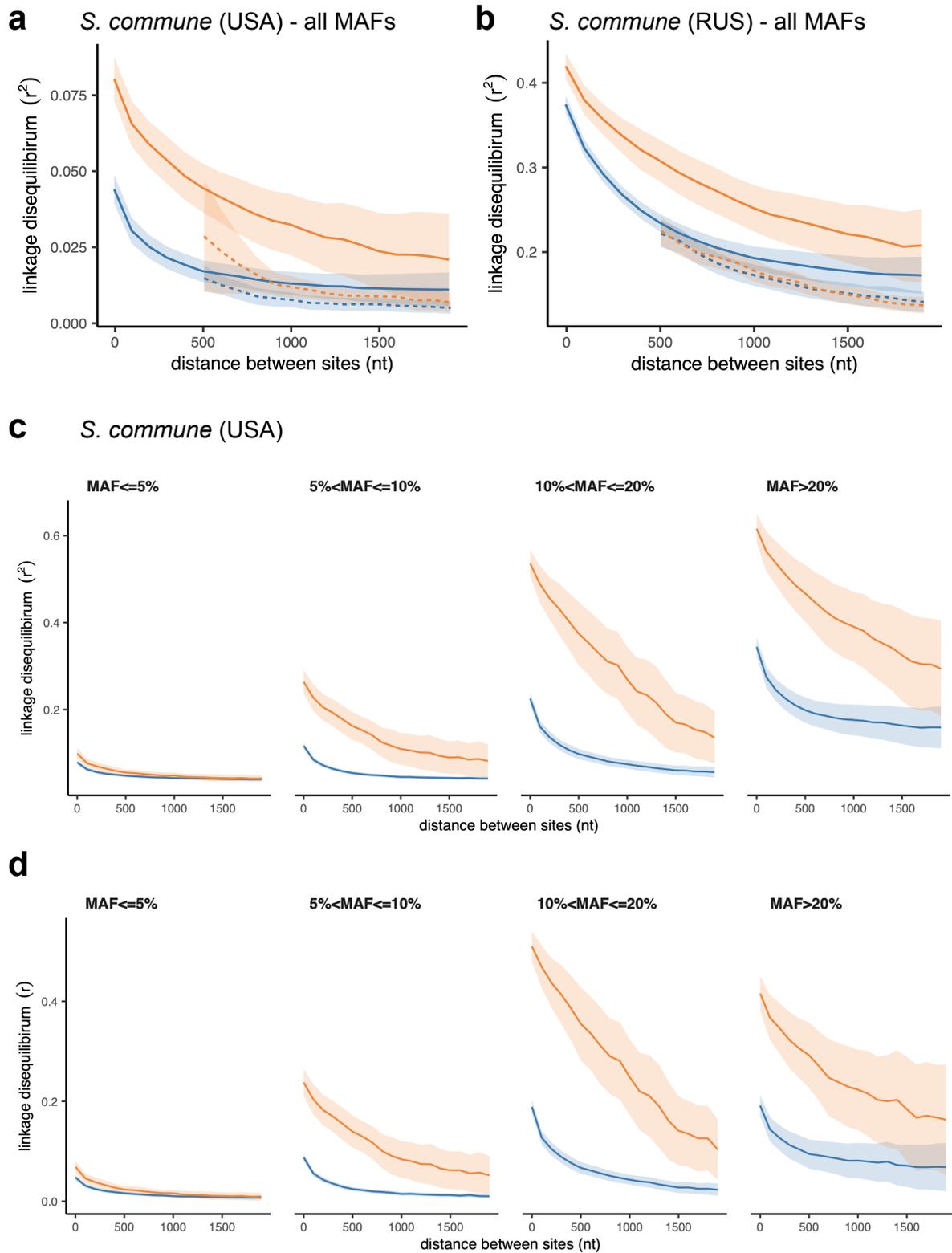


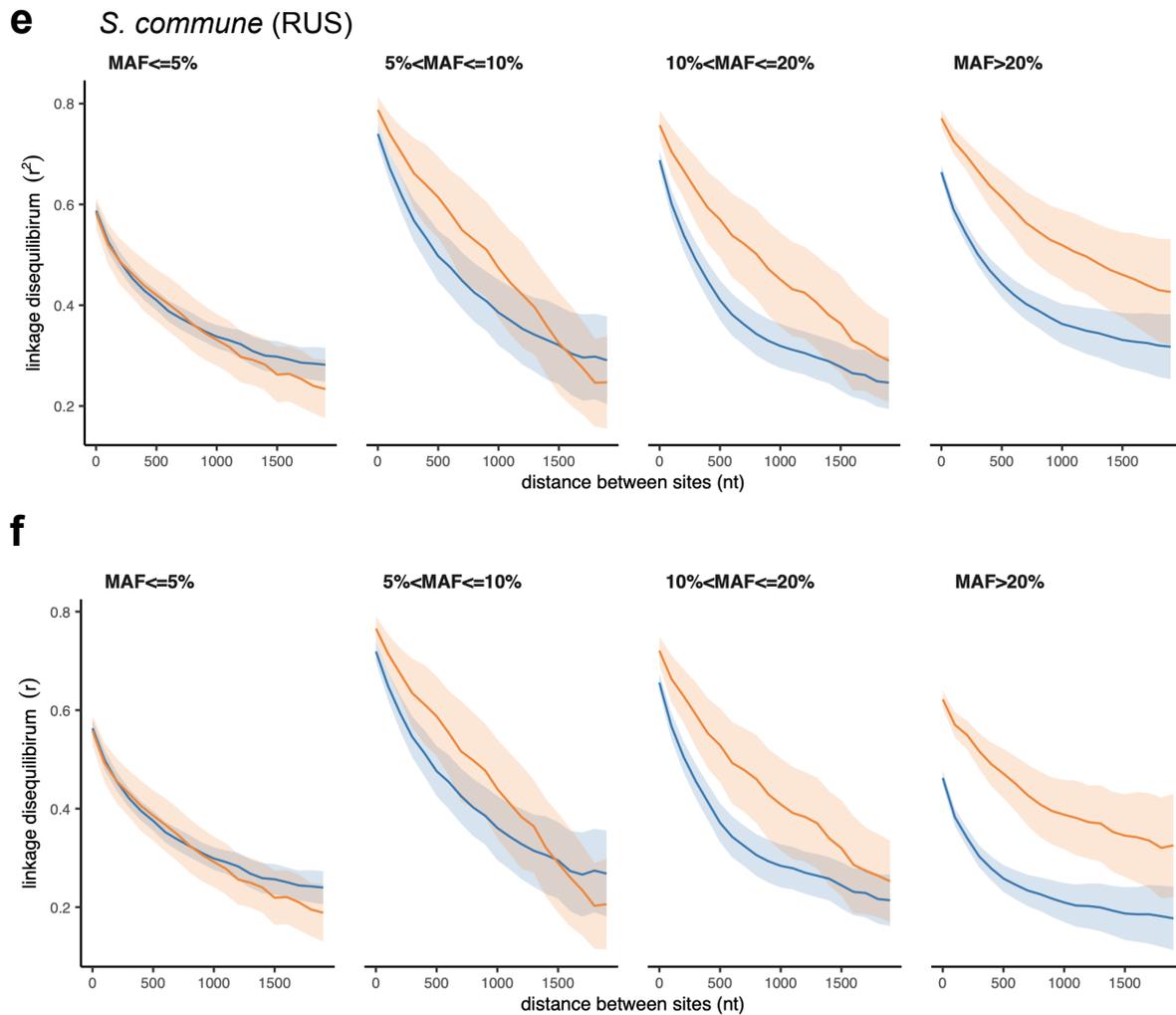
**Figure 3.5. Linkage disequilibrium in the Russian population of *S. commune* and EUR super-population of *H. sapiens*.** LD between nonsynonymous SNPs is shown in orange, and LD between synonymous SNPs is shown in blue. (a) Russian population of *S. commune*, (b) European super-population of *H. sapiens*. Solid lines indicate LD between pairs of SNPs located within the same gene; dashed lines correspond to pairs of SNPs located in different genes. Only SNPs with minor allele frequency > 0.05 are analysed. Filled areas indicate SE of LD calculated for each chromosome (for human) or scaffold (for *S. commune*) separately.



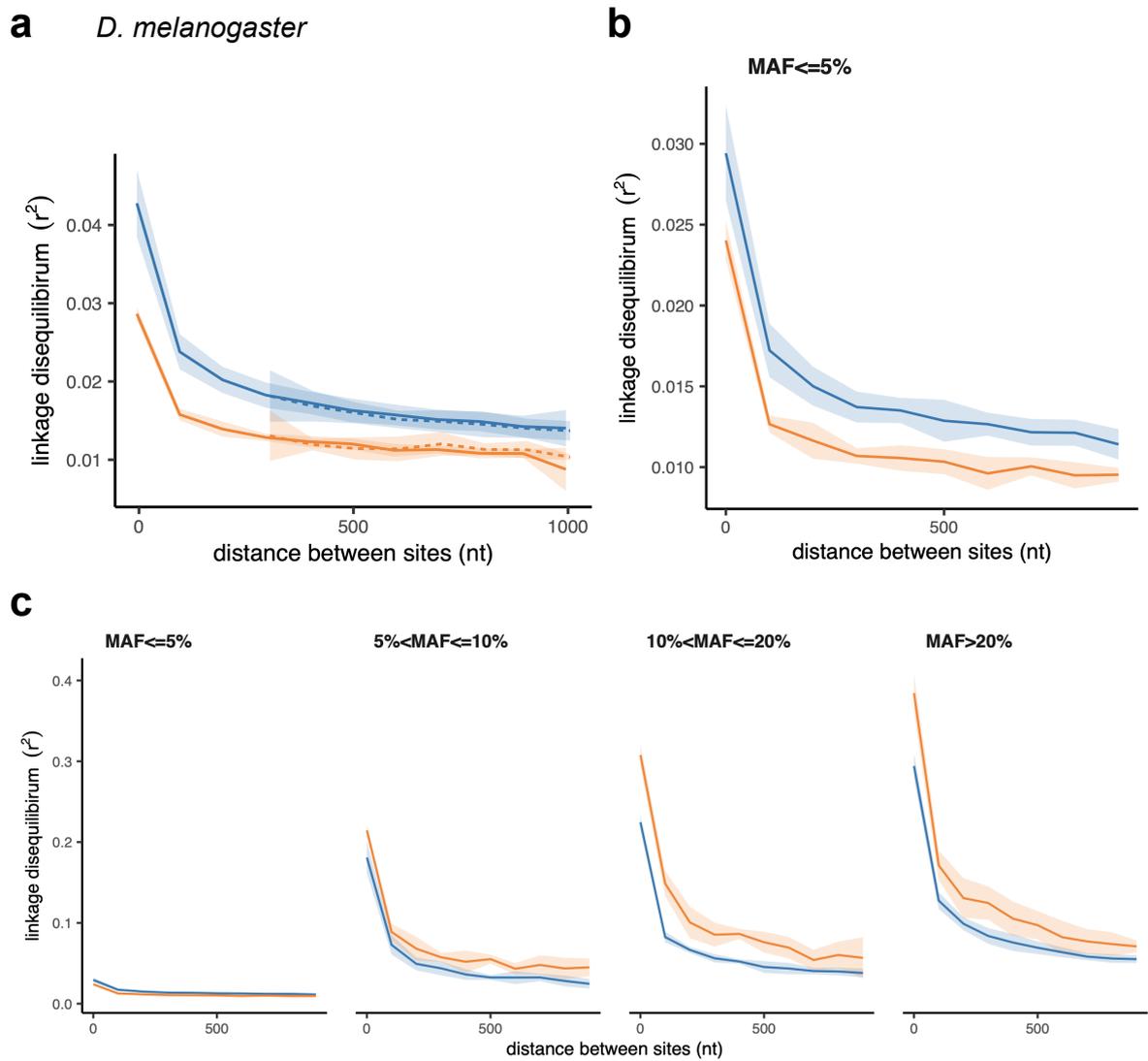
**Figure 3.6. Linkage disequilibrium within and between exons in *S. commune*.** LD between nonsynonymous SNPs is shown in orange, and LD between synonymous SNPs is shown in blue. Solid lines indicate LD between pairs of SNPs located within the same exon of the gene; dashed lines correspond to pairs of SNPs located in different exons of

the gene. **(a)** USA population of *S. commune*, **(b)** RUS population of *S. commune*. Only SNPs with minor allele frequency > 0.05 are analysed. Filled areas indicate SE of LD calculated for each scaffold separately.

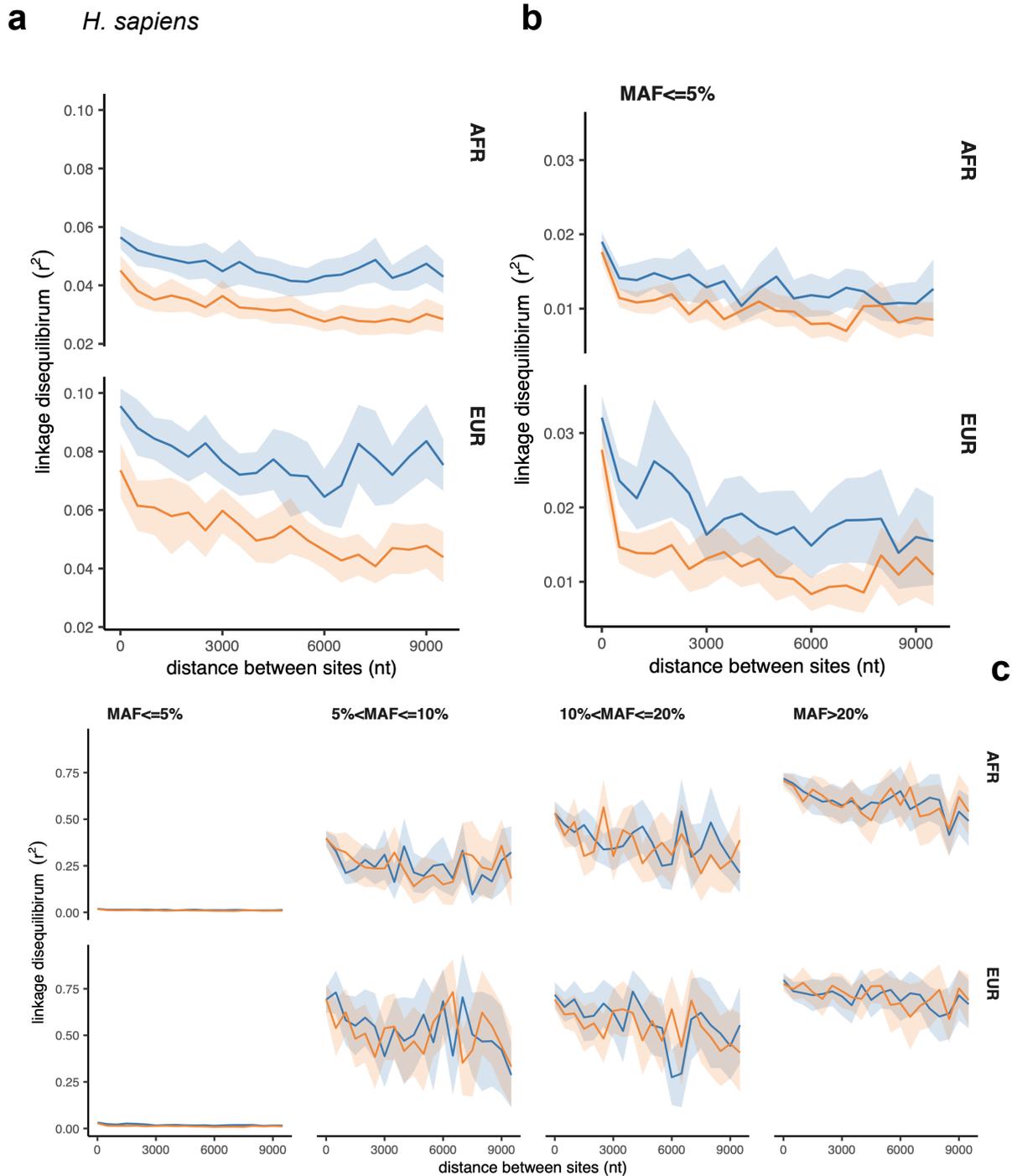




**Figure 3.7. LD between SNPs with different MAF in *S. commune*.** LD between nonsynonymous SNPs is shown in orange, and LD between synonymous SNPs is shown in blue. Filled areas indicate SE of LD calculated for each scaffold separately. **(a, b)** LD between all pairs of SNPs pooled together. Solid lines indicate LD between pairs of SNPs located within the same gene; dashed lines correspond to pairs of SNPs located in different genes. **(c-f)** LD for pairs of SNPs split by MAF. **(c, e)** LD measured as  $r^2$ , **(d, f)** LD measured as  $r$ .



**Figure 3.8. LD between SNPs with different MAF in *D. melanogaster*.** LD between nonsynonymous SNPs is shown in orange, and LD between synonymous SNPs is shown in blue. Filled areas indicate SE of LD calculated for each chromosome separately. **(a)** LD between all pairs of SNPs pooled together. Solid lines indicate LD between pairs of SNPs located within the same gene; dashed lines correspond to pairs of SNPs located in different genes. **(b)** Pairs of SNPs with  $MAF < 0.05$  (large scale). **(c)** Pairs of SNPs split by MAF.



**Figure 3.9. LD between SNPs with different MAF in *H. sapiens*.** LD between nonsynonymous SNPs is shown in orange, and LD between synonymous SNPs is shown in blue. Filled areas indicate SE of LD calculated for each chromosome separately. **(a)** LD between all pairs of SNPs pooled together. Solid lines indicate LD between pairs of SNPs located within the same gene; dashed lines correspond to pairs of SNPs located in different genes. **(b)** Pairs of SNPs with MAF < 0.05 (large scale). **(c)** Pairs of SNPs split by MAF.

Although we report LD between pairs of polymorphic sites as  $r^2$ , which is symmetric regarding the major or minor variants, the observed high values of  $r^2$  correspond to positive LD between minor alleles for both synonymous and nonsynonymous SNPs (Figure 3.7d,f). Thus,  $LD_{\text{nonsyn}} > LD_{\text{syn}}$  means that attraction between minor nonsynonymous alleles is stronger than between minor synonymous alleles. This pattern may seem to be surprising, because there are three factors that work in the opposite direction.

First, random drift, which affects nearly-neutral synonymous sites more than nonsynonymous sites which are mostly under negative selection, leads to attraction between minor alleles (Sandler, Wright, and Agrawal 2021). Second, negative selection at nonsynonymous sites causes repulsion between rare, deleterious alleles, due to Hill-Robertson interference, even if this selection does not involve any epistasis (W. G. Hill and Robertson 1966; Comeron, Williford, and Kliman 2008; Garcia and Lohmueller 2021). Third, there are data on negative epistasis in this selection, which also should lead to repulsion of deleterious alleles and, thus, negative LD between rare nonsynonymous alleles (Sohail et al. 2017; Garcia and Lohmueller 2021; Sandler, Wright, and Agrawal 2021). The first and the second factors are weak and can produce noticeable LD only between tightly linked loci, while the third factor may generate even long-range LD. By contrast,  $LD_{\text{nonsyn}} > LD_{\text{syn}}$  can be explained only by positive epistasis in selection at nonsynonymous sites.

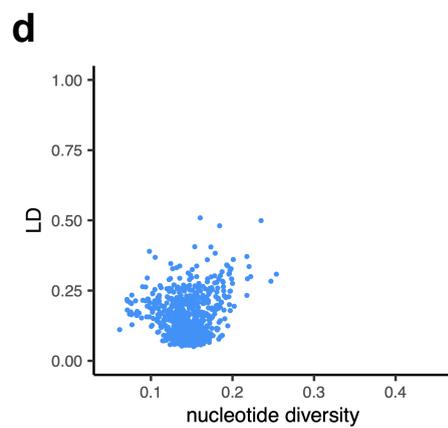
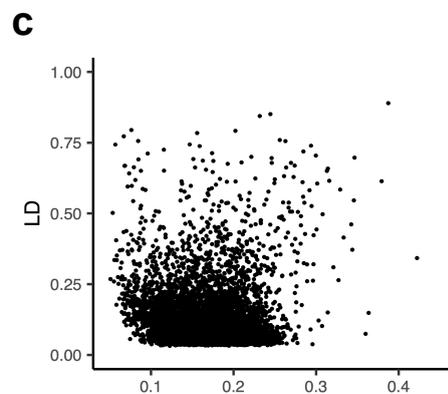
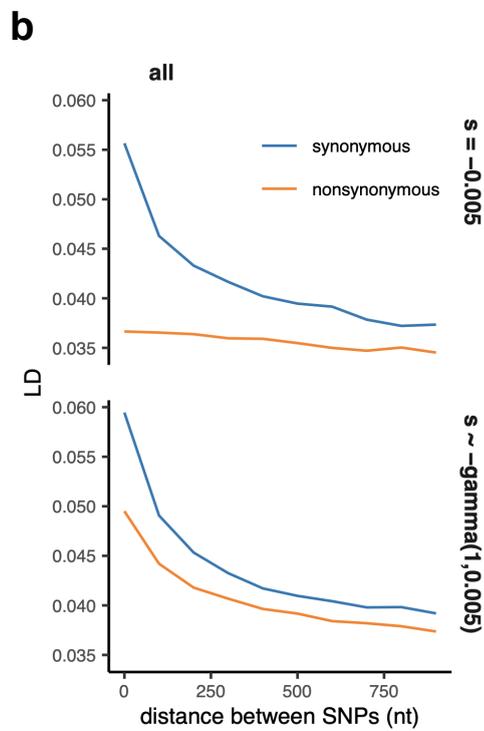
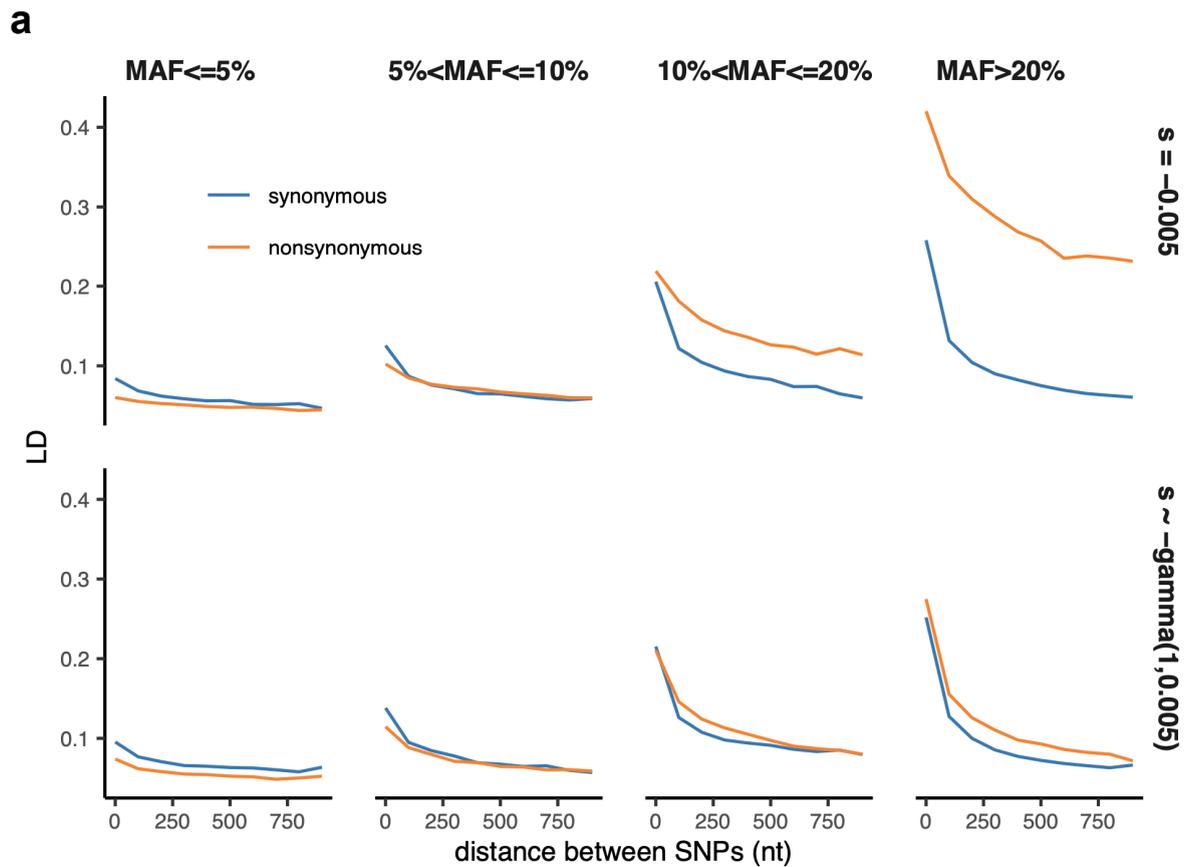
Although negative selection generally results in  $LD_{\text{nonsyn}} < LD_{\text{syn}}$ , our simulations demonstrated that Hill-Robertson interference without epistasis can produce attraction between minor alleles under a rather restrictive set of conditions. In these simulations, weakly deleterious polymorphisms can achieve high frequency only in regions of low recombination, leading to  $LD_{\text{nonsyn}} > LD_{\text{syn}}$  for extremely high MAF (Figure 3.10a). However, this effect doesn't hold if assuming unequal fitness effects of deleterious mutations or while merging SNPs of different frequencies together (Figure 3.10b-d).

This attraction can only appear due to positive epistasis between such alleles — higher-than-expected fitness of their combinations. Positive epistasis can be expected to cause stronger LD in more polymorphic populations (Figure 3.4d-f) and must be more

common for pairs of sites located within the same gene, which are more likely to interact with each other.

For *S. commune*, the excess  $LD_{\text{nonsyn}}$  holds under different minor allele frequency thresholds (Figure 3.7) However, in *D. melanogaster* and *H. sapiens*, rare nonsynonymous SNPs (with  $MAF < 0.05$ ) taken alone show the opposite trend: the LD between such SNPs is reduced compared to synonymous SNPs at the same nucleotide distance (Figures 3.8, 3.9). In human populations, the vast majority of SNPs are rare, leading to  $LD_{\text{nonsyn}} < LD_{\text{syn}}$  when all allele frequencies are considered (Figure 3.9), in line with recently published results (Garcia and Lohmueller 2021).

Decreased LD between negatively selected polymorphisms is expected due to Hill-Robertson interference between deleterious alleles (W. G. Hill and Robertson 1966; Roze and Barton 2006); this effect has been described previously for *H. sapiens* (Garcia and Lohmueller 2021) and *D. melanogaster* (Sandler, Wright, and Agrawal 2021) and is observed in our simulations (Figure 3.14). In addition, both allele frequencies and  $LD_{\text{nonsyn}}$  can be reduced by negative epistasis between deleterious alleles (Garcia and Lohmueller 2021), similarly to the negative LD detected among loss-of-function polymorphisms in humans, flies and plants (Sohail et al. 2017; Sandler, Wright, and Agrawal 2021).

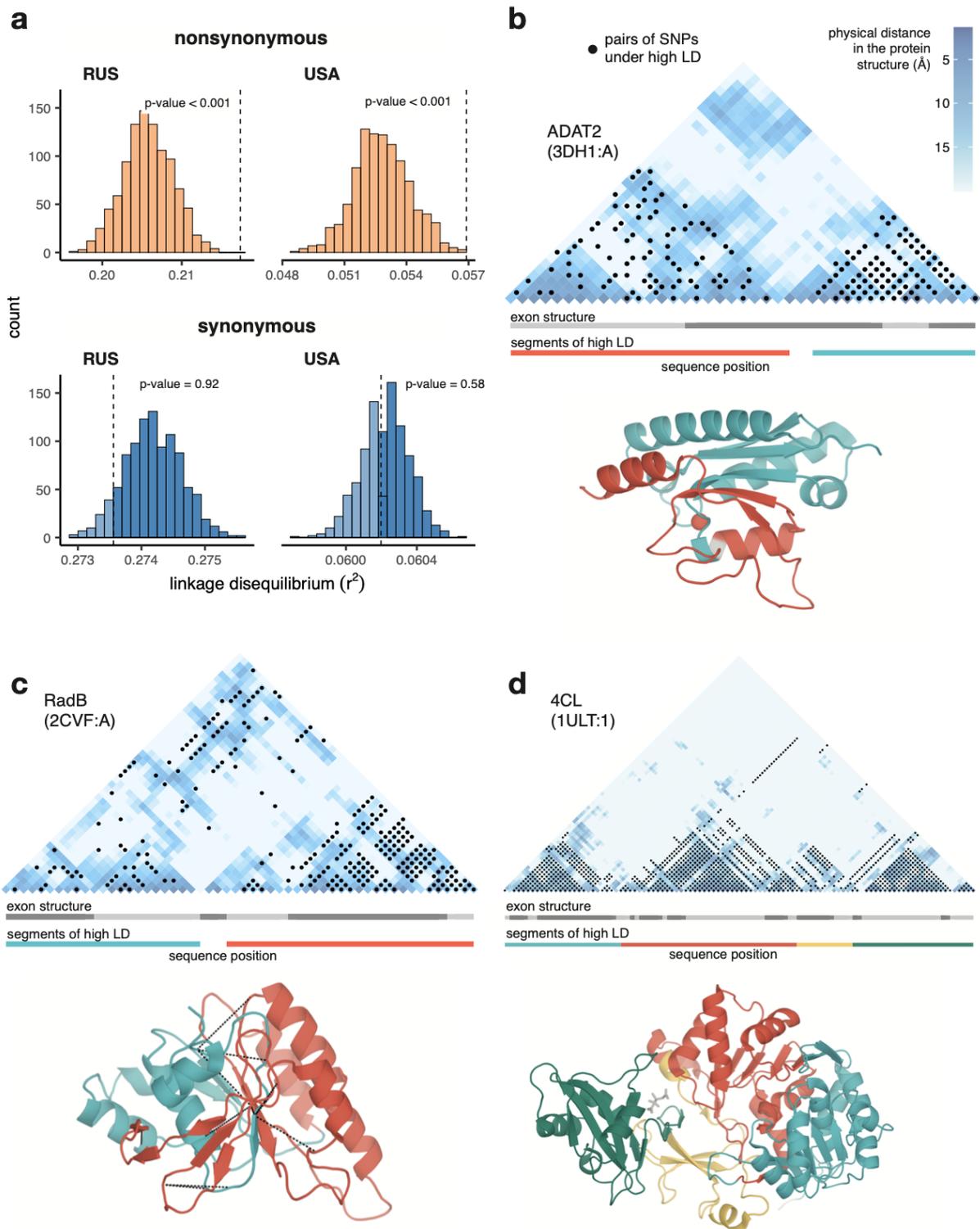


**Figure 3.10. Patterns of LD in simulations under negative selection.** (a) LD between nonsynonymous and synonymous pairs of SNPs split by MAF. (b) LD between all pairs of nonsynonymous and synonymous SNPs pooled together. (a-b) Haploid population size  $N = 2000$ , sequence length  $L = 1000$  bp. Top panels - selection coefficients of all nonsynonymous mutations are equal to  $-0.005$ ; bottom panels - selection coefficients of nonsynonymous mutations are gamma-distributed with parameters  $\text{rate}=1$ ,  $\text{scale}=0.005$ . (c) LD and nucleotide diversity within genes of the USA population of *S. commune* (each point represents one gene). (d) LD and nucleotide diversity obtained in simulations.

### Physically interacting amino acid sites are under stronger LD

Natural selection acting on physically interacting amino acids that are located close to each other within the three-dimensional structure of a protein is characterized by strong epistasis which leads to their coevolution at the level of between-species differences (Ovchinnikov, Kamisetty, and Baker 2014; Marks et al. 2011; Sjordt et al. 2018). The Extraordinary diversity of *S. commune* makes it possible to observe an analogous phenomenon at the level of within-population variation. In both *S. commune* populations, pairs of nonsynonymous SNPs are in stronger LD when they are located at codons encoding physically close (within  $10 \text{ \AA}$ ) than distant amino acids (Figure 3.11a; permutation test  $p\text{-value} < 1e\text{-}3$ ). This is not the case for pairs of synonymous SNPs (Figure 3.11a; permutation test  $p\text{-value} = 0.58$ ).

Hyperpolymorphism of *S. commune* allows us to identify individual proteins with significant associations between the patterns of LD and of physical interactions between sites. We identified 22 genes with pairs of adjacent sites having significantly higher LD in the USA population (out of 1,286 eligible genes in total), and 87 genes in the Russian population (out of 967) at a 5% FDR (Table A2); three examples are shown in Figure 3.11b-d. The alignment of ADAT2 protein contains two segments (Figure 3.11b, teal and red colors), characterized by high within-segment LD. The boundaries of these segments match that of structural units of the protein, but not the exon structure of its gene. In RadB protein, a similar pattern is observed, and LD is also elevated between pairs of SNPs from different segments on the interface of the corresponding structural units (Figure 3.11c). The alignment of 4CL protein can be naturally split into four high-LD segments, which also match its structure (Figure 3.11d).



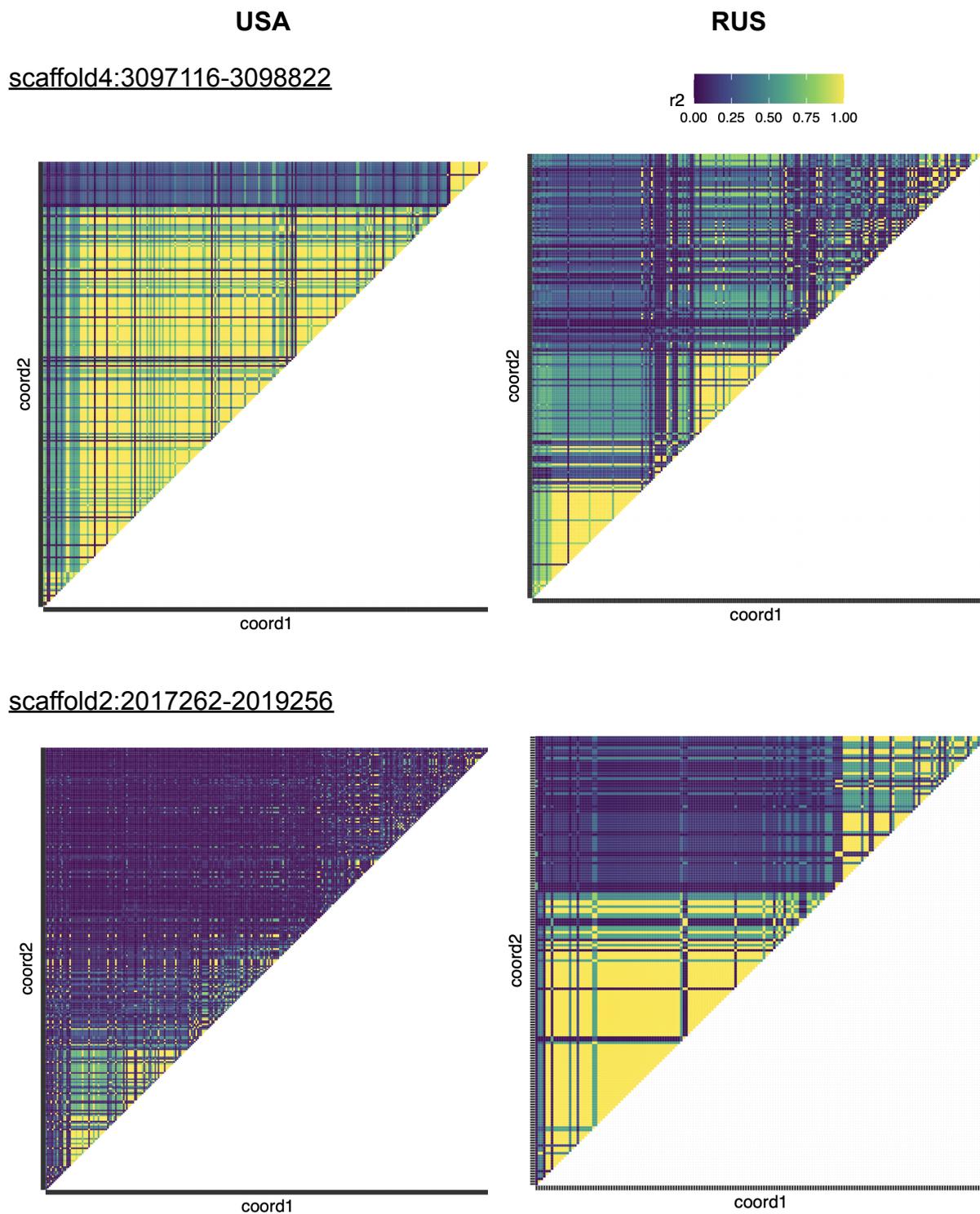
**Figure 3.11. Excessive LD between physically interacting protein sites.** (a) Within pairs of SNPs that correspond to pairs of amino acids that are colocalized within 10 Å in the protein structure, the LD is elevated between nonsynonymous, but not between synonymous, sites. Dashed lines show the average LD. Permutations were performed by randomly sampling pairs of non-interacting SNPs while controlling for genetic distance

between them, measured in amino acids; pairs of SNPs closer than 5 aa were excluded. **(b-d)** Examples of proteins with LD patterns matching their three-dimensional structures. The axis correspond to genomic positions, so that the diagonal corresponds to the nearest polymorphic sites; the heatmaps show the physical distance between each pair of polymorphic sites in the protein structure. Black dots correspond to pairs of sites with high LD ( $> 0.9$  quantile for the gene). Dashed lines show high LD between physically close SNPs from different segments of high LD. In these examples, LD is calculated in the Russian population of *S. commune*.

### **Excess of $LD_{\text{nonsyn}}$ is more pronounced in distinct regions of high LD**

The magnitude of LD varies widely along the *S. commune* genome. Visual inspection of the data shows a salient pattern of regions of relatively low LD, alternating with mostly short regions of high LD (haploblocks, Figure 3.12). We calculated LD along the genome in a sliding window of 250 nucleotides and regarded as a haploblock any continuous genomic region with LD values that belong to the heavy tail of its distribution.

In the USA population, 8.4% of the genome is occupied by 5,316 such haploblocks, 56% consist of regions with background LD level, and the rest cannot be analyzed due to poor alignment quality or low SNP density. 88% of the haploblocks are shorter than 1,000 nucleotides, although the longest haploblocks spread for several thousands of nucleotides. In the Russian population, there are 10,694 haploblocks, occupying 15.9% of the genome, and regions of background LD cover 39% of it. There is only a modest correlation between the USA and Russian haploblocks: the probability that a genomic position belongs to a haploblock in both populations is 2.3% instead of the expected 1.3%, indicating their relatively short persistence time in the populations.



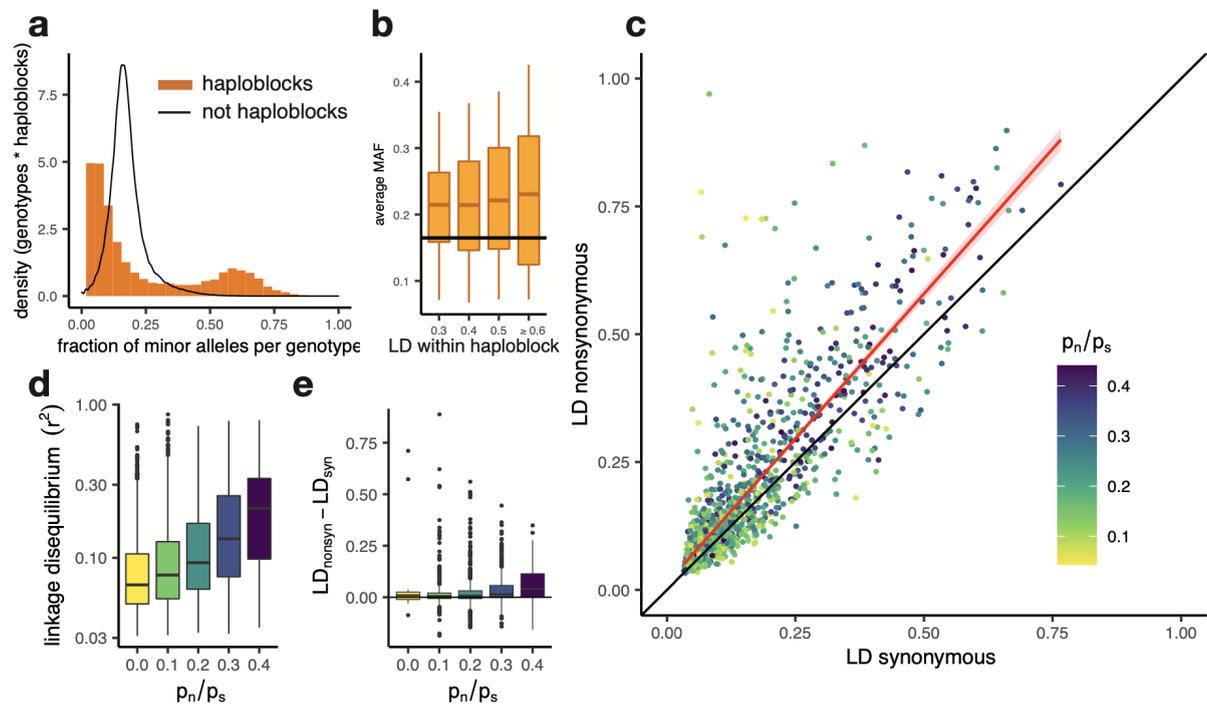
**Figure 3.12. Examples of haploblocks in two populations of *S. commune*.** The heatmaps show LD between polymorphic SNPs in the same genomic regions in the USA and RUS populations of *S. commune*. Only biallelic polymorphic sites with minor allele frequency > 1 are shown, the number of such can vary between

populations. The coordinates of the haploblocks are shown in the titles (the length is ~2000 nt); the number of SNPs in these regions differs between the populations.

LD within a haploblock is usually so high that most genotypes can be attributed to one of only two distinct haplotypes, which carry different sets of alleles. This results in a bimodal distribution of the fraction of minor alleles in a genotype within a haploblock, because some genotypes belong to the major haplotype and, thus, carry only a small fraction of minor alleles, and other genotypes belong to the minor haplotype and, thus, possess a high fraction of minor alleles (Figure 3.13a). Polymorphic sites within haploblocks are characterized by higher MAF than that at sites that reside in non-haploblock regions (t-test p-value < 2e-16 for both populations), and in the USA population MAFs within a haploblock are positively correlated with its strength of LD (Figure 3.13b, Pearson correlation estimate = 0.07, p-value < 2e-6).

There is no one-to-one correspondence between haploblocks and genes, which are, on average, longer. Still, different genes are covered by haploblocks to different extent, which leads to wide variation in the strength of LD and other characteristics among them. The excess of  $LD_{\text{nonsyn}}$  over  $LD_{\text{syn}}$  is also largely restricted to the genes with high LD, e.g. containing haploblocks (Figure 3.13c). As a result, because both haplotypes tend to be common in a haploblock (Figure 3.13), this excess is much stronger for loci with MAF > 0.05.

LD between alleles of all kinds is higher within genes with large  $p_n/p_s$  (Spearman correlation p-value < 2e-16, Figure 3.13d). The same is true for the excess of  $LD_{\text{nonsyn}}$  over  $LD_{\text{syn}}$  (Figure 3.13e, Spearman correlation p-value = 4.4e-17). Although positive correlation between  $p_n/p_s$  and LD is expected under Hill-Robertson interference, positive correlation between  $p_n/p_s$  and the excess of  $LD_{\text{nonsyn}}$  may be indicative of positive epistasis weakening negative selection acting on the nonsynonymous polymorphisms.



**Figure 3.13. Patterns of linkage disequilibrium in the USA population of *S. commune*.** (a) Distribution of the fraction of polymorphic sites that carry minor alleles in a genotype within haploblocks. Black line shows the distribution of fraction of minor alleles in genotypes in non-haploblock regions. In haploblocks (orange), the majority of genotypes carry either small or large number of minor alleles, since they represent one of the persisting haplotypes. (b) Distributions of the average MAF within a haploblock for haploblocks with different average values of LD. The average MAF in non-haploblock regions is shown as a horizontal black line for comparison; the average MAF expected under neutrality after exclusion of singletons is 0.17. (c) LD between nonsynonymous and synonymous SNPs within individual genes. Linear regression of  $LD_{\text{nonsyn}}$  on  $LD_{\text{syn}}$  is shown as the red line. To control for the gene length, only SNPs within 300 nucleotides from each other were analyzed. Genes with fewer than 100 such pairs of SNPs were excluded. (d,e) The positive correlation between  $p_n/p_s$  of the gene and its average LD (d) or the difference between  $LD_{\text{nonsyn}}$  and  $LD_{\text{syn}}$  (e). Here, the data on the USA population of *S. commune* are shown.

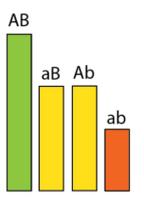
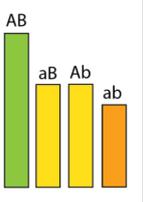
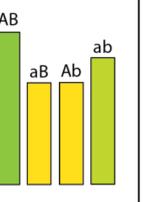
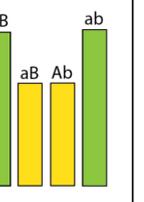
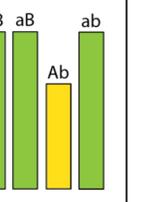
## Excess of $LD_{\text{nonsyn}}$ requires stable polymorphism

Our simulations show that positive epistasis alone cannot lead to an observed large excess  $LD_{\text{nonsyn}}$  over  $LD_{\text{syn}}$ , for which two extra conditions need to be satisfied (Figure 3.14). The general reason for this is simple: in order for a substantial LD between not-too-rare alleles to appear, these alleles must persist in the population for a long enough time.

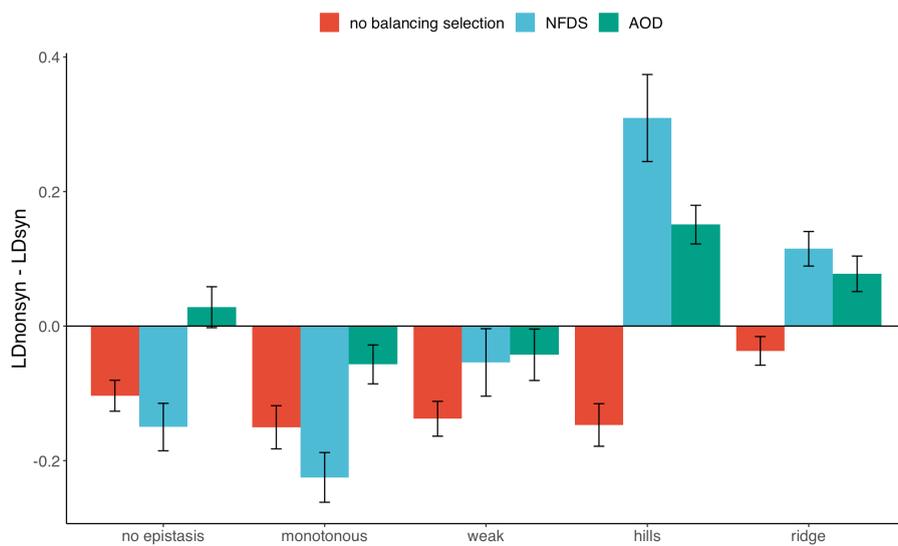
First, positive epistasis must lead to a full compensation of deleterious effects of individual alleles or, in other words, at least two genotypes that are present in the population at substantial frequencies must have (nearly) the same highest fitness (Figure 3.14). If this is not the case, selection favoring the only most-fit genotype leads to a too low level of genetic variation, which persists only due to recurrent mutation. The high-fitness genotypes can correspond either to isolated fitness peaks of equal heights or to a flat, curved ridge of high fitness. The available data are insufficient to distinguish these two options, although it is natural to assume that two major haplotypes that are common within a haplotype block correspond to high-fitness genotypes. Of course, with complete selective neutrality there is no reason for  $LD_{\text{nonsyn}} > LD_{\text{syn}}$ , so that at least some mixed genotypes, carrying alleles from different high-fitness genotypes, must be maladapted.

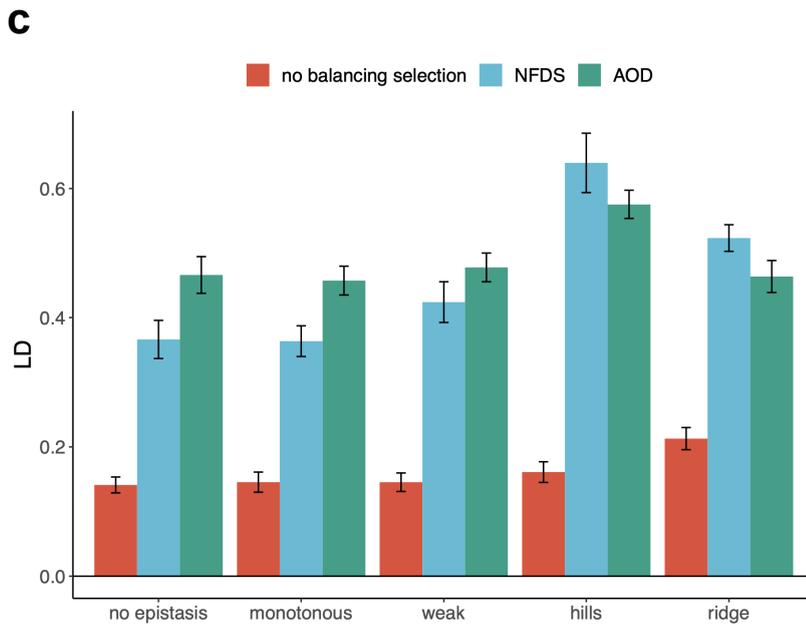
Second, there must be some kind of balancing selection that specifically works to maintain variation, because otherwise random drift does not allow genetic variation to persist for a long enough time even if some, or even all, genotypes are equally fit (Figure 3.14). Here, there are at least two options. On the one hand, a *bona fide* negative frequency-dependent selection (NFDS) can act either directly at loci that display high LD or at some other tightly linked loci. On the other hand, variation can be maintained due to associative overdominance (AOD), resulting from selection against recurrent deleterious mutations at linked loci (Ohta 1971; Zhao and Charlesworth 2016; Gilbert et al. 2020).

**a**

|                        | <br>no epistasis | <br>"monotonous"<br>epistasis | <br>"weak" sign<br>epistasis | <br>"hills" | <br>"ridge" |
|------------------------|---|--|---|---|--|
| no balancing selection | -   | -  | -   | -   | -  |
| NFDS                   | -   | -  | -   | +   | +  |
| AOD                    | -   | -  | -   | +   | +  |

**b**





**Figure 3.14. The excess of  $LD_{nonsyn}$  under pairwise epistasis and balancing selection.** (a) Five models of selection (additive selection and four types of pairwise epistasis) are simulated without balancing selection, under NFDS and AOD. The height of columns shows log fitness of the corresponding genotypes. (+) indicate simulations with  $LD_{nonsyn} > LD_{syn}$ ; (-) indicate simulations with  $LD_{nonsyn} < LD_{syn}$ . (b) The difference between  $LD_{nonsyn}$  and  $LD_{syn}$  in simulations under epistasis and balancing selection. (c) Average LD in the same simulations. Error bars in (b, c) indicate SE calculated based on 100 simulations.

Balancing selection is also a *sine qua non* for the presence of haploblocks, because a pair of divergent haplotypes can evolve only if they coexist for a considerable time. Although the simultaneous existence of two haplotypes may emerge in a finite population under neutrality, simulations without balancing selection didn't reproduce the abundant haploblocks with high LD similar to the ones observed in the data (Figure 3.14c). However, a single locus under NFDS is enough to maintain a haploblock comprising the region of the genome around it. If variation is maintained by AOD, it is more likely that selection against recessive mutations occurs at a number of tightly linked loci. Long coexistence of diverged haplotypes that comprise a haploblock enables accumulation of co-adapted combinations of alleles within them.

So it is not surprising that a pronounced excess of  $LD_{\text{nonsyn}}$  over  $LD_{\text{syn}}$  in *S. commune* is observed primarily within haploblocks.

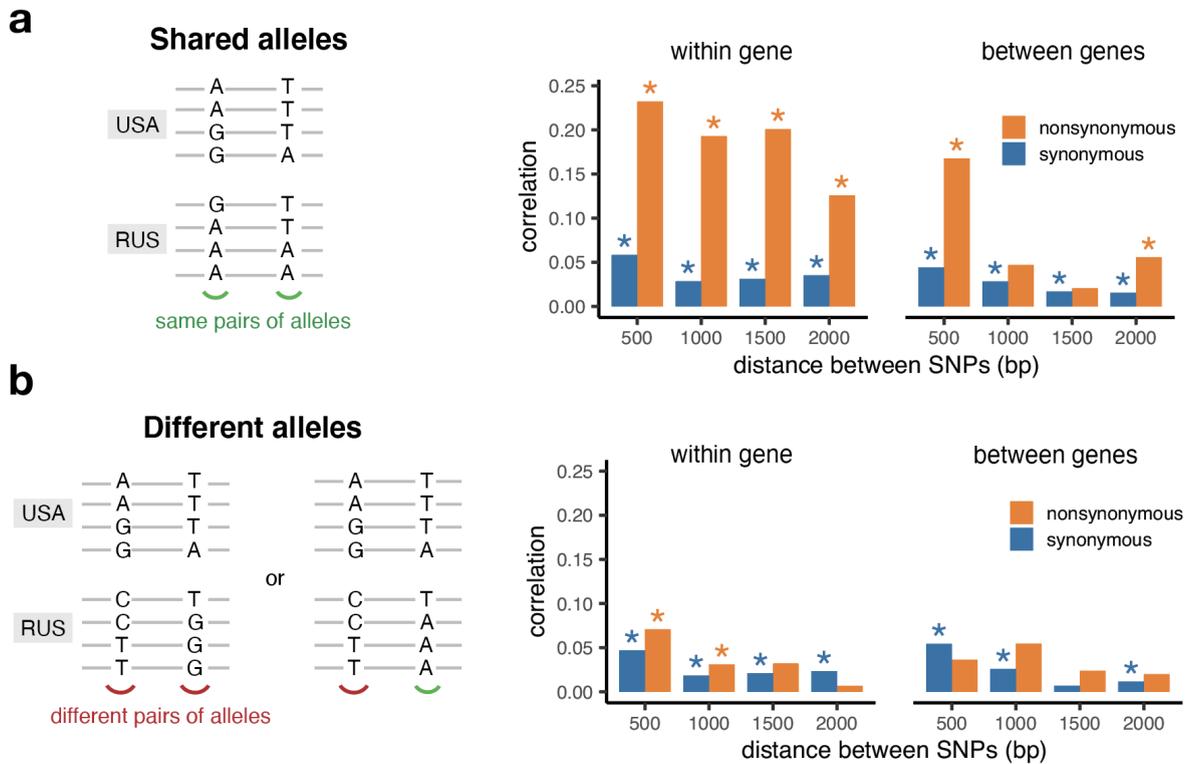
### Correlated LDs between shared SNPs in two populations

Although a high excess of  $LD_{\text{nonsyn}}$  is observed only within haploblocks, a signature of epistasis can also be seen outside of them in the form of a correlation between LDs in the two populations. This correlation can be high even if LDs *per se* are low.

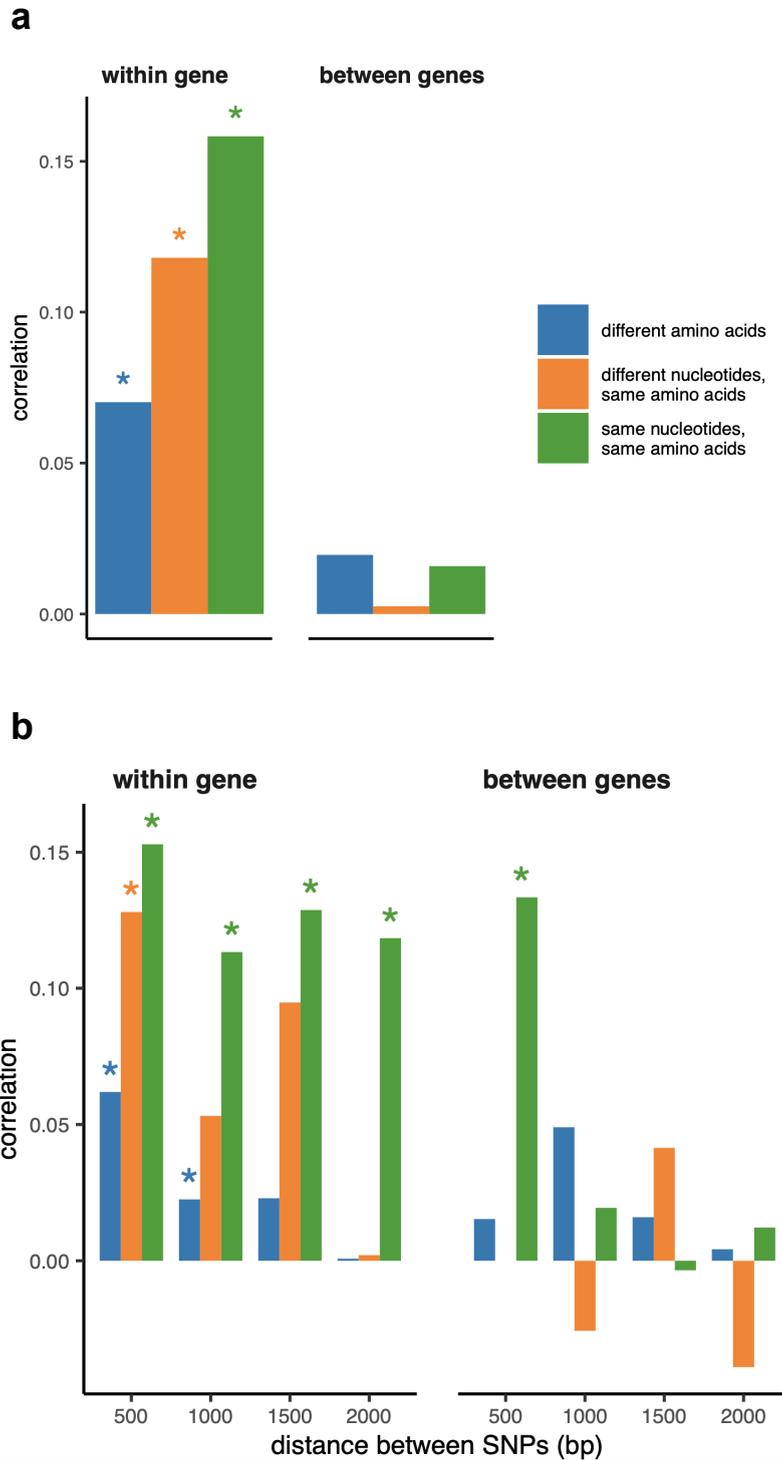
The USA and the Russian populations share a large proportion of their SNPs. Given the high divergence between the two populations, few such shared SNPs are expected to have common origin in the ancestral population, and instead they are likely to have arisen from recurrent mutation. The high prevalence of coincident SNPs is not surprising because SNPs comprise 0.28 and 0.13 of all the aligned nucleotide sites in the USA and Russian populations, respectively (Baranova et al. 2015), Figure 3.1). We identified pairs of shared biallelic SNPs located within 2kb from one another and calculated the LD between them in both populations. To avoid the effects of strong within-population linkage and the occasional co-occurrence of haploblocks between populations, we excluded SNPs located within haploblocks or within genes under high LD ( $> 0.8$  LD quantile for the corresponding population) in either population.

Values of LD in the two populations are strongly correlated only for pairs of nonsynonymous SNPs located within the same gene, and only if both populations carry the same pairs of alleles in the same sites (Figure 3.15). Correlation of LDs is the strongest if shared SNPs carry the same pairs of nucleotides, but is also observed if they encode the same amino acids by different nucleotides (Figure. 3.16). The contrast between correlations within pairs of sites that reside in the same vs. different genes cannot be explained by inheritance of LD from the common ancestral population. Moreover, synonymous SNPs are expected to be on average older than nonsynonymous ones, so that this mechanism should lead to a higher correlation of LDs for pairs of synonymous sites. Thus, the observed pattern can be explained only by epistatic selection shared between the two populations.

Correlation of LDs between SNPs located within haploblocks in both populations is high regardless of whether they reside in the same or different genes, apparently because of occasional coincidence of haploblocks between populations (Figure 3.17).

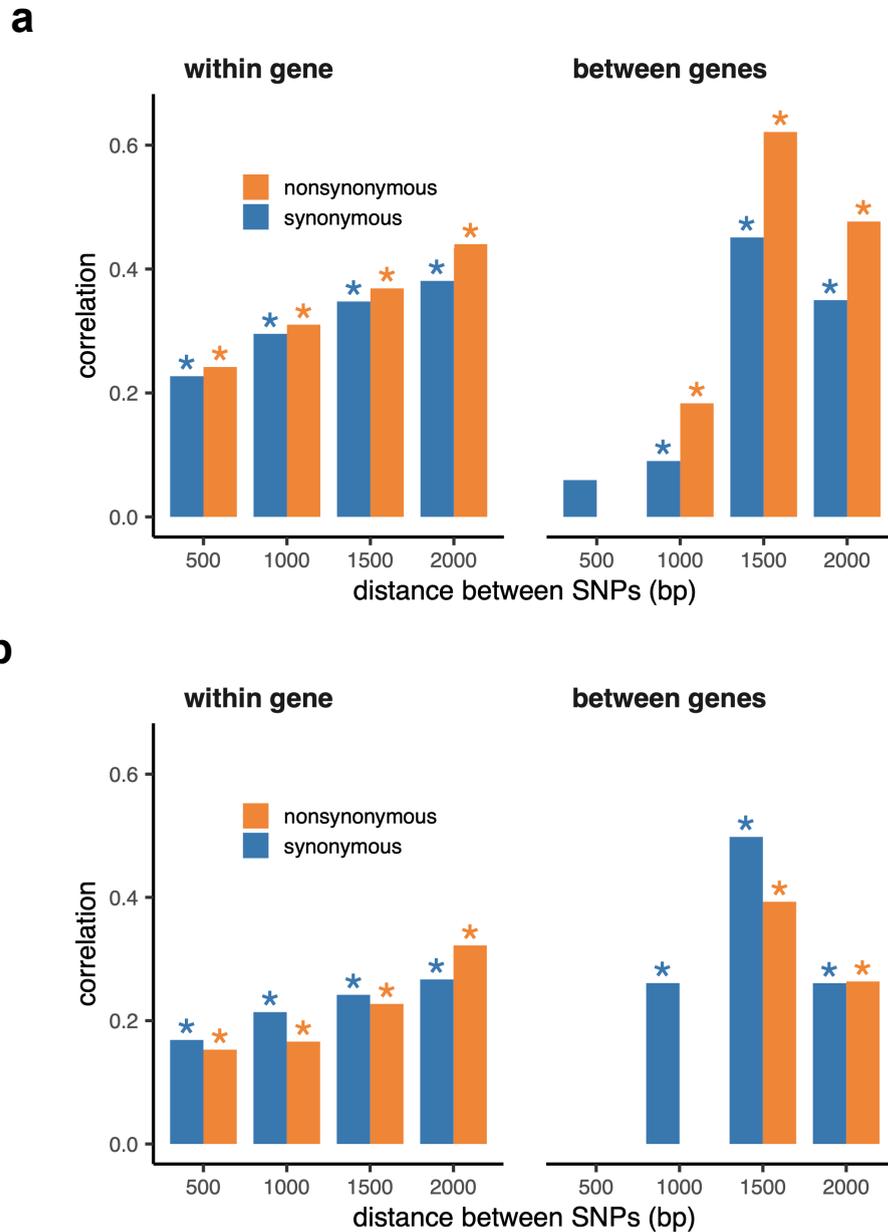


**Figure 3.15. Correlation of LD values between pairs of shared SNPs in the two *S. commune* populations. (a) Pairs of SNPs with the same alleles in both sites, (b) pairs of SNPs differing by at least one allele. Asterisks indicate Spearman correlation p-values < 0.001.**



**Figure 3.16. Association of LD values between pairs of shared nonsynonymous SNPs encoding the same amino acids in the two *S. commune* populations. (a) All pairs of SNPs pooled together. Pair of SNPs is considered to carry different alleles if at**

least one allele differs in at least one site. **(b)** Pairs of SNPs stratified by distance between them. Asterisks indicate Spearman correlation p-values < 0.01.



**Figure 3.17. Association of LD values between pairs of shared SNPs within haploblocks in the two *S. commune* populations. (a) Pairs of SNPs with the same major and minor alleles in both sites, (b) pairs of SNPs differing by at least one allele. Asterisks indicate Spearman correlation p-values < 0.001.**

## Discussion

On top of its most salient property, an exceptionally high  $\pi$ , genetic variation within *S. commune* possesses two other pervasive features. The first is a high prevalence of mostly short haploblocks, genome segments comprising two or occasionally three distinct haplotypes, which is a signature of balancing selection (DeGiorgio, Lohmueller, and Nielsen 2014; Leffler et al. 2013; Rasmussen et al. 2014). The overall fraction of the genome covered by haploblocks is  $\sim 10\%$ , which is about an order of magnitude higher than the fraction covered by detectable signatures of BS in genomes of other species.

The second feature is excessive attraction between rare nonsynonymous alleles polarized by frequency. This pattern is much stronger within haploblocks, indicating that they were shaped by both balancing and epistatic selection, so that amino acids common within a haplotype together confer a higher fitness. Polymorphisms that involve haplotypes that comprise many interacting genes, such as inversions (Theodosius Dobzhansky and Pavlovsky 1957; Brian Charlesworth and Charlesworth 1973; Singh 2008; Sturtevant and Mather 1938) and supergenes (Mather 1950; Joron et al. 2011; Kunte et al. 2014), are known from the dawn of population genetics, but here we are dealing with an analogous phenomenon at a much finer scale, because haploblocks are typically shorter than genes. Thus, instead of coadapted gene complexes (Theodosius Dobzhansky and Pavlovsky 1957), haplotypes represent coadaptive site complexes within genes.

In our simulations, equally high fitnesses of two or more genotypes was a necessary condition for a large excess of  $LD_{\text{nonsyn}}$ , because otherwise the polymorphism did not live long enough for any substantial LD to evolve. However, epistasis between loci responsible for real or apparent balancing selection and those involved in compensatory interactions probably abolished the need for this fine-tuning of fitnesses. For example, if each haploblock carries its own complement of partially recessive deleterious mutations, together with alleles engaged in compensatory interactions with each other which also make these recessive mutations less deleterious, AOD can be expected to cause stable coexistence of these alleles.

Why are haploblocks and positive LD between rare nonsynonymous alleles so common in *S. commune*, but not in other, less polymorphic, species? There may be several, not mutually exclusive, reasons for this. Regarding haploblocks, real or apparent balancing selection may be more common in *S. commune* due to its higher polymorphism. Also, the same balancing selection may protect polymorphism in a huge population of *S. commune*, but not in populations with lower  $N_e$ . Finally, an excess of haploblocks in *S. commune* may be at least due to better detection of signatures of balancing selection in a species with an extraordinary density of SNPs. The haploblocks are likely to be maintained in genomic regions with low recombination rate, however, low recombination alone can't explain the existence of the haploblocks of such strength and abundance like the ones we observe in *S. commune*. In the simulations in the absence of epistasis and balancing selection, we weren't able to reproduce high values of LD observed within haploblocks even if the recombination rate is low (and even zero) — it was possible only in simulations under balancing selection.

Excessive  $LD_{\text{nonsyn}}$  in *S. commune* is also likely to be due to its hyperpolymorphism which increases the probability that mutually compensating alleles at a pair of interacting sites achieve high frequency and encounter each other in the same haplotype before being eliminated by selection. In other words, even if the fitness landscape remains the same, it results in more epistatic selection and, thus, in stronger LD in a species whose genetic variation covers a larger chunk of this landscape (Figures 3.3, 3.4).

In a vast majority of species,  $\pi$  is a small parameter  $\ll 1$ . This imposes a severe constraint on operation of selection and obscures signatures of its particular modes. Thus, hyperpolymorphic species where  $\pi$  is  $\sim 1$  provide a unique opportunity to probe midrange properties of the fitness landscape.

## Chapter 4: Correlated positive selection leads to bursts of amino acid replacements

Evolution can occur both gradually and through alternating episodes of stasis and rapid changes. However, the prevalence and magnitude of fluctuations of the rate of evolution remain obscure. Detecting a rapid burst of changes requires a detailed record of past evolution, so that events that occurred within a short time interval can be identified. Here, we use the phylogenies of the Baikal Lake amphipods and of Catarrhini, which contain very short internal edges which make this task feasible. We detect six bursts of nonsynonymous substitutions in individual proteins during such short time periods, each involving between six and 39 substitutions. On average, in the course of a time interval required for one synonymous substitution per site, a protein undergoes a strong burst of rapid evolution with probability at least approximately 0.01.

## Introduction

(Non)uniformity of the rate of evolution is one of the oldest and most contentious issues in evolutionary biology. On the one hand, at both molecular and morphological levels evolution often occurs gradually, so that evolution of a trait doesn't experience drastical bursts of changes (A. P. Martin and Palumbi 1993; Sudhir Kumar and Blair Hedges 1998; dos Reis et al. 2012; Tamura et al. 2012; O'Meara et al. 2006). In particular, this is the case for selectively neutral segments of genomes, which evolve at rates equal to the corresponding mutation rates ("molecular clock", (Zuckermandl and Pauling 1965; M. Kimura and Ohta 1974). Instances of gradual adaptive evolution are also known (Barrick et al. 2009; Mahler et al. 2010).

On the other hand, evolution may also occur mostly through short bursts of changes alternating with long periods of stasis ("punctuated gradualism" or "punctuated equilibrium", (Gould and Eldredge 1993; Stanley 1998). Examples of punctuated equilibrium are provided by the evolution of mammalian body weight (Mattila and Bokma 2008), hominoid body size (Bokma 2002), several morphological traits of rockfish (Ingram 2011), intersexual signalling of cranes (Mooers et al. 1999), and many other data (Wolf et al. 2006; Hunt 2007, 2008; Strotz and Allen 2013; Bedford et al. 2014; Hunt, Hopkins, and Lidgard 2015; Voje 2016). Clearly, both gradual and burst-like evolution does happen, but their relative importance remains controversial (John H. Gillespie 1991; Pagel, Venditti, and Meade 2006; Venditti and Pagel 2008; Pennell, Harmon, and Uyeda 2014). Of course, there can be many causes for alternating episodes of stasis and evolution and of punctuation in molecular and morphological evolution.

There are several models of selection able to produce punctuated dynamics of adaptive evolution. One explanation is occasional drastic changes of the fitness landscape caused by ecological or environmental factors, disarranging the mutation-selection equilibrium and provoking positive selection (Wright 1932; Ville Mustonen and Lässig 2009; V. Mustonen and Lässig 2010). However, even on the static but rugged landscape, a population may experience long periods of stasis near the saddle points, punctuated by short episodes of selection (Bakhtin et al. 2021). Sampling from phylogenies of related species can also cause bias from the constant rate of accumulating of genetic differences

(J. H. Gillespie and Langley 1979). Correlated substitutions at multiple genomic sites can be caused by epistasis between positively selected mutations (Neverov et al. 2021, 2014; Schlosser and Wagner 2008).

In order to detect short bursts of changes, one needs to be able to identify evolutionary events that occurred during a short interval of time before they got averaged by periods of stasis or gradual changes. This is easy to accomplish if a very detailed paleontological record is available, such as those that exist for some marine invertebrates, e. g., Foraminifera (Malmgren, Berggren, and Lohmann 1983). However, such records are exceptions rather than the rule. Furthermore, paleontological data usually shed light only on the morphology of organisms and, thus, cannot reveal bursts of changes at the level of genomes. Fortunately, it may be possible to identify such bursts indirectly, through comparison of genomes of extant species, as long as their phylogenetic tree contains very short internal edges. Unfortunately, despite an avalanche of genomic data, the vast majority of the currently available phylogenetic trees do not satisfy this requirement. Still, we took advantage of two phylogenetic trees that contain such edges, those of the Lake Baikal amphipods (Naumenko et al. 2017) and of Catarrhini (Rosenbloom et al. 2015), UCSC 100 vertebrates multiple alignment), and investigated short bursts in the evolution of their proteins.

## Materials and methods

### Phylogenies of closely related species

In this work, we use two datasets representing multiple alignments of protein-coding sequences of closely related species. First, we consider transcriptomes-based clusters of orthologous genes (COGs) with exactly one ortholog represented in each species for five clades of the Lake Baikal amphipods (gammarids) (64 species and 3399 COGs in total), and the phylogeny based on them (Naumenko et al. 2017). The size of a clade varies from 6 to 24 species (Figure 4.1a). The search for bursts is performed for each clade separately.

Second, we consider a multiple alignment of protein-coding genes from 11 primate species obtained from the 100 vertebrates' genomes alignment of the UCSC Genome Browser together with the corresponding reconstructed phylogenetic tree (Rosenbloom et al. 2015) (Figure 4.1b). In total, there are 17,755 alignments of protein-coding genes of primates containing columns without gaps.

Only internal edges, i.e. segments of the phylogenetic tree ancestral to more than one species, are used in our analysis. For both datasets, we only consider internal edges of length  $< 0.005 dS$  units to focus on short bursts of evolution limited to these internal edges. We use *codeml* program of the PAML package (Ziheng Yang 2007) to reconstruct substitution histories of sequences and to estimate gene-specific  $dN/dS$  values. Only gapless alignment columns are considered.

Presumptive functions of amphipod genes are inferred from blast2GO predictions (Naumenko et al. 2017) and from the genome annotation of a related species *Hyalella azteca* (Poynton et al. 2018); functions of primate genes are inferred from the human genome annotation (hg38).

### Inference of bursts of nonsynonymous substitutions

A classic approach to infer selection acting on the protein-coding sites is to compare the rate of nonsynonymous substitutions (the ones leading to the amino acid replacement)  $dN$  to the rate to synonymous substitutions (not leading to the change of amino acid)  $dS$ ,

which are assumed to be neutral. Positive selection accelerates accumulation of nonsynonymous substitutions, resulting in  $dN/dS > 1$ , while negative selection in the absence of mutational biases restricts amino acid changes, reducing  $dN/dS$  (M. Kimura 1977; Z. Yang and Bielawski 2000; W. H. Li, Wu, and Luo 1985; Lawrie, Petrov, and Messer 2011).

In this work, we use the neutral null model, that assumes the rate of nonsynonymous substitutions equal to the rate of synonymous substitutions ( $dN = dS$ ). This approach isn't adjusted to the negative selection generally acting on the nonsynonymous mutations and leading to  $dN < dS$ , so it allows to detect only the strongest and the fastest bursts capable to overcome the effects of negative selection. We relate  $dN$  of each gene on a particular edge of the phylogenetic tree to the length of this edge, measured in the units of  $dS$  on the basis of all the available genes. We use this approach, instead of considering the  $dS$  value of only the gene that underwent a burst, because it is impossible to estimate the  $dS$  for an individual gene on a short edge with precision. For example, the expected number of synonymous substitutions in a gene encoding a 200 amino acid long protein on the edge of length 0.005  $dS$  is only  $\sim 1$ . The genes having  $dN$  significantly larger than  $dS$  within a particular edge are selected as candidate genes that experienced a short and strong burst of adaptive evolution. The p-values are calculated as the probability of observing this many or more nonsynonymous substitutions in a Poisson distribution with the parameter equal to the edge length ( $dS$ ) (or, for bursts spanning multiple edges in a row, the sum of their lengths) multiplied by the number of nonsynonymous sites in the gene (estimated according to (Nei and Gojobori 1986)). The bursts with Benjamini-Hochberg adjusted p-values  $< 0.05$  compose the primary set of putative bursts.

As a control, we also searched for statistically significant bursts of synonymous substitutions on the same set of short edges using an analogous approach.

### **Filtering of candidate bursts**

In order to eliminate the false positives and to compile a list of genes that experienced bursts of nonsynonymous substitutions, we use stringent filter criteria for the candidate genes.

First, alignments containing >50% of columns with gaps are excluded. Second, to ensure the precise phylogenetic positioning of all substitutions that constituted a burst, we exclude sites with *codeml* posterior probabilities for the reconstructed ancestral variant < 0.8 (which will be the case for discordant genes, e.g. caused by incomplete lineage sorting) and recalculate the statistics for the gene.

Third, to safeguard against contribution of anciently divergent paralogs or pseudogenes rather than orthologs to our findings, we apply additional filtering. We require that the *dS* value that characterized the edge of the putative burst obtained using the considered gene isn't higher than the *dS* value for this edge obtained using all genes (adjusted p-value > 0.001). Genes with substitutions in multiple repeated domains are excluded. Next, we require the absence of evidence for paralogs or duplications of the considered gene as follows. For each gene, we determine the pre-burst sequence as the sequence of the phylogenetic node immediately ancestral to the burst-carrying edge(s) reconstructed with *codeml*, and the post-burst sequence as the reconstructed sequence of the phylogenetic node immediately descendant to it. For gammarids, we map raw transcriptomic reads of the considered gene from all gammarid species from the same clade onto the pre-burst and post-burst sequences. If any reads from any of the species descendant to the edge of the provisional burst support the pre-burst variant, or if any reads from any of the species not descendant to the edge of the provisional burst support the post-burst variant, this gene is discarded. For primates, we align the pre-burst and post-burst sequences of this gene onto the assembled genomes from all species, and proceed analogously.

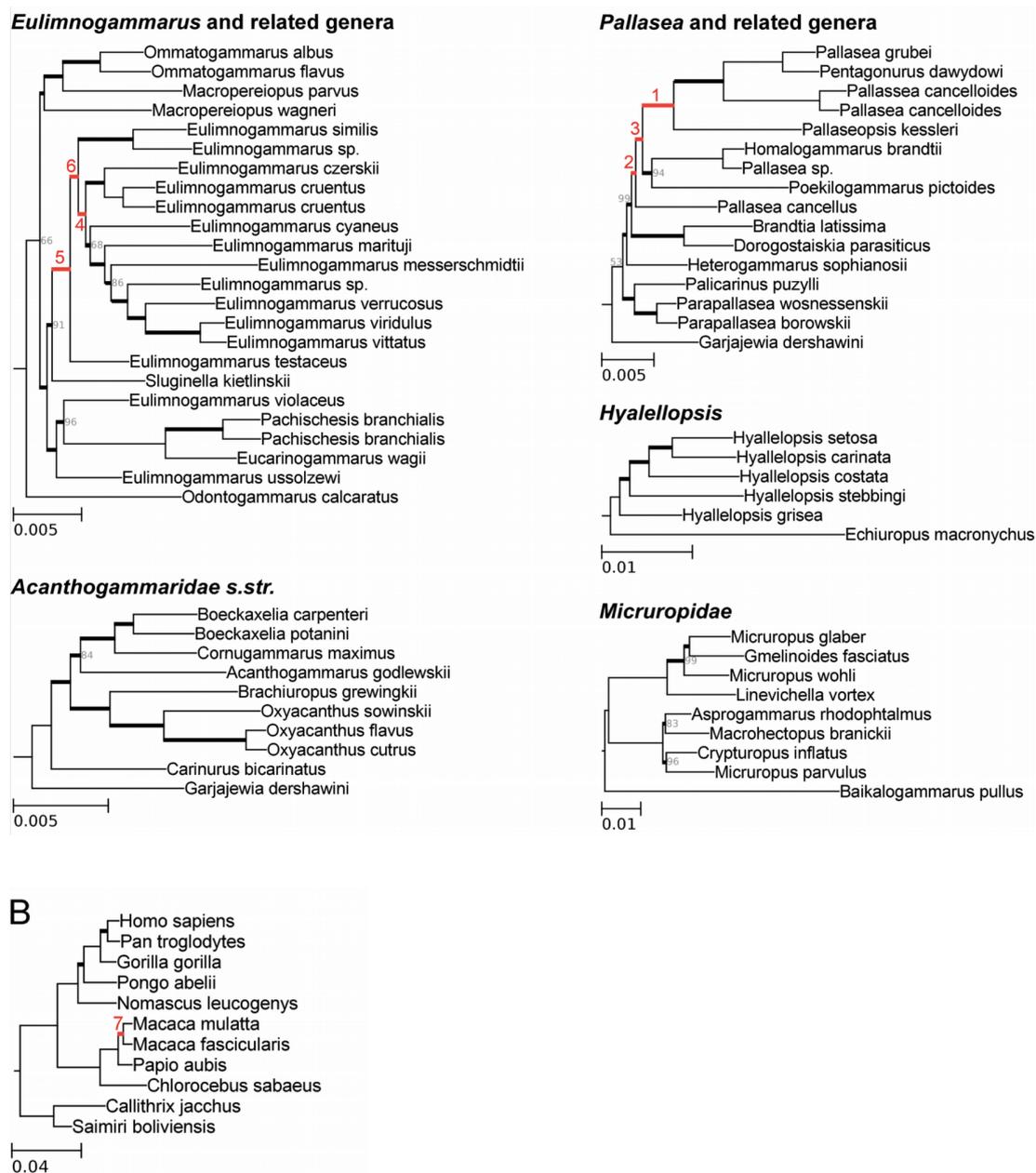
Finally, the burst-containing alignments that survive these filters are curated manually for any evidence for alignment errors, low complexity and unexpected patterns in substitutions. If most substitutions constituting a burst fall into regions of poor alignment or are located in the very beginning or the very end of the gene, the corresponding putative burst is discarded.

## Results

We search for bursts of amino acid substitutions (“bursts”) within internal edges of phylogenetic trees that are shorter than 0.005 *dS*. Suitable edges are present in 5 clades of the phylogenetic tree of gammarids from the Lake Baikal (Naumenko et al. 2017): *Eulimnogammarus* and related genera (18 edges), *Pallasea* and related genera (10), *Hyalleloopsis* (3), *Acanthogammaridae s. str.* (7), and *Micruropidae* (4); as well as within the Catarrhini clade (3 edges) of the tree of vertebrates (Figure 4.1) (Rosenbloom et al. 2015). A burst consists of several amino acid substitutions which occurred in a protein within such an edge or, perhaps, within several successive edges of combined length below 0.005 *dS*.

In gammarids, we identified 5 statistically significant bursts that occurred in 5 proteins within 2 clades. 3 of them occurred over the time period corresponding to a short individual edge of the phylogeny, while the remaining 2 spanned two very short adjacent edges. In Catarrhini, there is 1 significant burst satisfying the filtering criteria (Table 4.1). Each burst consists of between 6 and 38 amino acid substitutions, or between 6 and 39 nonsynonymous substitutions (as some amino acid sites underwent multiple nonsynonymous substitutions), scattered throughout the protein (example shown in Figure 4.2). All edges that harbor bursts have 100% bootstrap support. Unfortunately there is no accepted dating for amphipods diversification, so we can't estimate the time of the corresponding bursts with precision; all that we know is that they have occurred after Baikal has originated in the Miocene, i.e., < 30 Ma years ago. Genes that harbor bursts are enriched in proteins located in mitochondria: they constitute 3 of the 5 such genes, although only 14% of the initial set of COGs are annotated as components of mitochondria (binomial test, *p*-value = 0.02) (Naumenko et al. 2017). No significant bursts of synonymous substitutions at short internal edges are observed.

A

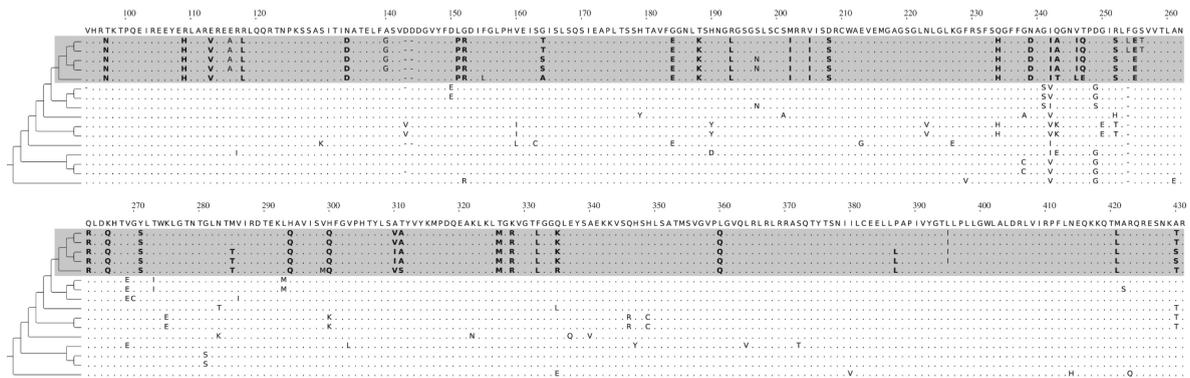


**Figure 4.1. Short internal edges on the reconstructed phylogenies of closely related species.** Internal edges shorter than 0.005  $dS$  are in bold and edges that harbor bursts of evolution (Table 4.1) are in red and numbered. Bootstrap values lower than 100 are shown, branch lengths are measured in units of  $dS$ . (a) Five Baikal gammarids clades (Naumenko et al. 2017). (b) Catarrhini (Rosenbloom et al. 2015).

| clade                                      | edge number | edge length (dS) | gene name      | description of the protein              | overall dN/dS for the gene (excluding the edge with the burst) | substitutions during burst |     | adjusted p-value |
|--|-------------|------------------|----------------|---|--|----------------------------|-----|------------------|
|  |             |                  |                |   |  | nonsyn                     | syn |                  |
| <i>Pallasea</i> and related genera         | 1           | 0.0040           | <i>DNAJC11</i> | DnaJ-like protein subfamily c member 11 | 0.41   | 39                         | 1   | 1.01e-25         |
|  | 1           | 0.0040           | <i>MRPL22</i>  | mitochondrial ribosomal protein L22     | 0.55   | 10                         | 1   | 0.012            |
|  | 2+3         | 0.0007 + 0.0008  | <i>NOP16</i>   | nucleolar protein 16-like               | 0.35   | 6                          | 2   | 0.046            |
| <i>Eulimnoga mmarus</i> and related genera | 4           | 0.0011           | <i>MRPS25</i>  | mitochondrial ribosomal protein S25     | 0.31   | 6.5                        | 0.5 | 0.0011           |
|  | 5+6         | 0.0024 + 0.0010  | <i>AKR1</i>    | aldo-keto reductase                     | 0.57   | 10                         | 1   | 0.032            |
| primates                                   | 7           | 0.0043           | <i>PKR</i>     | interferon-induced protein kinase R     | 1.27   | 18                         | 2   | 0.0010           |

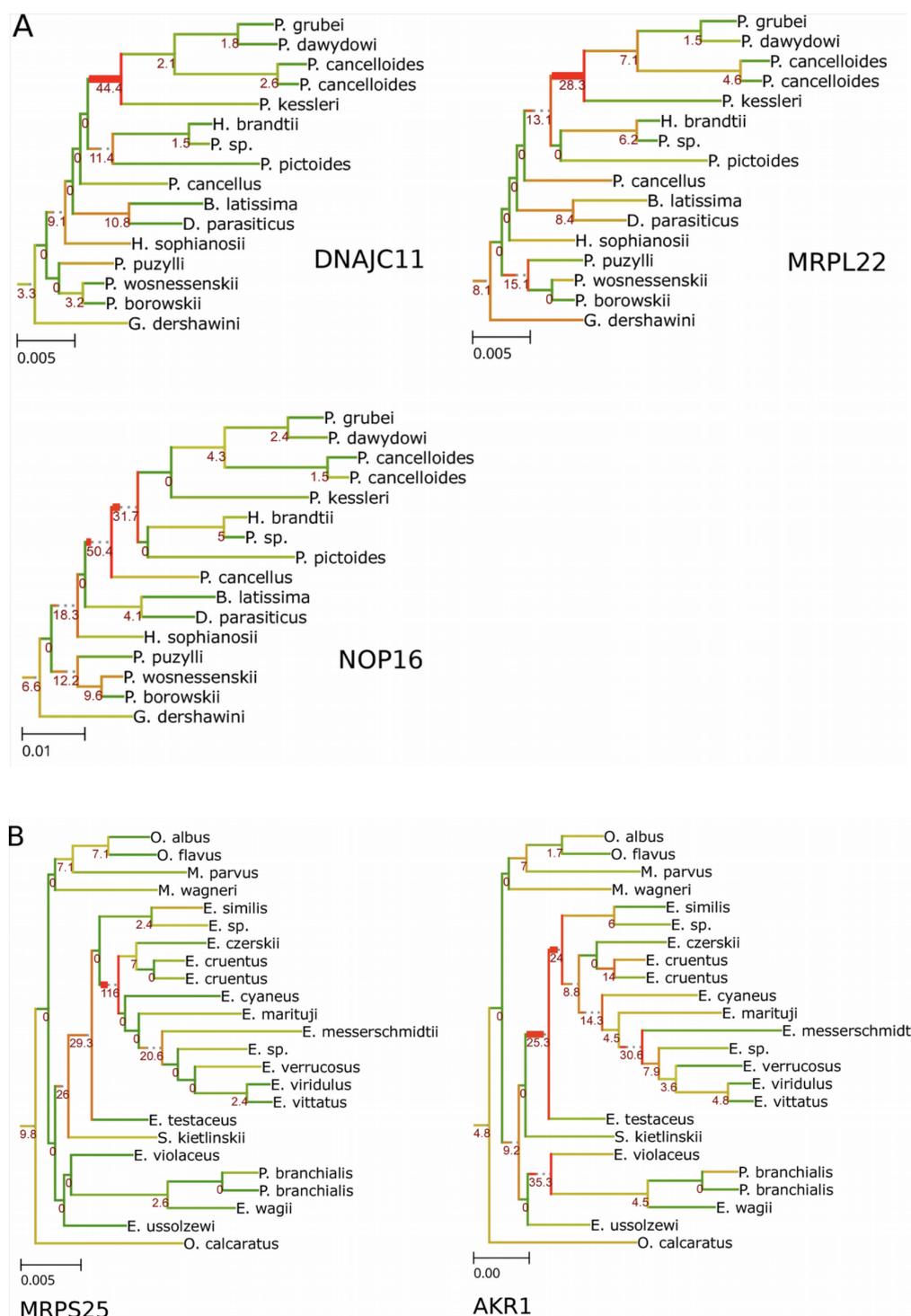
**Table 4.1. Bursts of amino acid substitutions in evolution of proteins of the Baikal gammarids and Catarrhini.**

The most remarkable burst involving 39 nonsynonymous substitutions occurred in the mitochondrial chaperone gene (*DNAJC11*) on an edge of length 0.004 *dS* in the *Pallasea* clade (Figure. 4.2). This edge also harbored another burst in a protein located in mitochondria (L22 ribosome protein) (Table 4.1).

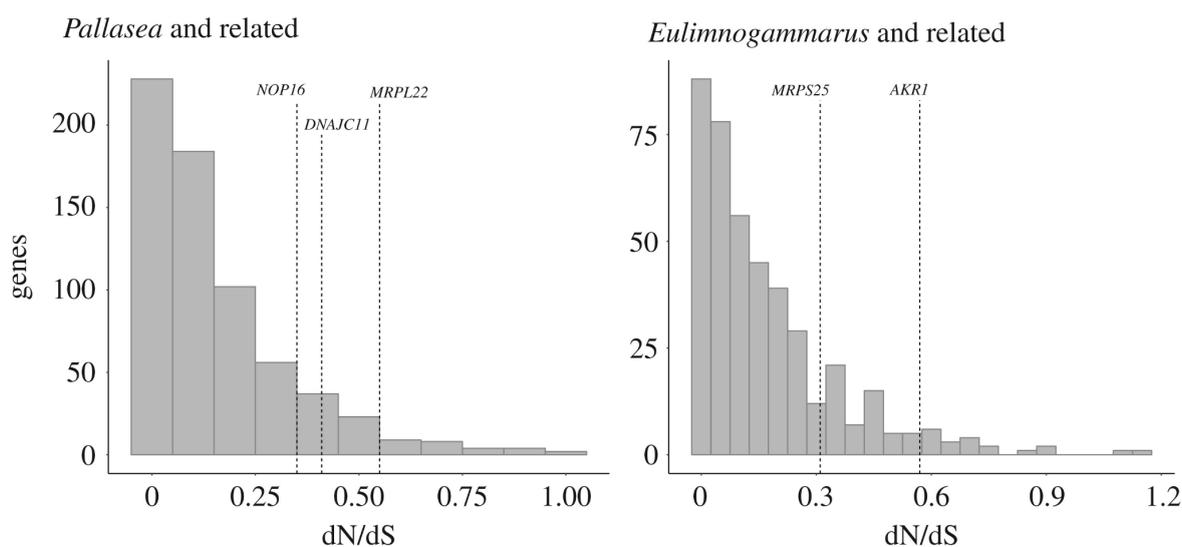


**Figure 4.2. A fragment of the alignment of orthologous *DNAJC11* genes of *Pallasea gammarids*.** Alleles derived in the adaptive burst are shown in bold, sequences originated from the bursts are shown in grey. Total length of the alignment is 1677 nucleotides, or 1515 nucleotides without gaps.

Bursts that are confined to one edge of the phylogenetic tree do not extend to preceding and/or successive edges (p-values for such edges > 0.27) (Figure 4.3). Hence, the characteristic duration of a burst is short,  $\sim 10^{-3}$  *dS*. The overall rate of evolution of some burst-carrying genes was somewhat higher than the average (Table 4.1, Figure. 4.4). Still, after multiple testing correction, there remains no genes with more than one statistically significant burst (p-values > 0.018, adjusted p-values = 1).



**Figure 4.3.  $dN/dS$  values for genes containing adaptive bursts on phylogenetic trees.** Numbers and colors correspond to the edge-specific, gene-specific  $dN$  divided by the length of the edge measured in  $dS$ . Burst-carrying edges are in bold. **(a)** Bursts detected in *Pallasea* and related genera, **(b)** Bursts detected in *Eulimnogammarus* and related genera.



**Figure 4.4. Distribution of  $dN/dS$  in genes of Baikal Amphipods in the clades carrying bursts.** Genes with confirmed bursts are shown with dashed lines.

Phylogenetic tree of 11 species of Catarrhini has only three internal edges shorter than 0.005  $dS$  (Figure 4.1b), and only one statistical significant burst has been detected. This small number may be due to several reasons: Catarrhini species are more distant from each other than gammarids, which results in longer phylogenetic edges and less confident ancestral state reconstruction; moreover, a larger initial dataset leads to a more substantial multiple testing correction.

The detected burst occurred on the internal edge ancestral to two macaque species. Based on the divergence time estimates we assume that the burst occurred approximately 8-3 Ma years ago (Perelman et al. 2011). The gene with the burst encodes the PKR protein (also known as EIF2AK), which is the eukaryotic translation initiation factor 2 kinase activated during viral infection. As in gammarids, the substitutions are scattered along the sequence (Figure. 4.5). Primate PKR contains two dsRNA binding motifs (DRBMs 1 and 2) and C-terminal catalytic kinase domain. The kinase domain carries 14 amino acid substitutions on the selected edge, and DRBM1 the remaining 4. Most substitutions lie in  $\alpha D$ ,  $\alpha G$  and  $\alpha H$  helices or nearby, which have been shown to be enriched in positively selected sites (Rothenburg et al. 2009).  $\alpha G$  helix and specifically positions with amino acid substitutions on the selected edge are involved in



## Discussion

Allele replacements driven by positive selection are the fundamental genetic mechanism of adaptive evolution. These replacements can occur independently of each other or be correlated (Bazykin et al. 2004; Neverov et al. 2014, 2021; Bakhtin et al. 2021). *A priori*, there is a continuum of possibilities, from fully independent individual substitutions to bursts of adaptive evolution, each consisting of multiple substitutions that occurred over a short period of time. We searched for such bursts within individual proteins, taking advantage of two phylogenetic trees, of the lake Baikal gammarids (Naumenko et al. 2017) and of Catarrhini (Rosenbloom et al. 2015).

Using only internal edges is essential because in this case the derived sequence is observed in more than one species, so that rare sequencing and alignment errors would not lead to false discovery of bursts. Our criteria for detection of bursts were rather stringent: conservative filtering of alignments, a neutral null model ( $dN = dS$ ), and multiple testing correction. Unfortunately, it is hard to define the correct null model which takes into account all possible features of protein evolution, for example, those resulting from non-uniformity of the mutation rate along the genomes. Thus, our p-values should be viewed with caution. Even neutrally evolving sequences can carry an increased number of substitutions because of variation in the evolution rate over large time scales (overdispersed molecular clock, (Ohta and Kimura 1971; J. H. Gillespie 1984; Cutler 2000)). However, the aim of our work was not to identify consistent deviations from the Poisson expectation, but to find the most radical outliers. The 6 bursts that we have found are likely to be “real”, in the sense of being caused by simultaneous or near-simultaneous action of positive selection at multiple sites within a protein.

Multiple studies have used methods similar to ours to find episodes of accelerated evolution which could have non-adaptive explanations, such as biased gene conversion (Berglund, Pollard, and Webster 2009; Galtier et al. 2009; Pollard et al. 2006; Brand, Wright, and Presgraves 2019). In contrast, we used stringent criteria for detection of bursts, and those that we found are likely to occur due to positive selection. In particular, the most radical bursts in Galtier et al. involve both synonymous and nonsynonymous substitutions; while the bursts described in our paper, in particular in the novel

amphipod dataset, are limited to nonsynonymous sites. Unlike Galtier and Berglund, we see no GC bias in the bursts-composing substitutions: *e.g.*, the number of AT->GC and GC->AT substitutions in the strongest burst are 14 and 21, correspondingly. This is inconsistent with gene conversion, and instead further supports the adaptive explanation.

Mutations that initiated amino acid substitutions that together constitute a burst are extremely unlikely to appear simultaneously as parts of one complex mutational event. Thus, a burst is likely to involve substitutions that were not precisely synchronous. In other words, a burst lasts longer than a substitution. Still, the bursts that we detected are quite short at the evolutionary time scale. Indeed, four bursts were confined to just one internal edge shorter than 0.005 *dS*. The remaining two bursts, involving 6 and 9 amino acid substitutions, each occurred on two successive internal edges, of lengths 0.0007 and 0.0008, and 0.0024 and 0.0010 *dS*, suggesting that ~3-4 nonsynonymous substitutions occurred per 0.001 *dS* of evolutionary time. Among the substitutions involved in such composite bursts, neither occurred on both edges after a cladogenesis (multiple non-synonymous substitutions did not occur on edges leading to *Eulimnogammarus testaceus* or *Pallasea cancellus*).

Of course, the two phylogenetic trees which we studied almost certainly contained other bursts of positive selection-driven amino acid substitutions which we could not detect with certainty. This would be the case for any burst that occurred within an external edge of a tree, or within an internal edge that is not short enough, or even within a short internal edge as long as the burst itself involved only a small number of substitutions. Unfortunately, we cannot estimate the number of such real but not confidently detectable bursts.

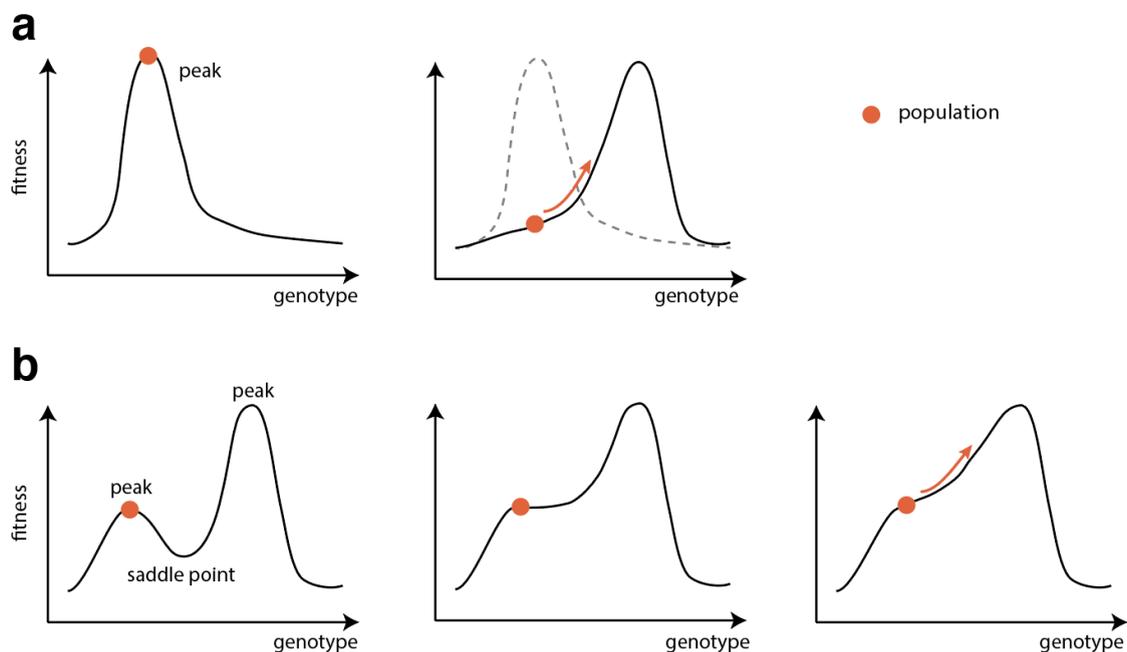
Obviously, our ability to detect a burst depends on the length of the internal edge. Roughly speaking, all bursts that involve at least ~8 amino acid substitutions within a protein of <300 amino acids can be detected within internal edges of length below 0.005 *dS*. Because we investigated 3411 proteins and the total length of all such edges in the gammarid tree was 0.15 *dS*, 5 bursts that we found in them imply that during an interval of time required for 1 synonymous substitution to occur per site, a protein undergoes a

strong burst of adaptive evolution with probability  $\sim 0.01$  (the estimate derived from primates is similar). If so, such bursts are not uncommon.

What can we say about the genes that underwent bursts? Not much: they evolve faster than an average gene, but only marginally so. Unexpectedly, 3 out of 5 genes encode proteins that are located in mitochondria: a mitochondrial chaperone and two mitochondrial ribosome proteins. This observation is hard to explain. Mitochondrial genomes of Baikal Lake gammarids have been shown to undergo intensive rearrangement, which in combination with mito-nuclear discordance and epistatic interactions between mitochondrial proteins coded in nuclear and mitochondrial genomes might lead to this phenomenon (Romanova et al. 2016).

Because multiple nucleotide substitutions that constitute a burst occur very close to each other, making recombination between them negligible, Hill-Robertson interference (W. G. Hill and Robertson 1966) can be expected to impede their fixations. Let us consider the most extreme burst comprising 39 nonsynonymous substitutions on the internal edge of length  $0.004 dS$ . Assuming the per nucleotide per generation mutation rate  $\mu \sim 10^{-8}$ , as in a number of animals (Michael Lynch 2010), this edge corresponds to  $\sim 400,000$  generations, leaving  $\sim 10,000$  generations per each substitution, if they occurred without overlaps. Is this feasible? Every generation,  $2N\mu$  mutations occur at a site, where  $N$  is the census population size and  $\mu$  is the mutation rate. An advantageous mutation will eventually reach fixation with probability  $2sN_e/N$ , where  $N_e$  is the effective population size and  $s$  is the selection advantage of a heterozygous mutation (Ferrière, Dieckmann, and Couvet 2004; Motoo Kimura 1983). Thus, the per generation probability of fixation of a particular advantageous mutation is  $4N_e\mu s$ . Under assumptions of  $N_e = 10^5$  (limited data indicate nucleotide diversity  $\sim 0.01$  in several lake Baikal amphipods),  $\mu = 10^{-8}$ , and  $s = 10^{-2}$ , this probability becomes  $4 \times 10^{-5}$ , which is not very different from  $10^{-4}$ . Thus, successive accumulation of substitutions that constitute a burst, which makes them immune to the Hill-Robertson interference, cannot be ruled out. This would be especially the case if during the whole course of a burst selection favors mutations at all the 39 sites that constitute it, or, in other words, the order in which the substitutions occur is not prescribed (F. A. Kondrashov and Kondrashov 2001b). If so, the target for advantageous mutation at a particular moment of time

consists of all sites where substitutions did not yet occur, and their order depends on the order in which mutations appear.



**Figure 4.6. Possible scenarios of fitness landscape changes driving bursts of adaptive evolution.** (a) Drastic changes of the landscape may cause the extinction of current adaptive peaks and emergence of new peaks. (b) Smooth changes of the landscape (the elimination of the intervening local minimum) may make the path to the nearby adaptive peak accessible.

Correlated positive selection at multiple sites that leads to a burst may emerge due to a variety of mechanisms. One possibility, of course, is a sudden, drastic change of the adaptive landscape of a protein, driving the adaptation of the population towards the newly established adaptive peak (Figure 4.6a). However, a burst can also occur as a result of only a small change of the landscape, if it is caused by a fold bifurcation which eliminates a fitness peak initially occupied by a protein and makes it possible to cross the former adaptive valley (Dodson and Hallam 1977; Steinberg and Ostermeier 2016) (Figure 4.6b). This mechanism is compatible with the fact that all genes with bursts in gammarids show a low overall  $dN/dS$  ratio on the entire phylogenetic tree ( $<0.57$ ), implying that these bursts of evolution affected genes that usually evolve slowly. By contrast, the PKR gene of primates possessed a high  $dN/dS$  ratio ( $>1$ ), implying that the

burst in this gene involved an episode of additional acceleration of evolution which was generally fast. Hopefully, the number of available dense phylogenetic trees will soon become much larger, which will make it possible to study bursts of rapid evolution in more detail.

## Chapter 5: Changes of single-position fitness landscapes affect evolution of amino acid sites

Amino acid propensities at a site change in the course of protein evolution. This may happen for two reasons. Changes may be triggered by substitutions at epistatically interacting sites elsewhere in the genome. Alternatively, they may arise due to environmental changes that are external to the genome. Here, we design a framework for distinguishing between these alternatives. We show that they cause opposite dynamics of the fitness of the allele currently occupying the site. Epistasis leads to the entrenchment of the current allele (the increase of its fitness with time since its origin), while random landscape changes cause its senescence (the decrease its fitness). Using large phylogenies of mitochondrial proteins, we identify 21 significantly entrenched and 28 senescing alleles. By analysing phylogenetic distribution of substitutions in the genomes of vertebrates and insects, we show that the amino acids originating at negatively selected sites experience strong entrenchment, while the amino acids originating at positively selected sites experience senescence.

## Introduction

The description of the shape of fitness landscapes is necessary to fully understand adaptive evolution and speciation (Wright 1932; Maynard Smith 1970; Pál and Papp 2017; Fragata et al. 2019). Unfortunately, the large dimensionality of even the landscapes of individual proteins makes them impossible to measure comprehensively in a direct experiment (Sergey Gavrilets 2004; de Visser and Krug 2014). Still, methods of comparative genomics can be used to assess the integral features of fitness landscapes. The simplest informative unit of landscape structure is the single-position fitness landscape (SPFL) (Bazykin 2015), *i.e.*, a vector of fitness values of all possible alleles at an individual genomic position. SPFLs change with time (Rogozin et al. 2008; Povolotskaya and Kondrashov 2010; A. S. Kondrashov et al. 2010; Usmanova et al. 2015; Goldstein et al. 2015; Zou and Zhang 2015; Klink and Bazykin 2017; Klink, Golovin, and Bazykin 2017); this may affect the optimality of the allele that is currently prevalent at this site, influencing subsequent evolution.

One factor causing changes of SPFL is substitutions at other sites of the genome. For this to be the case, these substitutions need to affect the relative fitness of different variants at the considered site, *i.e.*, these sites have to be involved in epistatic interactions. Epistasis has been postulated to be a prevalent factor of protein evolution and divergence across species (Maynard Smith 1970; D. D. Pollock, Taylor, and Goldman 1999; A. S. Kondrashov, Sunyaev, and Kondrashov 2002; Dimmic et al. 2005; Povolotskaya and Kondrashov 2010; Kryazhimskiy et al. 2011; de Visser, Cooper, and Elena 2011; Breen et al. 2012; McCandlish et al. 2013; de Visser and Krug 2014; Neverov et al. 2014; John H. Gillespie 1991). One expected manifestation of genome-wide epistasis is entrenchment, or the evolutionary Stokes shift (David D. Pollock, Thiltgen, and Goldstein 2012; Goldstein et al. 2015; Shah, McCandlish, and Plotkin 2015) — a phenomenon whereby the relative fitness of the allele currently prevalent at the site increases as substitutions at interacting sites accumulate. The reason for this increase is the constraint imposed by the site in consideration onto epistatically interacting sites. The evolution of the remaining sequence is constrained to preserve the high fitness of the resident allele, and may even increase it; at the same time, this sequence is free to

evolve to become less compatible with other variants not currently present at the site. Over time, this leads to an increase in the fitness of the current allele relative to other alleles, including those that resided at this site earlier.

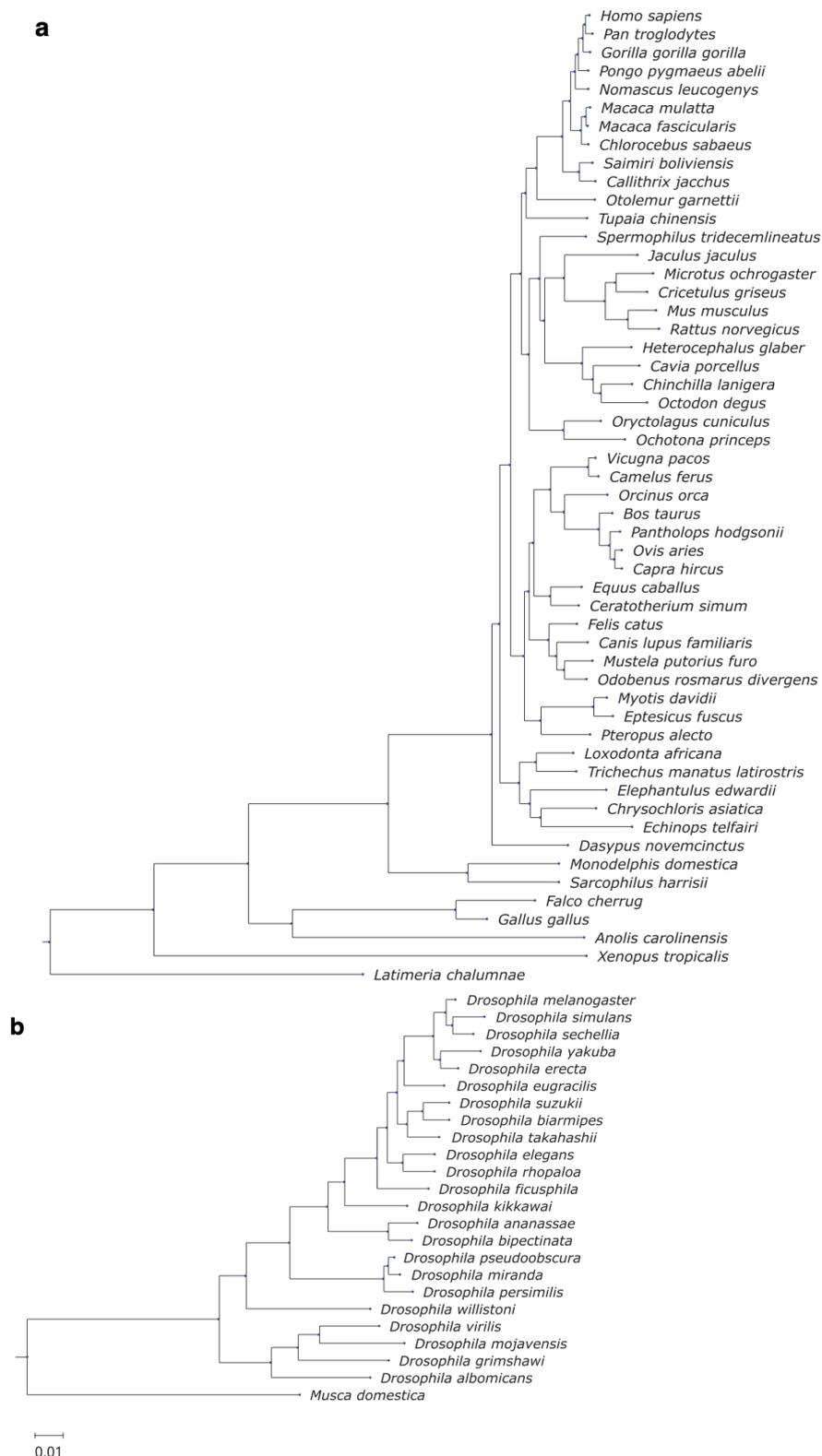
Entrenchment was demonstrated both in simulated protein evolution (David D. Pollock, Thiltgen, and Goldstein 2012; Goldstein et al. 2015; Shah, McCandlish, and Plotkin 2015) and in evolution of real-life proteins. For example, it was shown that reversals of past substitutions follow the phylogenetic distribution indicative of entrenchment: their rate declines with time, indicating that they become more deleterious, *i.e.*, that the current allele becomes more preferable compared to the previous one (Soylemez and Kondrashov 2012; Naumenko, Kondrashov, and Bazykin 2012; Risso et al. 2015; Goldstein and Pollock 2017). The decline in the rate of reversals is caused both by the increase in the fitness of the current allele and the decrease in the fitness of the replaced allele (Naumenko, Kondrashov, and Bazykin 2012).

However, the SPFL may change due to environmental changes even in the absence of epistasis. If such changes are recurrent, the fitness landscape becomes a time-dependent “seascape” (J. Gillespie 1973; Takahata, Ishii, and Matsuda 1975; Huerta-Sanchez, Durrett, and Bustamante 2008; Ville Mustonen and Lässig 2008, 2009; V. Mustonen and Lässig 2010; John H. Gillespie 1991). This leads to recurrent positive selection (fluctuating selection) in favor of the newly beneficial alleles and to adaptive evolution (V. Mustonen and Lässig 2007; Ville Mustonen and Lässig 2009; Eyre-Walker and Keightley 2009; Bengner and Sella 2013; Cvijovic et al. 2015). Nowadays, the way fluctuating selection shapes the dynamics of the relative fitness of the current allele remains poorly studied. Here, we characterize the effects of epistasis and of fluctuating selection on SPFL changes and estimate the contribution of these forces in past evolution.

## Materials and methods

### Multiple alignments of protein-coding sequences

We use multiple alignments of exons of vertebrates and insects from the UCSC Genome Browser database together with the corresponding phylogenies (Figure 5.1) (Rosenbloom et al. 2015). Columns with gaps are excluded. From these alignments, we reconstruct the alleles in the internal nodes of phylogenetic trees with *codeml* (Ziheng Yang 2007). We re-estimate the lengths of individual branches as the average frequency of amino acid substitutions per site on this branch. Based on site-specific  $dN/dS$  ( $\omega$ ) values we classify codon sites as negatively selected ( $\omega < 1$ ), neutral ( $\omega = 1$ ) or positively selected ( $\omega > 1$ ) using Bayes empirical Bayes (BEB) method as implemented in the *PAML* package (Ziheng Yang, Wong, and Nielsen 2005), and use the estimate of  $\omega$  to classify all sites based on the substitution rate. The size of the datasets and the number of sites and substitution subtrees in each bin are shown in Table 5.1. The mitochondrial dataset consists of the amino acid alignment of five proteins for several thousand metazoan species (Klink and Bazykin 2017).



**Figure 5.1. The phylogenies used for the inference of current allele fitness change.** The phylogenies of 53 species of vertebrates (a) and of 24 species of insects (b) from UCSC Genome Browser database. The branch lengths are given in *dS*.

## Simulations of amino acid evolution on dynamic landscapes

To perform simulations of amino acid sequence evolution we use the *SELVa* simulator (Nabieva and Bazykin 2019). *SELVa* is a forward-time Markov chain simulator that allows the user to model sequence evolution along a predefined phylogenetic tree on static or dynamic SPFLs. The user can specify both the shape of SPFL (i.e., the vector of allele fitnesses for a single position in the genome) and the rule for its change. In this work, we use three types of SPFLs for amino acid sites with 20 possible alleles (*SELVa* doesn't support codon models): flat SPFL corresponding to neutral sites (no substitution leads to change of fitness, log fitness vector is  $(0, 0, \dots, 0)$ ); rugged SPFL (one allele is highly preferable over the other ones, log fitness vector is  $(10, 0, \dots, 0)$ ) and gamma-distributed SPFL, where the log fitness values for alleles are randomly chosen from the gamma distribution with user-defined parameter (shape = rate =  $\alpha$ ); fitnesses used by *SELVa* are relative: since they are defined as log fitness, simulation results remain the same if we add any number to vector element; the vector  $(0, 0, \dots, 0)$  is equal to  $(x, x, \dots, x)$  for any  $x$ .

We use two modes of SPFL change. In the random change mode, the SPFL changes are a Poisson process with a user-defined rate. In this case, fitness values are either reshuffled between alleles (for the flat or rugged SPFL) or redrawn from the same distribution (for gamma-distributed fitnesses). In the current allele-dependent mode, the log fitness of the current allele increases or decreases linearly with time. The user can define the rate of this change ( $k$ ) and the length of the time interval between changes ( $\Delta t$ ). The log fitness of the current allele  $B$  at time  $t + \Delta t$  is then set to

$$f_B(t + \Delta t) = f_B(t) + k * \Delta t.$$

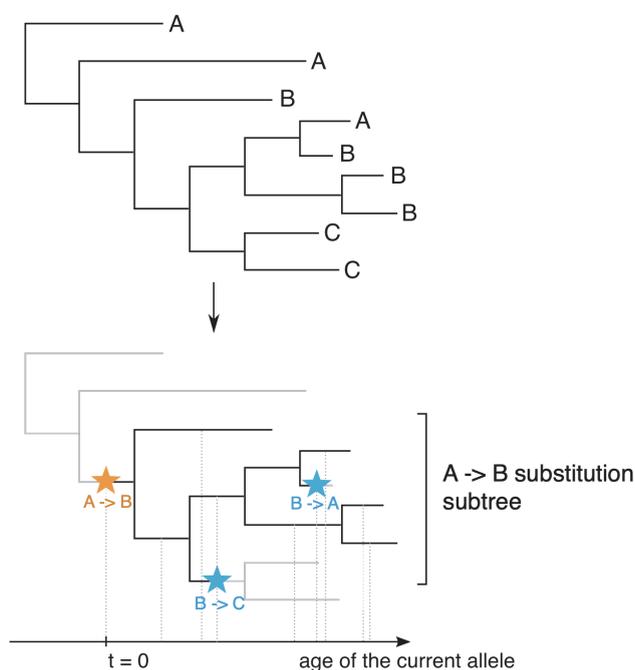
The fitness values for other alleles remain unchanged. Positive values of  $k$  correspond to entrenchment of the current allele, which means its fitness increases with time, and negative ones, to senescence, so that its fitness decreases. When a substitution occurs, i.e. the current allele is replaced with another one, the fitness of the replaced allele stops changing, and the fitness of the new allele starts to change at rate  $k$ . In this work, we used  $\Delta t = 0.01 dS$ , where  $dS$  is the length of time required for a single substitution at a neutral site; this is small enough to simulate the gradual change of allele fitness.

Individual sites are simulated independently, and interactions between them are not modeled directly. We also do not explicitly model the heterogeneity in substitution rates between alleles, although such heterogeneity arises from the differences in the fitness values of individual alleles drawn from the underlying distribution.

For ABC-based inference of the dynamics of the current allele fitness in the evolution of real-life protein sequences, we used the corresponding phylogenies of vertebrates and insects (Figure 5.1).

### Substitution subtrees

For every replacement  $A \rightarrow B$  on any internal branch of the phylogeny, we can define the corresponding substitution subtree, namely, the contiguous segment of the phylogeny where every internal node carries the derived variant  $B$  (Figure 5.2). Within a substitution subtree, the current allele  $B$  can be replaced by the ancestral variant  $A$  or some other variant  $C$  in the course of allele loss(es).



**Figure 5.2. Substitution subtrees.** For every individual amino acid site, ancestral state reconstruction can be used to infer allele substitution history at this site. For every replacement  $A \rightarrow B$  (orange) on any internal branch of the phylogeny, we can define the corresponding substitution subtree — the contiguous segment of the phylogeny where every internal node carries the derived variant  $B$ . It can be replaced by the ancestral variant  $A$  or some other variant  $C$  in the course of allele losses (blue).

For such subtree,  $A$  is the ancestral allele and  $B$  is the current allele.

A genomic position can carry no substitution subtrees if it is fully conservative, or carry one or more substitution subtrees; the number of substitution subtrees equals the number of substitutions on the internal phylogenetic branches in this position. This means that rapidly evolving sites carry more substitution subtrees than conservative ones.

We define a statistic  $s_{branch}$  — the frequency at which the allele  $B$  that has occupied the considered genomic position at the origin of a specific branch has been replaced at this branch. If a single site is analyzed (as in the analysis of mitochondrial genes),  $s_{branch}$  can take the values of 0 or 1; if multiple sites are pooled,  $s_{branch}$  for different substitution subtrees are considered separately, and  $s_{branch}$  can also take values between 0 and 1.  $s_{branch}$  is determined by the SPFL and the overall substitution rate of allele  $B$ . If the mutation rate is assumed to be constant, changes in  $s_{branch}$  with time since the origin of  $B$  within the substitution subtree can be used to detect changes in SPFL. We estimate the rate of replacement of  $B$  ( $s_{branch}$ ) as a function of its age, i.e. the evolutionary time since it was gained.

### **Inference of senescence or entrenchment for groups of alleles**

An individual substitution subtree usually does not provide enough data to infer SPFL changes. To identify such changes with confidence, we have to pool data across subtrees and sites. However, pooling data on different subtrees creates a spurious signal of entrenchment due to heterogeneity of evolution rate and unevenness of SPFLs (see Results) (Naumenko, Kondrashov, and Bazykin 2012; McCandlish, Shah, and Plotkin 2016). One approach to adjust for these confounding factors is to estimate the mean substitution rate of a subtree, which combines the mutation rate of the site and the fitness of the current allele and to use this value for model fitting, *e.g.* in the maximum likelihood (ML) framework. However, our phylogenies are not deep enough to perform ML estimates: the number of substitutions per subtree is too low, while the variance of branch lengths is too large. Instead, to account for confounding effects, we use the approximate Bayesian computation approach (ABC).

The age-dependent patterns of substitutions are sensitive to data heterogeneity and the shape of SPFLs (for details, see Results), so we can't directly measure the rate at which the fitness of the current allele changes. To estimate the strength and abundance of senescence and entrenchment in the evolution of protein sequences, we use rejection ABC with ridge regression adjustment as implemented in the *abc* package for R (Csilléry, François, and Blum 2012). ABC is a popular method used for parameter inference if the likelihood function is not known, using simulations to infer the posterior distributions of estimated parameters (Csilléry et al. 2010).

To produce the simulations for ABC prior we also use *SELVa*. Importantly, rather than simulating the full phylogenetic trees, we simulate individual substitution subtrees. For each dataset of interest, we extract the list of subtrees generated by substitutions (allele gain events) in this dataset. For each substitution subtree, we then run *SELVa* with the given parameters, assuming that the number of sites in the simulation equaled the number of cases when this subtree appeared in the data. The results are pooled across subtrees, and summary statistics were calculated. This approach has two advantages in comparison to simulations based on full phylogenetic trees. First, our summary statistics are based on subtrees only, and using the list of substitution subtrees from the data is more informative than simply the number of sites: this way, we don't have to wait until the ancestral substitution occurs in a simulation, but can start the simulation at the moment we know it has occurred in the data. Second, since the subtrees are smaller, the simulations run faster.

We use two model functions for ABC. The first one is based on the assumption that all sites in the dataset are susceptible to senescence or entrenchment of the same strength (two-parameter model). It requires two parameters: *alpha* rate parameter for the gamma distribution of alleles' fitness values (as described above) and the rate of change of the fitness of the current allele *k*.

The second model represents a mixture of two categories of sites: those with a static SPFL ( $k = 0$ ) and those under senescence ( $k < 0$ ) or entrenchment ( $k > 0$ ). It takes three parameters as input: in addition to *alpha* and *k*, it uses the fraction of substitution subtrees (which corresponds to the fraction of alleles) under senescence or entrenchment with rate *k*. The simulated values for *k* were distributed uniformly from

-100 to 100; the fraction of alleles under senescence or entrenchment was also distributed uniformly from 0% to 100%; and  $\alpha$  was distributed log-uniformly from -1.5 to 1.

The number of amino acid sites in the datasets with different site-specific  $\omega$  values, and the number of substitution subtrees at these sites, vary between  $10^3$  and  $10^6$  (Table 5.1). To account for the variance in summary statistics for smaller datasets, the ABC model function takes a list of subtrees and their counts in the given dataset as input and generates simulations of the same size (but not larger than 100 000 substitution subtrees due to runtime restrictions). For all datasets, we use ridge regression algorithm for parameter estimation as implemented in the *abc* package.

### Summary statistics

After evaluating a range of possible summary statistics for ABC, we ended up using two statistics based on the dynamics of allele replacement. All branches across all subtrees in the simulation are pooled together and used to calculate the following linear regression:

$$s_{branch} = a * length_{branch} + b * age_{branch} + c ,$$

where  $s_{branch}$  is the frequency at which the current allele is lost on the given branch,  $length_{branch}$  is the average length of this branch across all substitution subtrees, and  $age_{branch}$  is the age of the current allele, i.e. the distance from the root of the substitution subtree to the branch. As summary statistics, we use the values of  $a$  and  $b$ .

### ABC validation

We validated the ABC pipeline for parameter inference using *SELVa* simulations based on the reconstructed phylogeny of 53 vertebrates and the *abc* package for  $R$  (Csilléry, François, and Blum 2012). To cross-validate ABC performance under different tolerance rates and to evaluate the accuracy of parameter estimation for both two-parameter and three-parameter models, we calculated prediction error for parameters based on 100 randomly chosen simulations with the cross-validation function of the *abc* package. The prior size is  $10^4$  simulations for both models. Cross-validation tests for parameter

inference with our ABC pipeline showed that we can accurately estimate the parameters of both models using the selected set of summary statistics. Based on their results, we selected the tolerance level of 0.01 for both models.

Next, we asked whether our method is sensitive to changes in the overall rate of evolution. For each model, we generated the testing set of 100 simulations with randomly chosen parameters with normal and twofold increased substitution rate and then used ABC to infer the parameters. We demonstrate that, although the magnitude of  $k$  was overestimated for simulations with accelerated evolution rate, the estimates were not biased in any direction (t-test p-value = 0.53). While the fraction of senescing or entrenched alleles in the three-parameter model was overestimated for simulations with accelerated evolution rate, the magnitude of the bias was not large (on average 0.10, t-test p-value =  $3e-10$ ).

We also checked whether our method allows us to confidently distinguish between senescence and entrenchment. Applying the method to the same testing set of simulations shows that the frequency of misclassification is 0% for the two-parameter model and 1% for the three-parameter model, and cases of misclassification were only observed in simulations with low  $k$  ( $< 1$ ). Furthermore, in the few erroneously classified cases, the 95% probability interval for  $k$  overlapped with zero.

*SELVa* stores the sequences of internal nodes of the phylogenetic tree (the ancestral sequences) so that the history of simulated amino acid replacements is known exactly. However, for real data, we use *codeml* to reconstruct the ancestral sequences, and this reconstruction can be erroneous. To make sure that the ancestral state reconstruction does not affect the accuracy of parameter estimation, we reconstructed the ancestral sequences generated by *SELVa* on the basis of the sequences of terminal nodes in the same way as it was done for the actual data, and used ABC to estimate the parameters using the same procedure as above. We found that ancestral states reconstruction slightly biased both  $k$  (the difference between the true and estimated values  $\sim 3.8$ , t-test p-value =  $6e-4$ ) and the fraction of alleles with changing fitness (by  $\sim 0.08$ , t-test p-value  $< 2e-16$ ) upwards, but the confusion frequency remained low (0% for the two-parameter model, 0.5% for the three-parameter model), and the only erroneously classified simulation had a low fraction of entrenched alleles (0.09).

We also used *evolver* to simulate datasets under different modes of selection to test whether the artifactual signal of senescence or entrenchment can occur in the stationary model of evolution (Ziheng Yang 2007) (Figure 5.8a,c).

The simulated prior distributions used in the current study, the summary statistics calculated from the genomic datasets of vertebrates and insects and source code for the analysis are available at <https://github.com/astolyarova/senescence-ABC>.

## Results

### **Environmental fluctuations decrease the fitness of the current allele**

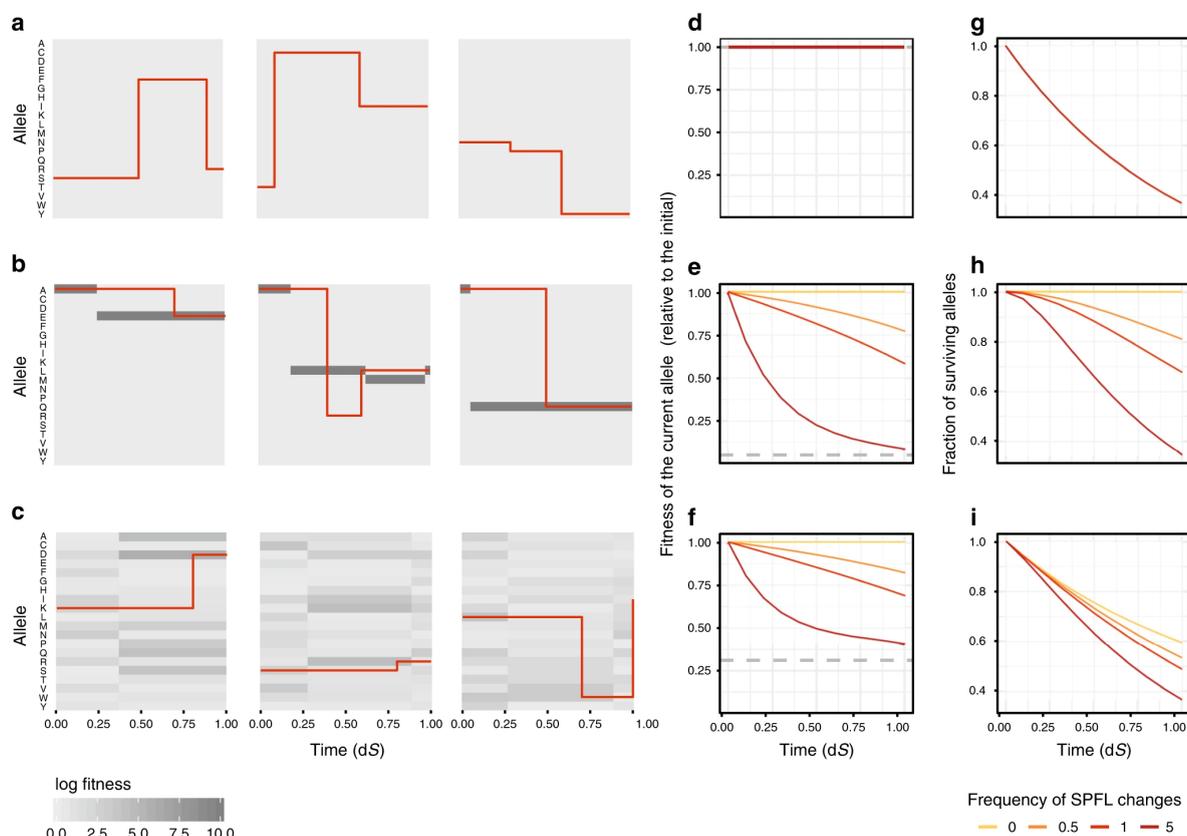
First, we ask how fluctuating selection affects the relative fitness of different alleles at a site. If changes of the SPFL are random with regard to the identity of the allele currently residing at the site, we expect that they, on average, will reduce its relative fitness. This is because, in mutation-selection equilibrium, the relative fitness conferred by the current variant is, on average, higher than that of a random variant at this site. An episode of positive selection triggered by this change may then cause the spread of a novel variant which would confer high fitness till the next SPFL change.

To illustrate this, we simulate amino acid evolution on a randomly changing fitness landscape. In this simulation, the fitness values for each of the 20 possible amino acids are drawn from a predefined distribution, and the amino acid substitutions occur with probabilities determined by the corresponding selection coefficients. At random moments of time, fitness values are redrawn from the same distribution (Figure. 5.3a-c).

As a result of selection, the fitness of the current allele is on average higher than that of other alleles (Figure 5.3b-c); in particular, if selection is strong, the site is typically occupied by the best-possible allele (Figure. 5.3b). However, as the landscape changes randomly, the fitness of this original allele, on average, decreases with time, gradually approaching the mean fitness across all possible variants (Figure. 5.3e,f). We call this process senescence of the current allele (Popova et al. 2019). This effect is more pronounced for the rugged landscape, when one allele is highly more beneficial than others (Figure. 5.3e), and less pronounced when selection is weaker (Figure. 5.3f).

The decline in fitness of the current allele due to fluctuating selection leads to an increase in the rate at which it is lost (Figure. 5.3h,i), in line with the quenched theory of fluctuating selection (V. Mustonen and Lässig 2007; Ville Mustonen and Lässig 2008). We assume that most fluctuation-induced substitutions occur when the SPFL change frequency is lower than the rate of evolution or comparable to it, so our model is still suitable to study evolution under fluctuating selection: if the fluctuations in the SPFL are very rapid, the resulting landscape will be “quasi-neutral”. In this case, the substitution

rate will be reduced, and not increased, by further increase in the fluctuation rate, ultimately reaching the neutral value (J. Gillespie 1973; M. Kimura 1954; Takahata, Ishii, and Matsuda 1975; V. Mustonen and Lässig 2007).



**Figure 5.3. Random changes of SPFL reduce the fitness of the current allele.**

(a–c) Examples of how simulated random changes in SPFLs of different shapes provoke allele substitutions. Each of the nine plots shows the history of one simulated amino acid site; red lines represent the current allele at the site, with vertical red lines indicating substitutions. SPFL changes randomly at the average rate of one change per time required for one neutral substitution ( $1 dS$ ). (a) All 20 possible alleles have the same fitness (flat SPFL), so that all substitutions are neutral. (b) One allele is substantially more beneficial than others (rugged SPFL); most of the observed substitutions are positively selected. (c) Log fitness values are drawn from a gamma distribution. (d–f) Changes in the average fitness of the current allele with evolutionary time under random SPFL changes. The mean fitness across all possible alleles is shown with a dashed line. (g–i) The fraction of surviving ancestral alleles as a function of time since the beginning of the simulation. (d,g) flat SPFL, (e,h) rugged SPFL, (f,i) gamma-distributed SPFL. For d–i, 95% confidence bands based on ten repeats are plotted (but too narrow to be seen).

## Senescence and entrenchment result in opposite substitution patterns

Therefore, the two different modes of change of the SPFL are expected to produce the opposite dynamics of the fitness of the allele that currently occupies the site. If the current allele is favored by epistatic interactions with other sites, it will be entrenched, i.e. its fitness, compared to that of other alleles at this site, is expected to increase with time. By contrast, random SPFL changes that occur without regard to the identity of the allele currently occupying the site are expected to decrease its fitness, leading to senescence.

We propose that this dichotomy can be used to distinguish between these two modes of SPFL changes. To infer the changes in the relative fitness of an allele with time, we study the differences in the rate at which it is lost in the course of evolution. Indeed, the relative fitness of an allele specifies the probability that it is substituted by another allele per unit time (Motoo Kimura 1983).

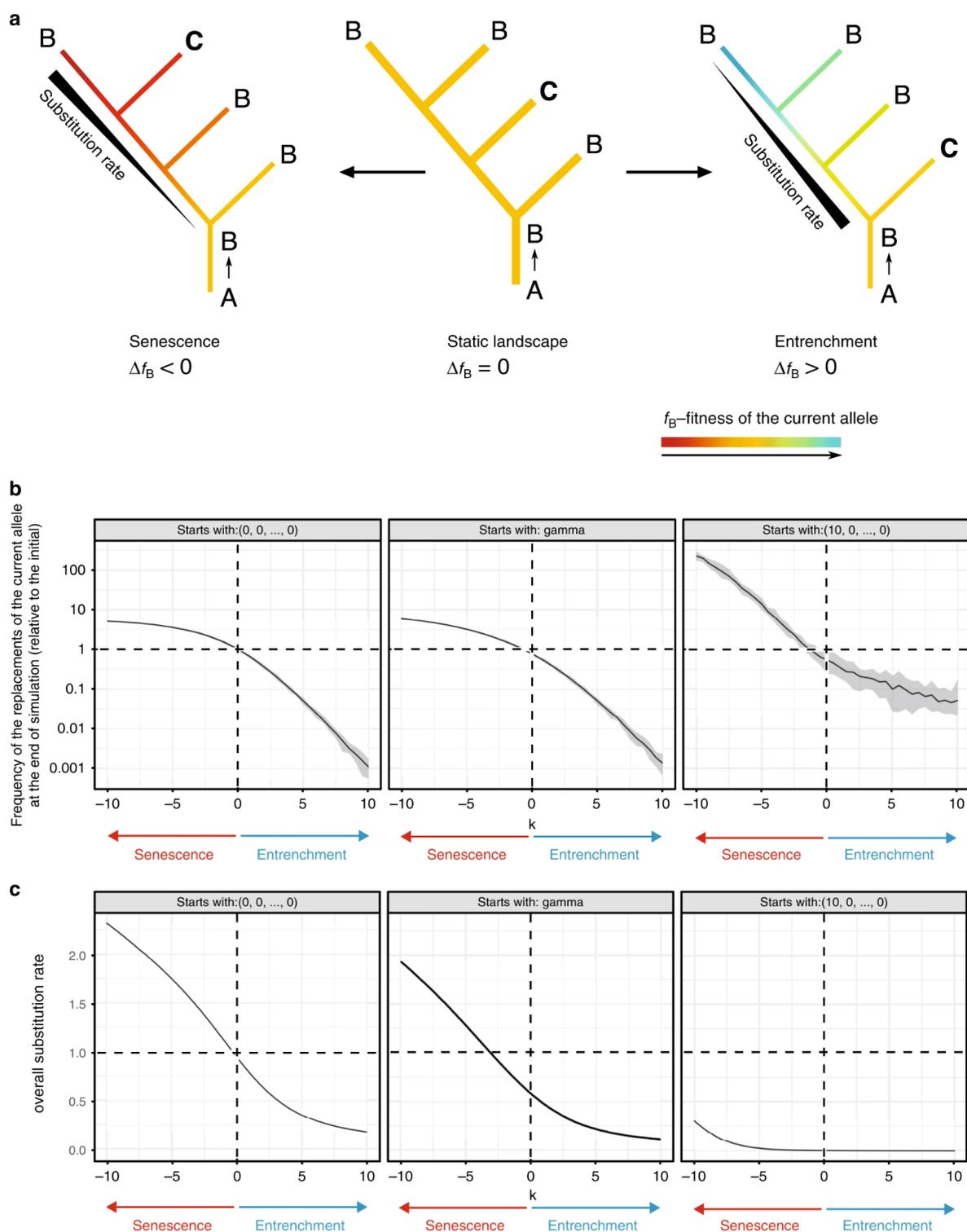
Let us assume that a substitution of an ancestral variant  $A$  for another variant  $B$  (allele gain) has occurred at some internal branch of the phylogenetic tree, and this current allele  $B$  has been preserved in several extant species (Figure 5.4a). In other species, it could be lost, for example, as a result of a reversal to  $A$  or a substitution for some other allele  $C$ . If the SPFL for this site has remained static (the fitness of the current allele  $B$  has not changed,  $\Delta f_B = 0$ ), the probability of replacement of  $B$  per unit time is independent of the time elapsed since its gain.

Under senescence, the fitness conferred by  $B$  decreases with its age ( $\Delta f_B < 0$ ), and the probability of its replacement increases with it. In this case, we will observe a higher rate of substitutions on the branches originating much later than the allele gain, compared to the branches leading to close descendants (Figure 5.4a, left). By contrast, under entrenchment, the fitness of the current allele increases ( $\Delta f_B > 0$ ), so the rate at which  $B$  is lost declines with its age (Figure 5.4a, right).

To test the validity of this approach, we simulate molecular evolution at individual sites assuming that the fitness of the allele currently residing at the site changes with time. Specifically, we assume that the log fitness of the current allele is initially drawn from a predefined distribution, and then changes with time linearly with rate  $k$ . Positive values

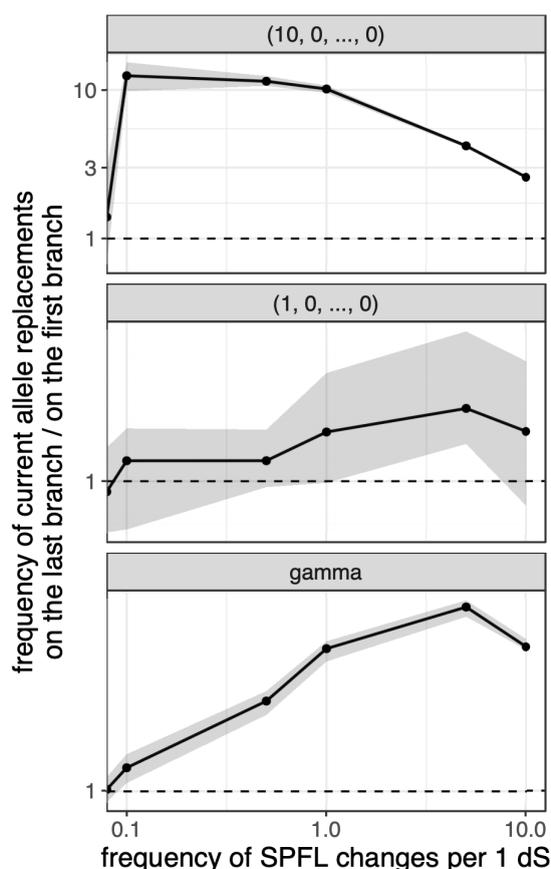
of  $k$  correspond to an increase in the fitness of the current allele, *i.e.*, entrenchment, while negative values correspond to a decrease in its fitness, *i.e.*, senescence. As expected, entrenchment ( $k > 0$ ) results in a high rate of substitutions immediately after the ancestral substitution, but a reduced rate later on. By contrast, under senescence ( $k < 0$ ), the rate of substitutions increases with time since the allele gain (Figure 5.4b). An alternative mode of simulation of senescence, whereby random changes in SPFL and the molecular evolution caused by them are modeled explicitly, gives the same results (Figure. 5.5).

Besides the phylogenetic distribution of substitutions, the mode and rate of SPFL change also affect the overall rate of molecular evolution (Figure 5.4c). Compared to a static landscape of the same shape, entrenchment reduces the substitution rate, as the time-averaged fitness of the current allele is higher, and therefore it is replaced less frequently. Conversely, under senescence, many of the substitutions of the current allele are advantageous, increasing the overall rate of evolution. Importantly, senescence doesn't necessarily result in an overall evolution rate exceeding the neutral rate, which is a hallmark of positive selection. Indeed, if an allele is strongly preferred, a drop in its fitness over the course of senescence may still leave it the optimal one, so that negative selection will still maintain it (as observed for the rugged SPFL, Figure 5.4c right).



**Figure 5.4. Replacement patterns of the current allele reflect changes in its fitness.** (a) On the static landscape, the probability that B is replaced per unit time does not depend on the time since its gain  $A \rightarrow B$ . Under entrenchment, B becomes more favorable with time ( $\Delta f_B > 0$ ); therefore, the  $A \rightarrow B$  substitution rate declines and there are fewer substitutions observed on “late” branches of the phylogeny. Under senescence, the fitness of B decreases ( $\Delta f_B < 0$ ), leading to an increase in the rate of its

loss with time. **(b, c)** In simulated evolution, changes in fitness of the current allele affect the dynamics of its replacements (calculated as the ratio between  $B$  substitution rate at the terminal branch of the substitution subtree to  $B$  substitution rate on the root branch, **b**) and the overall substitution rate on the whole tree (**c**). Simulations were started with SPFLs of different shapes: flat SPFL, rugged SPFL and gamma-distributed SPFL. Over the course of simulation, the log fitness of the current allele was linearly changing with time at rate  $k$ . For **b, c**, mean values and 95% confidence bands based on ten repeats are shown.

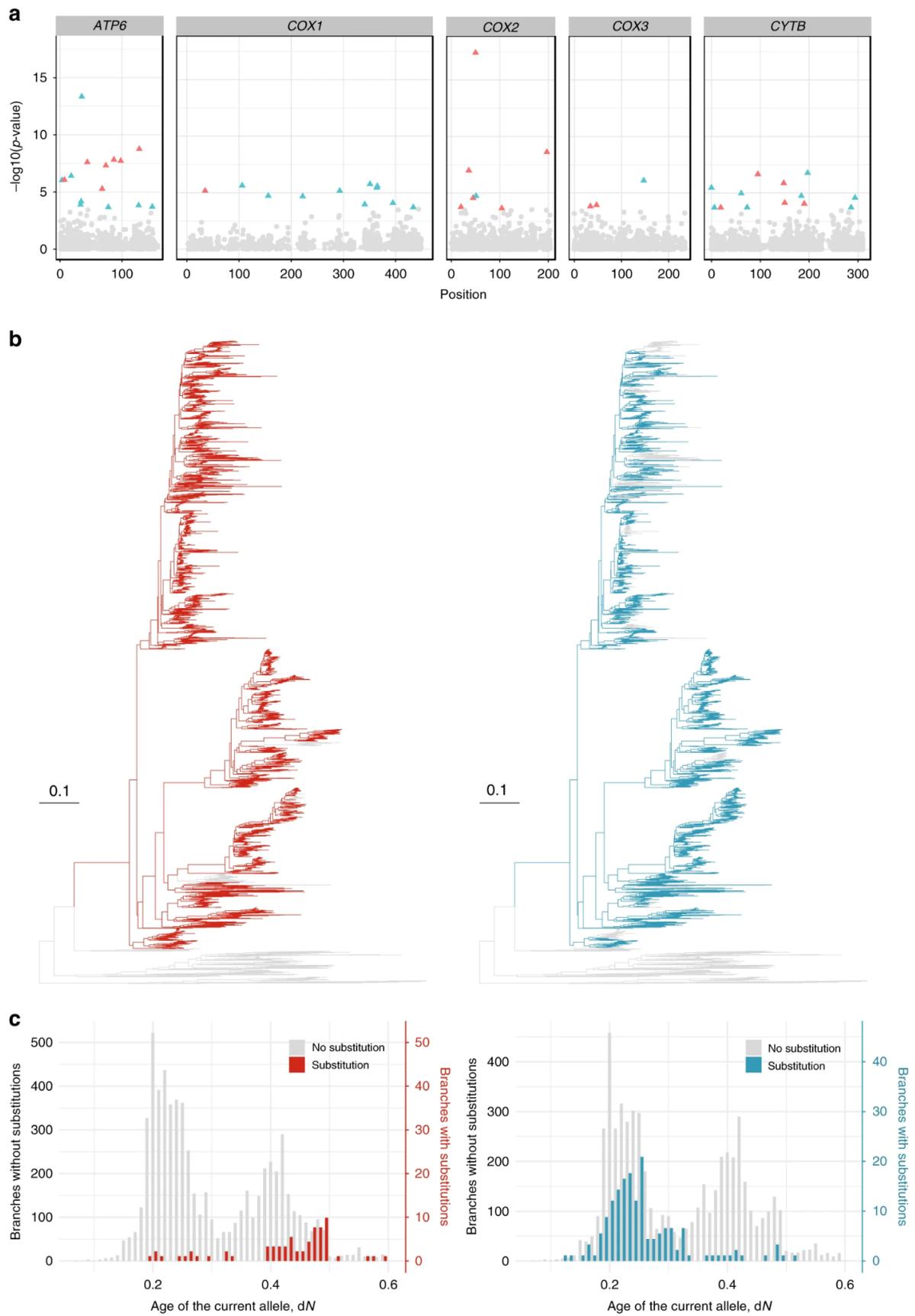


**Figure 5.5. Growth of the rate of the current allele replacements since its origin driven by random SPFL changes.** For the static SPFLs of different shapes, i.e. if initial fitness values don't change with time, the frequency of the replacements of the current allele remain the same for "early" and "late" branches. If fitness values are redrawn from the same distribution with some frequency comparable with the evolution rate (*i. e.* SPFL changes randomly with time), the rate of current allele replacements on the "late" branch increases. The mean values and 95% confidence bands obtained in 10 simulation repeats are shown.

## Senescence and entrenchment at single-allele resolution

Large phylogenies allow detecting changes in substitution frequencies for individual alleles. Each originating allele, e.g. an amino acid arising at a specific site from an ancestral amino acid substitution, can be inherited by multiple descendant lineages leading to different extant species. Ancestral state reconstruction can then be used to infer the lineages at which this allele has been lost due to a reversion or substitution to a different amino acid. If enough such lineages are available, this allows us to trace the decline or increase in the rate of allele substitution since its origin, *i.e.*, entrenchment or senescence.

We apply binomial logistic regression to detect changes in substitution frequencies with the age of the current allele along the phylogeny for five mitochondrial genes of Metazoa (Klink and Bazykin 2017). The regression is performed separately for each allele  $B$  with a known time of origin (corresponding to allele gain  $A \rightarrow B$ ) at each site. Among the 42,637 such alleles, we identified 28 alleles for which the frequency of replacement significantly increased with time since their origin (*i.e.* senescing alleles), and 21 alleles where it decreased (*i. e.* entrenched alleles) at 5% false discovery rate (Figure 5.6a, Table A3). The examples of phylogenies indicating allele replacements at senescing and entrenched alleles are shown in Figure 5.6b-c. Despite the opposite time-dependent dynamics of the substitution rate in the sites containing entrenched and senescing alleles, the overall number of substitutions occurring along the phylogeny is similar between them (sign test  $p$ -value = 0.84), indicating that the overall substitution rate is insufficient for distinguishing between these scenarios.



**Figure 5.6. Senescence and entrenchment of individual alleles in the mitochondrial genes of Metazoa.** (a) Manhattan plot of senescing and entrenched alleles. Only the alleles with a known phylogenetic position of origin, i.e., those that were not yet present in the tree root, were analyzed; a single genomic site can contain zero, one or several alleles. *P*-values are calculated using binomial logistic regression. The alleles demonstrating significant senescence under 5% FDR are shown in red; the alleles demonstrating entrenchment are shown in green. No amino acid sites contained more than one significantly senescing or entrenched alleles. (b) Examples of senescing (COX2 position 56, red) and entrenched (ATP6 position 71, blue) alleles. The contiguous segment of the phylogeny carrying the derived allele is shown in color. (c) Distribution of substitutions along the lifetime of alleles shown in (b). For the senescing allele, the phylogenetic branches corresponding to allele replacements (red) originate later than the branches without replacements (gray). Conversely, the entrenched allele is more frequently replaced soon after its origin (green).

## Heterogeneity of alleles leads to an artifactual signal of entrenchment

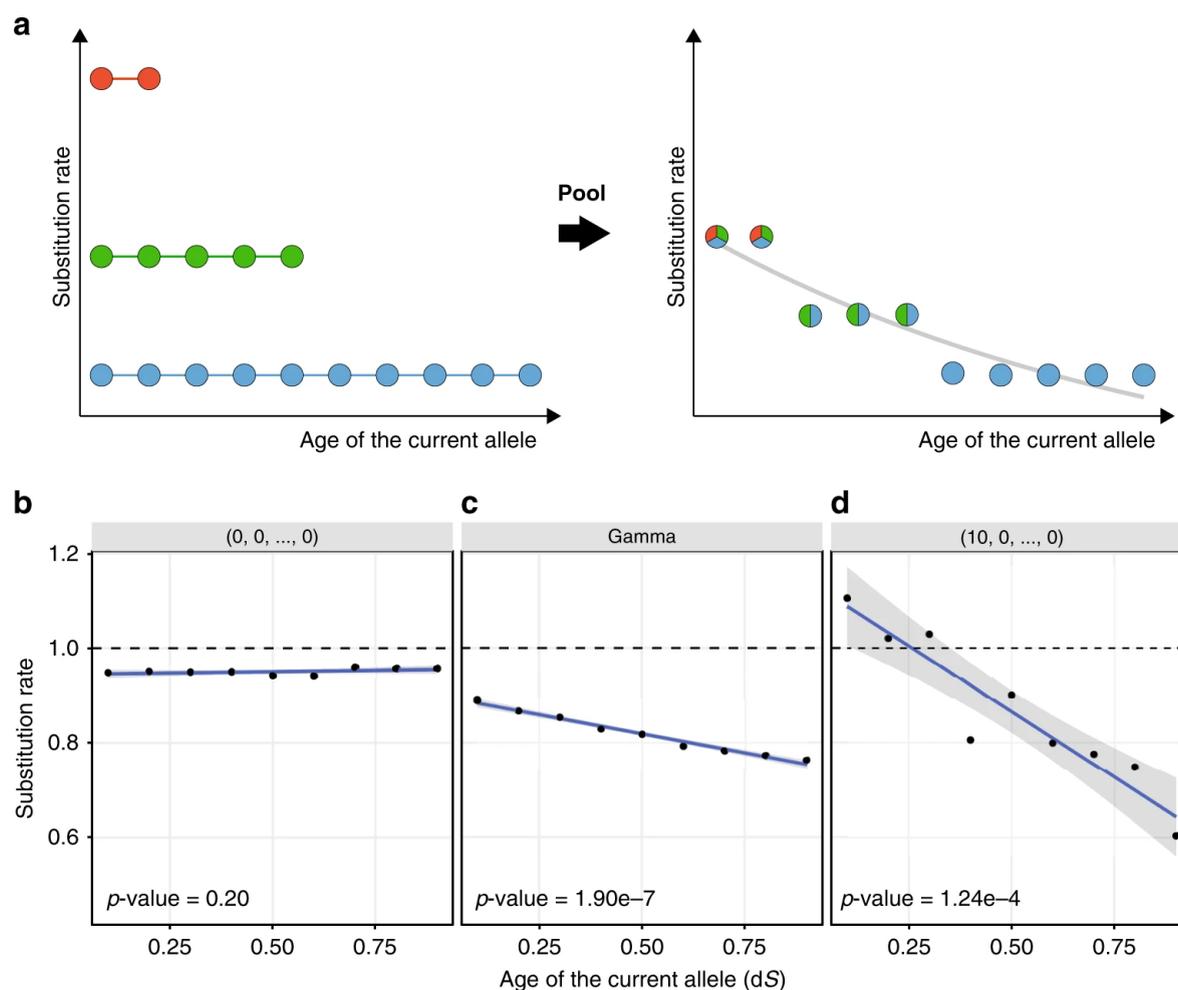
While phylogenies spanning hundreds and thousands of species, like those available for mitochondrial proteins, allow us to measure the changes in the substitution rate for individual alleles, in smaller phylogenies, the number of substitutions experienced by an allele can be insufficient for such an analysis. Still, it may be possible to identify the prevailing patterns of substitutions by pooling alleles together. However, such pooling can be problematic: even in the absence of SPFL changes, the rate of substitution can appear to change with time since allele origin if the pooled alleles have different time-invariant substitution rates, confounding inference of SPFL changes.

Indeed, consider a set of alleles, each characterized by its own substitution rate that is stationary (constant in time) but differs between alleles. While the replacement rate may be constant for each allele, so that the time to replacement is characterized by an exponential distribution, it will not, in general, be exponentially distributed in the resulting heterogeneous dataset. Instead, the frequency of substitution will appear to decline with time (Figure 5.7a), making it non-stationary and mimicking entrenchment of the current allele. The problem of data heterogeneity leading to decreasing hazard function is well known in demographic inference (Proschan 1963; Vaupel, Manton, and Stallard 1979), and has been previously appreciated in the inference of substitution rates dynamics in molecular evolution (Naumenko, Kondrashov, and Bazykin 2012;

McCandlish, Shah, and Plotkin 2016). Notably, no mixture of stationary processes can give rise to an increase in the substitution rate, i.e., senescence (Proschan 1963).

It is obvious that heterogeneity of substitution rates arises from pooling of different amino acid sites with varying substitution rates. More subtly, it also arises within individual sites as a result of differences between rates of substitution of different alleles. Substitution rate is the product of mutation rate and fixation probability, and this heterogeneity will arise due to any differences in either of these factors between alleles. For example, consider a single site which is non-neutral, i.e., such that different alleles confer different fitness. Such alleles will be characterized by different replacement rates (lower for high-fitness alleles, and higher for low-fitness alleles), and pooling over different alleles over the course of evolution of this site (or, identically, over different independent and identically distributed sites) would lead to heterogeneity of substitution rates and to an apparent decline in substitution rates with time since allele origin.

To show this, we simulate molecular evolution on static SPFLs of different shapes. If all alleles have the same fitness, i.e., if all substitutions are neutral (“flat” landscape), the substitution rate is independent of time since allele origin (Figure 5.7b). By contrast, if the fitness values of alleles are drawn from a gamma distribution, so that different alleles have different fitness, the substitution frequency decreases with the age of the current allele, even though the SPFL doesn’t change (Figure 5.7c). On a more rugged SPFL, when one allele is much more fit than all others, this decline is even sharper (Figure 5.7d).



**Figure 5.7. Heterogeneity of substitution rates among alleles on static SPFLs imitates entrenchment.** (a) Consider three classes of alleles, characterized by different constant substitution rates: fast (red), moderate (green), and slow (blue) alleles. For each allele, we can calculate the substitution frequency on the phylogenetic branches located at different evolutionary distances from the gain of that allele. If the dynamics of replacement of these alleles are analyzed separately for each substitution-rate class, no spurious signal of entrenchment or senescence is observed (left). However, if alleles from different classes are pooled together, the substitution frequency appears to decrease with time, mimicking the signal of entrenchment (right). (b) No artificial signal of entrenchment is observed on a static flat SPFL. (c) On a gamma-distributed SPFL, the heterogeneity of fitness of the alleles produces entrenchment-like decline of replacement rate of the current allele with time, although the SPFL remains static. The effect is even more pronounced on a more rugged SPFL (d). The p-values are obtained with linear regression.

## **Inferring senescence and entrenchment from phylogenetic distribution of substitutions**

To address the problem of heterogeneity of alleles, we account for differences between alleles in mutation rates and “baseline” selection while inferring the SPFL dynamics for pooled datasets. In the absence of prior information about the distribution of these characteristics, it is impossible to reconstruct the explicit likelihood function for the substitution rates. Instead, we used the approximate Bayesian computation (ABC) (Pritchard et al. 1999) approach to obtain the posterior distribution of the rate of current allele fitness change per unit time  $k$  (positive values of  $k$  corresponding to entrenchment, and negative to senescence). ABC depends on a set of summary statistics to evaluate the difference between the simulation results and the data. We use two summary statistics, each aggregating over all individual alleles, which reflect the age-dependent dynamics of substitution rates (see Methods).

We use two models for parameter inference. Under the two-parameter model, we assume that log fitness values for individual alleles were drawn from a gamma distribution with rate and shape parameters denoted by  $\alpha$ , and the log fitness of the current allele at all sites changed linearly with rate  $k$ . Under the three-parameter model, the fitness changed linearly only for a fraction of alleles, while the fitness of the remaining alleles was invariant.

Simulations show that both models perform well in identifying senescence and entrenchment under a broad range of parameters, and are robust to overall substitution rate, phylogeny shape, pooling of sites with diverse characteristics and errors in ancestral state reconstruction (see Methods for details).

### **Positively selected sites show strong senescence**

We apply the developed ABC approach to protein sequences of vertebrates and insects (Figure 5.1). To understand how the direction of fitness change depends on the overall conservation of an amino acid site, in both datasets, we roughly classify all codon sites by the type of selection acting at them, on the basis of the ratio of nonsynonymous and synonymous substitutions per site ( $dN/dS$ , or  $\omega$ ): negatively selected ( $\omega < 1$ ), neutral

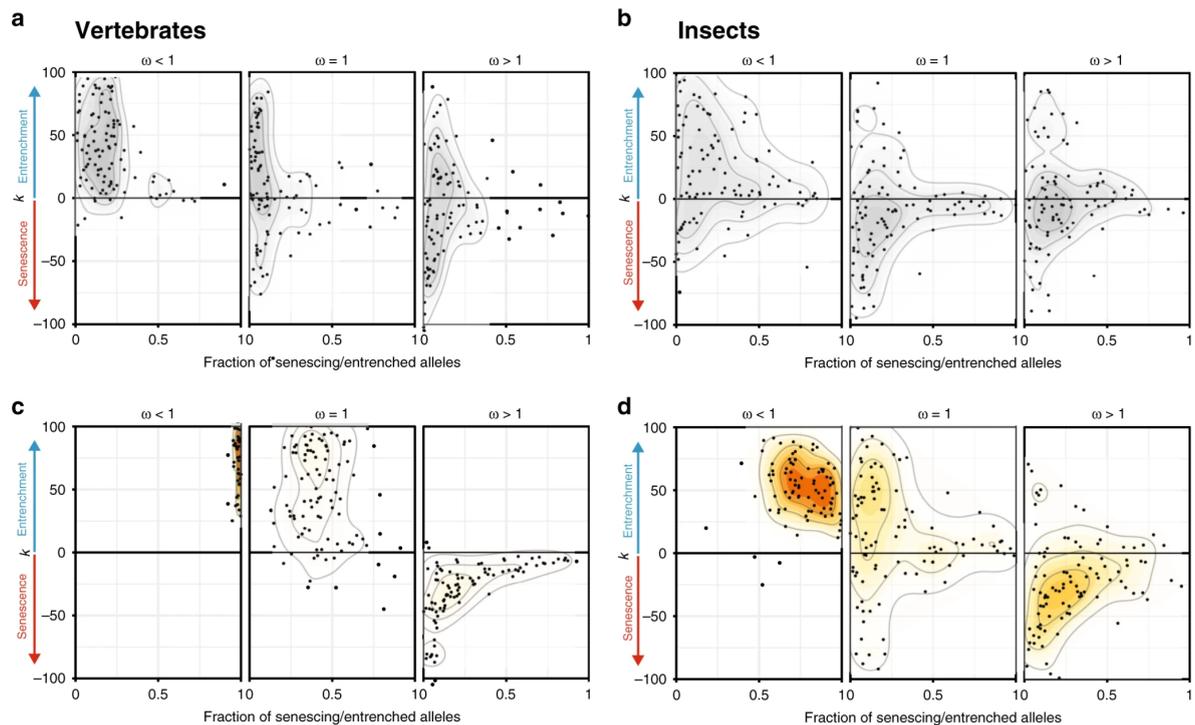
( $\omega = 1$ ) or positively selected ( $\omega > 1$ ) sites, and analyze them independently (Table 5.1). The three-parameter model provides a better fit to the data than the two-parameter model (posterior probability of the three-parameter model  $> 0.74$  for all datasets except sites with  $\omega < 1$  in insects with posterior probability of 0.46), so we use the former for the analysis.

|              | # of sites              | # of alleles            | $k$                 | fraction of senescing or entrenched alleles |
|--------------|-------------------------|-------------------------|---------------------|---|
| Vertebrates  |                         |                         |                     |   |
| $\omega < 1$ | 2 598 701 (87.8%)       | 574 137 (36.4%)         | 62.3 (25.1, 100.0)  | 0.99 (0.95, 1.00)                           |
| $\omega = 1$ | 348 047 (11.8%)         | 917 340 (58.1%)         | 45.3 (-36.8, 98.9)  | 0.42 (0.15, 0.90)                           |
| $\omega > 1$ | 13 189 (0.4%)           | 86 852 (5.5%)           | -23.3 (-94.7, -3.5) | 0.18 (0.0, 0.86)                            |
| <b>Total</b> | <b>2 959 937 (100%)</b> | <b>1 578 329 (100%)</b> |                     |   |
| Insects      |                         |                         |                     |   |
| $\omega < 1$ | 2 699 432 (93.2%)       | 429 800 (48.8%)         | 47.8 (-1.1, 87.9)   | 0.81 (0.47, 1.00)                           |
| $\omega = 1$ | 185 829 (6.4%)          | 413 101 (46.8%)         | 9.4 (-80.2, 91.4)   | 0.16 (0.01, 0.92)                           |
| $\omega > 1$ | 9 698 (4.4%)            | 38 577 (4.4%)           | -22.6 (-88.7, 52.0) | 0.25 (0.03, 0.92)                           |
| <b>Total</b> | <b>2 894 959 (100%)</b> | <b>881 478 (100%)</b>   |                     |   |

**Table 5.1. The analyzed datasets, the corresponding estimates of the rate of entrenchment ( $k > 0$ ) or senescence ( $k < 0$ ) and the fraction of alleles that experience these processes.** The values show the median of the ABC posterior distribution of parameter values; numbers in parentheses represent the 95% posterior probability intervals.

In vertebrates, both the fraction of senescing or entrenched alleles and the value of  $k$  for them depend on the mode of selection acting at the site. The 36% of alleles originating at negatively selected ( $\omega < 1$ ) sites demonstrate strong evidence for entrenchment: we estimate that all of them are entrenched, indicating that the fitness of the current variant increases with time since its origin (Figure 5.8c left panel; Table 5.1). By contrast, of the 6% of alleles arising at positively selected sites ( $\omega > 1$ ), 18% experience senescence (Figure 5.8c right panel), indicating a decrease in the fitness of the current allele. While we are unable to distinguish robustly between a low fraction of alleles undergoing strong senescence and a high fraction of alleles undergoing weak senescence, the 95% posterior probability interval does not include  $k = 0$ , rejecting stationarity. The neutral sites demonstrate an intermediate signal with little evidence for entrenchment or senescence (Figure 5.8c middle panel). A similar pattern is observed in phylogenies of insects (Figure 5.8d).

While senescence is observed at sites that undergo rapid substitution ( $\omega > 1$ ), it is distinct from an increase in the overall substitution rate. Similarly, while entrenchment is observed at constrained sites ( $\omega < 1$ ), it is distinct from a reduction in substitution rate. To illustrate this, we simulate evolution under different substitution rates but constant SPFL on the phylogenies of vertebrates and insects, using the same distribution of  $\omega$  values as in the real data, and estimated  $k$  using the ABC pipeline (Figure 5.8a,c). No senescence or entrenchment is detected for the datasets simulated using the stationary model under  $\omega = 1$  or  $\omega > 1$ . A weak spurious signal of entrenchment is detected for the simulated datasets with  $\omega < 1$ , resulting from the high heterogeneity in evolution rates; however, it is much weaker than that observed in the data for this category of sites.



**Figure 5.8. Senescence and entrenchment in protein sequences of vertebrates and insects.** Plots provide ABC estimates of the rate of senescence or entrenchment  $k$  and the fraction of alleles with changing fitness for protein sequences of vertebrates and insects. For each dataset, the posterior distribution of parameters under 1% acceptance threshold after local ridge regression adjustment is shown. **(a, b)** Simulated data under the stationary model using the vertebrate **(a)** or insect **(b)** phylogeny and distribution of  $\omega$  values. **(c, d)** In real genomic data of vertebrates **(c)** or insects **(d)**, sites under negative selection show strong entrenchment, neutral sites demonstrate the intermediate signal, and positively selected sites are senescing.

## Discussion

While the direction of changes in fitness in the course of evolution is unpredictable for an individual allele, there are certain statistical regularities. Previous works have shown that, at a site involved in epistatic interactions with other sites, the relative fitness conferred by the incumbent allele as compared to other alleles possible in this site is expected to increase with time since its origin (A. S. Kondrashov, Sunyaev, and Kondrashov 2002; Povolotskaya and Kondrashov 2010; Breen et al. 2012; McCandlish et al. 2013). Acting alone, *i.e.*, if the overall fitness landscape is static, this process of entrenchment should make the propensities existing at individual sites more pronounced. This, in turn, should limit the level of divergence between highly divergent sequences, although reaching this level may take a very long time (F. A. Kondrashov and Kondrashov 2001b; Weinreich, Watson, and Chao 2005; Weinreich et al. 2006; Povolotskaya and Kondrashov 2010; A. S. Kondrashov et al. 2010; Ferretti et al. 2018).

Here, we consider the dynamics of allele fitness due to changes in the overall fitness landscape itself. We show that, if the direction of these changes is independent of the current position of the population in the genotype space, the expected mean dynamics — senescence — is opposite to that of entrenchment. We design a method to distinguish the two patterns from the phylogenetic distribution of substitutions and find that entrenchment is ubiquitous at negatively selected sites, while senescence is prevalent at sites undergoing adaptive evolution under positive selection.

While senescence underlies positive selection in evolution of vertebrates and insects (Figure 5.8), these phenomena are distinct. Indeed, firstly, weak senescence of a highly beneficial allele currently occupying an amino acid site can result in weakening of the negative selection restricting the fixation of other alleles; however, if this allele remains the optimal one, the direction of selection remains the same, and no positive selection in favor of a different allele starts to act (see examples in Figure 5.4b and c, right panel). Secondly, there are models of positive selection that don't imply senescence, including variations of “stairway to heaven” (STH) landscapes without the finite fitness peak (Gerrish and Lenski 1998; Desai and Fisher 2007; Kryazhimskiy, Tkacik, and Plotkin 2009). Importantly, the existing models of positive selection do not imply senescence,

and senescence is not observed in them, as evidenced by the simulated datasets with positive selection and no senescence (Figure 5.8a,c). However, the observed concordance between the direction of non-stationarity and selection mode imply that these phenomena are closely related, and that the ongoing decrease of the fitness of the current allele causes many of the adaptive substitutions.

What causes allele senescence? Firstly, it can sometimes result from negative epistasis with an allele arising at an interacting site. While this could result in senescence occasionally, on average epistasis results in entrenchment (David D. Pollock, Thiltgen, and Goldstein 2012; Shah, McCandlish, and Plotkin 2015). Secondly, senescence could result from changes in selection pressure external to the genome. Previously, the acceleration of substitutions over the course of allele lifetime in the evolution of influenza A virus was described (Popova et al. 2019). This pattern has been mainly observed at sites associated with avoidance of the host immune system pressure, and was interpreted as evidence for negative frequency-dependent selection actively disfavoring the current allele (Popova et al. 2019). Here, we show that this type of selection is not a prerequisite for senescence. Instead, senescence is expected whenever selection changes without regard to the identity of the current allele. How much of senescence, and positive selection resulting from it, is due to random changes of the fitness landscape, and how much is due to systematic selection against the current allele or negative epistatic interactions, can be a subject of further research.

## Chapter 6: Conclusions

In this work, we used comparative genomics methods to analyze patterns of within-population variation and between-species divergence evident on epistasis and fitness landscape changes:

- excess of linkage disequilibrium between rare nonsynonymous alleles within hyperpolymorphic populations of *S. commune*, possibly indicative of epistatic selection maintaining coadapted combinations of alleles;
- short bursts of nonsynonymous replacements between closely related species, caused by correlated positive selection;
- time-dependent changes of the substitution rate in the course of species divergence, which are evidence for epistasis-driven entrenchment of negatively selected alleles and environment-driven senescence of positively selected alleles.

While in almost all species nucleotide diversity is a small parameter  $\ll 1$ , this is not the case for *S. commune* with two randomly sampled genotypes differing by about 20% at silent sites. As a result, patterns of genetic variation in *S. commune* reveal properties of natural selection, not accessible through data on other species, making it a promising model for population genomics studies. In this work we observe pervasive signatures of positive epistasis due to mutual compensation of deleterious effects of individual alleles, particularly pronounced in genomic regions where nucleotide diversity is maintained by balancing selection (Chapter 3). We believe that this is the first evidence of abundant epistasis affecting polymorphism within natural variation.

The totality of patterns observed in the populations of *S. commune* (the excess  $LD_{\text{nonsyn}}$  within genes, correlated LD between shared polymorphisms in two divergent populations, and the increased LD between physically interacting sites) is indicative of positive epistasis favoring combinations of coadapted alleles segregating within this species. In the presence of epistasis, disruption of co-evolved combinations of alleles by recombination is expected to be disadvantageous. Indeed, the patterns evident of

epistasis which we observed in *S. commune* are limited to the short-range interactions (*i.e.* between sites located within the same gene or to a lesser extent between neighboring genes). Similarly, the high density of polymorphisms in *S. commune* makes it possible to detect short haploblocks — signatures of balancing selection maintaining the co-existence of diverged haplotypes in the linked genomic region, covering ~10% of the *S. commune* genome.

While studying within-population variation, we examine a small part of the full fitness landscape covered by a population. By comparing the genomes of diverged species, we address a much larger region of the landscape; moreover, the landscapes by themselves can differ between species due to environmental or ecological changes. Patterns of interspecific differences can reflect the impact of strong selection, which cannot be observed within standing variation. By analyzing the dynamics of genetic differences accumulated in the course of species divergence, we show that selective constraints shaping these differences are not static. Changes of selection can promote adaptive evolution, resulting in faster evolution rates, as well as conserving the current genomic state, slowing the rate of divergence.

Looking at the genomic data of recently diverged species, we conclude that interspecific differences can accumulate non-linearly on short evolutionary scales (Chapter 4). The observed bursts of nonsynonymous replacements found in the phylogenies of Lake Baikal amphipods and primates are evidence of the impact of sudden positive selection. Moreover, fast fixation of such a high number of nonsynonymous mutations under such a short time is not expected if they are selected independently and can be explained by correlated positive selection. The burst-like accumulation of nonsynonymous differences within proteins can be caused by the opening of a new highly epistatic adaptive path which was inaccessible before the landscape change.

By comparing evolutionarily distant species, we can infer the long-term changes in selection guiding the accumulation of genetic differences (Chapter 5). We observe two opposite trends in the evolution of nuclear genomes of vertebrates and insects and in the mitochondrial genomes of Metazoa. The majority of fixed alleles get entrenched, *i.e.* become more favorable with time, while other alleles demonstrate senescence, *i.e.* their fitness declines with time. The oppositely directed dynamics of substitutions allow us to

distinguish between two possible causes of fitness change. Entrenchment is expected to arise as a consequence of coevolution of epistatically interacting sites. On the contrary, senescence may result from random changes of the landscape, *i.e.* environmental fluctuations, or from selection negatively correlated with the current genome content, *e.g.* ecological interactions with other species.

In contrast to the commonly used methods of inferring selection based on the estimation of the overall substitution rates, our method utilizes its derivative by addressing the time-dependent patterns of substitution rate. However, we found that negatively selected sites more often demonstrate entrenchment of the current allele, while alleles under positive selection are senescing. Therefore, we link the mode of selection acting on the site at some moment to the ongoing process of fitness change and thereby to the evolutionary mechanisms of such change.

Based on these observations, we conclude that patterns of genomic differences show the imprint of epistasis on various evolutionary scales. The non-independent evolution of genomic sites reflects the high complexity of fitness landscapes. Given the complexity and instability of the fitness landscapes, it's practically impossible to describe them comprehensively. Moreover, the way the structure of the fitness landscape defines the evolutionary paths in natural populations is non-trivial and is by itself a matter of study. In this work, we look for indirect evidence of epistasis and fitness landscape changes, using methods relying on the average statistics of the patterns of variation. The power of these methods is insufficient to identify specific epistatic interactions between co-evolving sites or draw conclusions about the molecular mechanisms underlying the revealed features of the fitness landscapes. Nevertheless, the genome-wide patterns of variation within natural populations and between species can be used to uncover general properties of fitness landscapes.

# References

- 1000 Genomes Project Consortium, Adam Auton, Lisa D. Brooks, Richard M. Durbin, Erik P. Garrison, Hyun Min Kang, Jan O. Korb, et al. 2015. "A Global Reference for Human Genetic Variation." *Nature* 526 (7571): 68–74.
- Aita, T., H. Uchiyama, T. Inaoka, M. Nakajima, T. Kokubo, and Y. Husimi. 2000. "Analysis of a Local Fitness Landscape with a Model of the Rough Mt. Fuji-Type Landscape: Application to Prolyl Endopeptidase and Thermolysin." *Biopolymers* 54 (1): 64–79.
- Altschuh, D., A. M. Lesk, A. C. Bloomer, and A. Klug. 1987. "Correlation of Co-Ordinated Amino Acid Substitutions with Function in Viruses Related to Tobacco Mosaic Virus." *Journal of Molecular Biology* 193 (4): 693–707.
- Arbeithuber, Barbara, Andrea J. Betancourt, Thomas Ebner, and Irene Tiemann-Boege. 2015. "Crossovers Are Associated with Mutation and Biased Gene Conversion at Recombination Hotspots." *Proceedings of the National Academy of Sciences of the United States of America* 112 (7): 2109–14.
- Arnold, Brian, Mashaal Sohail, Crista Wadsworth, Jukka Corander, William P. Hanage, Shamil Sunyaev, and Yonatan H. Grad. 2020. "Fine-Scale Haplotype Structure Reveals Strong Signatures of Positive Selection in a Recombining Bacterial Pathogen." *Molecular Biology and Evolution* 37 (2): 417–28.
- Ayala, F. J., and C. A. Campbell. 1974. "Frequency-Dependent Selection." *Annual Review of Ecology and Systematics* 5 (1): 115–38.
- Bajić, Djordje, Jean C. C. Vila, Zachary D. Blount, and Alvaro Sánchez. 2018. "On the Deformability of an Empirical Fitness Landscape by Microbial Evolution." *Proc. Natl. Acad. Sci. U. S. A.* 115 (44): 11286–91.
- Bakhtin, Yuri, Mikhail I. Katsnelson, Yuri I. Wolf, and Eugene V. Koonin. 2021. "Evolution in the Weak-Mutation Limit: Stasis Periods Punctuated by Fast Transitions between Saddle Points on the Fitness Landscape." *Proc. Natl. Acad. Sci. U. S. A.* 118 (4).
- Bank, Claudia, Ryan T. Hietpas, Alex Wong, Daniel N. Bolon, and Jeffrey D. Jensen. 2014. "A Bayesian MCMC Approach to Assess the Complete Distribution of Fitness Effects of New Mutations: Uncovering the Potential for Adaptive Walks in Challenging Environments." *Genetics* 196 (3): 841–52.
- Bank, Claudia, Sebastian Matuszewski, Ryan T. Hietpas, and Jeffrey D. Jensen. 2016. "On the (un)predictability of a Large Intragenic Fitness Landscape." *Proc. Natl. Acad. Sci. U. S. A.* 113 (49): 14085–90.
- Bankevich, Anton, Sergey Nurk, Dmitry Antipov, Alexey A. Gurevich, Mikhail Dvorkin, Alexander S. Kulikov, Valery M. Lesin, et al. 2012. "SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing." *Journal of Computational Biology: A Journal of Computational Molecular Cell Biology* 19 (5): 455–77.
- Baranova, Maria A., Maria D. Logacheva, Aleksey A. Penin, Vladimir B. Seplyarskiy, Yana Y. Safonova, Sergey A. Naumenko, Anna V. Klepikova, et al. 2015. "Extraordinary Genetic Diversity in a Wood Decay Mushroom." *Molecular Biology and Evolution* 32 (10): 2775–83.
- Barrett, J. A., J. Antonovics, Bryan Campbell Clarke, Linda Partridge, Alan Robertson, Bryan Campbell Clarke, and Linda Partridge. 1988. "Frequency-Dependent Selection in Plant-Fungal Interactions." *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 319 (1196): 473–83.
- Barrick, Jeffrey E., and Richard E. Lenski. 2013. "Genome Dynamics during Experimental Evolution." *Nature Reviews. Genetics* 14 (12): 827–39.
- Barrick, Jeffrey E., Dong Su Yu, Sung Ho Yoon, Haeyoung Jeong, Tae Kwang Oh, Dominique Schneider, Richard E. Lenski, and Jihyun F. Kim. 2009. "Genome Evolution and Adaptation in a Long-Term Experiment with *Escherichia Coli*." *Nature* 461 (7268): 1243–47.

- Barton, John P., Nilu Goonetilleke, Thomas C. Butler, Bruce D. Walker, Andrew J. McMichael, and Arup K. Chakraborty. 2016. "Relative Rate and Location of Intra-Host HIV Evolution to Evade Cellular Immunity Are Predictable." *Nature Communications* 7 (May): 11660.
- Barton, N. H. 1995. "A General Model for the Evolution of Recombination." *Genetical Research* 65 (2): 123–45.
- . 1998. "The Effect of Hitch-Hiking on Neutral Genealogies." *Genetics Research* 72 (2): 123–33.
- . 2010. "Genetic Linkage and Natural Selection." *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 365 (1552): 2559–69.
- Barton, N. H., and B. Charlesworth. 1998. "Why Sex and Recombination?" *Science* 281 (5385): 1986–90.
- Bateson, W. 1909. "Heredity and Variation in Modern Lights." In *Darwin and Modern Science*, edited by A. C. Seward, 85–101. Cambridge University Press.
- Bazykin, Georgii A. 2015. "Changing Preferences: Deformation of Single Position Amino Acid Fitness Landscapes and Evolution of Proteins." *Biology Letters* 11 (10).
- Bazykin, Georgii A., Fyodor A. Kondrashov, Aleksey Y. Ogurtsov, Shamil Sunyaev, and Alexey S. Kondrashov. 2004. "Positive Selection at Sites of Multiple Amino Acid Replacements since Rat-Mouse Divergence." *Nature* 429 (6991): 558–62.
- Bedford, Trevor, Marc A. Suchard, Philippe Lemey, Gytis Dudas, Victoria Gregory, Alan J. Hay, John W. McCauley, Colin A. Russell, Derek J. Smith, and Andrew Rambaut. 2014. "Integrating Influenza Antigenic Dynamics with Molecular Evolution." *eLife* 3 (February): e01914.
- Beissinger, T. M., M. Gholami, M. Erbe, S. Weigend, A. Weigend, N. de Leon, D. Gianola, and H. Simianer. 2016. "Using the Variability of Linkage Disequilibrium between Subpopulations to Infer Sweeps and Epistatic Selection in a Diverse Panel of Chickens." *Heredity* 116 (2): 158–66.
- Belinky, Frida, Itamar Sela, Igor B. Rogozin, and Eugene V. Koonin. 2019. "Crossing Fitness Valleys via Double Substitutions within Codons." *BMC Biology* 17 (1): 105.
- Benger, Etam, and Guy Sella. 2013. "Modeling the Effect of Changing Selective Pressures on Polymorphism and Divergence." *Theoretical Population Biology* 85 (May): 73–85.
- Benjamini, Yoav, and Yosef Hochberg. 1995. "Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing." *Journal of the Royal Statistical Society* 57 (1): 289–300.
- Bergland, Alan O., Emily L. Behrman, Katherine R. O'Brien, Paul S. Schmidt, and Dmitri A. Petrov. 2014. "Genomic Evidence of Rapid and Stable Adaptive Oscillations over Seasonal Time Scales in *Drosophila*." *PLoS Genetics* 10 (11): e1004775.
- Berglund, Jonas, Katherine S. Pollard, and Matthew T. Webster. 2009. "Hotspots of Biased Nucleotide Substitutions in Human Genes." *PLoS Biology* 7 (1): e26.
- Bershtein, Shimon, Michal Segal, Roy Bekerman, Nobuhiko Tokuriki, and Dan S. Tawfik. 2006. "Robustness–epistasis Link Shapes the Fitness Landscape of a Randomly Drifting Protein." *Nature* 444 (7121): 929–32.
- Bershtein, Shimon, Adrian Wr Serohijos, and Eugene I. Shakhnovich. 2017. "Bridging the Physical Scales in Evolutionary Biology: From Protein Sequence Space to Fitness of Organisms and Populations." *Current Opinion in Structural Biology* 42 (February): 31–40.
- Bertram, Jason, and Joanna Masel. 2020. "Evolution Rapidly Optimizes Stability and Aggregation in Lattice Proteins Despite Pervasive Landscape Valleys and Mazes." *Genetics* 214 (4): 1047–57.
- Bezmenova, Aleksandra V., Elena A. Zvyagina, Anna V. Fedotova, Artem S. Kasianov, Tatiana V. Neretina, Aleksey A. Penin, Georgii A. Bazykin, and Alexey S. Kondrashov. 2020. "Rapid Accumulation of Mutations in Growing Mycelia of a Hypervariable Fungus *Schizophyllum Commune*." *Molecular Biology and Evolution*, April.
- Bhatt, Samir, Edward C. Holmes, and Oliver G. Pybus. 2011. "The Genomic Rate of Molecular Adaptation of the Human Influenza A Virus." *Molecular Biology and Evolution* 28 (9): 2443–51.
- Biswas, Avik, Allan Haldane, Eddy Arnold, and Ronald M. Levy. 2019. "Epistasis and

- Entrenchment of Drug Resistance in HIV-1 Subtype B." *eLife* 8 (October).
- Blanchette, Mathieu, W. James Kent, Cathy Riemer, Laura Elnitski, Arian F. A. Smit, Krishna M. Roskin, Robert Baertsch, et al. 2004. "Aligning Multiple Genomic Sequences with the Threaded Blockset Aligner." *Genome Research* 14 (4): 708–15.
- Bloom, Jesse D. 2014. "An Experimentally Determined Evolutionary Model Dramatically Improves Phylogenetic Fit." *Molecular Biology and Evolution* 31 (8): 1956–78.
- Bloom, Jesse D., Lizhi Ian Gong, and David Baltimore. 2010. "Permissive Secondary Mutations Enable the Evolution of Influenza Oseltamivir Resistance." *Science* 328 (5983): 1272–75.
- Bokma, F. 2002. "Detection of Punctuated Equilibrium from Molecular Phylogenies." *Journal of Evolutionary Biology* 15 (6): 1048–56.
- Borghans, José A. M., Joost B. Beltman, and Rob J. De Boer. 2004. "MHC Polymorphism under Host-Pathogen Coevolution." *Immunogenetics* 55 (11): 732–39.
- Boyrie, Léa, Corentin Moreau, Florian Frugier, Christophe Jacquet, and Maxime Bonhomme. 2021. "A Linkage Disequilibrium-Based Statistical Test for Genome-Wide Epistatic Selection Scans in Structured Populations." *Heredity* 126 (1): 77–91.
- Brand, Cara L., Lori Wright, and Daven C. Presgraves. 2019. "Positive Selection and Functional Divergence at Meiosis Genes That Mediate Crossing Over Across the Drosophila Phylogeny." *G3* 9 (10): 3201–11.
- Breen, Michael S., Carsten Kemena, Peter K. Vlasov, Cedric Notredame, and Fyodor A. Kondrashov. 2012. "Epistasis as the Primary Factor in Molecular Evolution." *Nature* 490 (7421): 535–38.
- Bridgham, Jamie T., Eric A. Ortlund, and Joseph W. Thornton. 2009. "An Epistatic Ratchet Constrains the Direction of Glucocorticoid Receptor Evolution." *Nature* 461 (7263): 515–19.
- Burger, Lukas, and Erik van Nimwegen. 2010. "Disentangling Direct from Indirect Co-Evolution of Residues in Protein Alignments." *PLoS Computational Biology* 6 (1): e1000633.
- Campos, José L., Daniel L. Halligan, Penelope R. Haddrill, and Brian Charlesworth. 2014. "The Relation between Recombination Rate and Patterns of Molecular Evolution and Variation in *Drosophila Melanogaster*." *Molecular Biology and Evolution*.  
<https://doi.org/10.1093/molbev/msu056>.
- Carius, H. J., T. J. Little, and D. Ebert. 2001. "Genetic Variation in a Host-Parasite Association: Potential for Coevolution and Frequency-Dependent Selection." *Evolution; International Journal of Organic Evolution* 55 (6): 1136–45.
- Chan, Yvonne H., Sergey V. Venev, Konstantin B. Zeldovich, and C. Robert Matthews. 2017. "Correlation of Fitness Landscapes from Three Orthologous TIM Barrels Originates from Sequence and Structure Constraints." *Nature Communications* 8 (March): 14614.
- Charlesworth, B. 1996. "Background Selection and Patterns of Genetic Diversity in *Drosophila Melanogaster*." *Genetical Research* 68 (2): 131–49.
- Charlesworth, B., M. T. Morgan, and D. Charlesworth. 1993. "The Effect of Deleterious Mutations on Neutral Molecular Variation." *Genetics* 134 (4): 1289–1303.
- Charlesworth, Brian. 1990. "Mutation-Selection Balance and the Evolutionary Advantage of Sex and Recombination." *Genetical Research* 55 (3): 199–221.
- . 2012. "The Effects of Deleterious Mutations on Evolution at Linked Sites." *Genetics* 190 (1): 5–22.
- Charlesworth, Brian, and Deborah Charlesworth. 1973. "Selection of New Inversions in Multi-Locus Genetic Systems." *Genetics Research* 21 (2): 167–83.
- Charlesworth, D., and B. Charlesworth. 1975. "Theoretical Genetics of Batesian Mimicry II. Evolution of Supergenes." *Journal of Theoretical Biology* 55 (2): 305–24.
- Charlesworth, Deborah. 2006. "Balancing Selection and Its Effects on Sequences in Nearby Genome Regions." *PLoS Genetics* 2 (4): e64.
- Cheverud, J. M., and E. J. Routman. 1995. "Epistasis and Its Contribution to Genetic Variance Components." *Genetics* 139 (3): 1455–61.
- Comeron, J. M., A. Williford, and R. M. Kliman. 2008. "The Hill–Robertson Effect: Evolutionary Consequences of Weak Selection and Linkage in Finite Populations." *Heredity* 100 (1): 19–31.

- Conover, D. O., and D. A. Van Voorhees. 1990. "Evolution of a Balanced Sex Ratio by Frequency-Dependent Selection in a Fish." *Science* 250 (4987): 1556–58.
- Corbett-Detig, Russell B., Daniel L. Hartl, and Timothy B. Sackton. 2015. "Natural Selection Constrains Neutral Diversity across a Wide Range of Species." *PLoS Biology* 13 (4): e1002112.
- Coyne, Jerry A., Nicholas H. Barton, and Michael Turelli. 1997. "PERSPECTIVE: A CRITIQUE OF SEWALL WRIGHT'S SHIFTING BALANCE THEORY OF EVOLUTION." *Evolution; International Journal of Organic Evolution* 51 (3): 643–71.
- Crow, James F. 1987. "Muller, Dobzhansky, and Overdominance." *Journal of the History of Biology* 20 (3): 351–80.
- . 2010. "On Epistasis: Why It Is Unimportant in Polygenic Directional Selection." *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 365 (1544): 1241–44.
- Crow, James Franklin, and Motoo Kimura. 1970. *An Introduction to Population Genetics Theory*. New York, Evanston and London: Harper & Row, Publishers.
- Csilléry, Katalin, Michael G. B. Blum, Oscar E. Gaggiotti, and Olivier François. 2010. "Approximate Bayesian Computation (ABC) in Practice." *Trends in Ecology & Evolution* 25 (7): 410–18.
- Csilléry, Katalin, Olivier François, and Michael G. B. Blum. 2012. "Abc: An R Package for Approximate Bayesian Computation (ABC)." *Methods in Ecology and Evolution* 3 (3): 475–79.
- Cutler, D. J. 2000. "Understanding the Overdispersed Molecular Clock." *Genetics* 154 (3): 1403–17.
- Cutter, Asher D., Richard Jovelin, and Alivia Dey. 2013. "Molecular Hyperdiversity and Evolution in Very Large Populations." *Molecular Ecology* 22 (8): 2074–95.
- Cvijovic, Ivana, Benjamin H. Good, Elizabeth R. Jerison, and Michael M. Desai. 2015. "The Fate of a Mutation in a Fluctuating Environment." *Proc. Natl. Acad. Sci. U. S. A.* 112 (36): E5021–28.
- DeGiorgio, Michael, Kirk E. Lohmueller, and Rasmus Nielsen. 2014. "A Model-Based Approach for Identifying Signatures of Ancient Balancing Selection in Genetic Data." *PLoS Genetics* 10 (8): e1004561.
- De Leonardis, Eleonora, Benjamin Lutz, Sebastian Ratz, Simona Cocco, Rémi Monasson, Alexander Schug, and Martin Weigt. 2015. "Direct-Coupling Analysis of Nucleotide Coevolution Facilitates RNA Secondary and Tertiary Structure Prediction." *Nucleic Acids Research* 43 (21): 10444–55.
- Delph, Lynda F., and John K. Kelly. 2014. "On the Importance of Balancing Selection in Plants." *The New Phytologist* 201 (1): 45–56.
- Desai, Michael M., and Daniel S. Fisher. 2007. "Beneficial Mutation–Selection Balance and the Effect of Linkage on Positive Selection." *Genetics* 176 (3): 1759–98.
- Desai, Michael M., Daniel Weissman, and Marcus W. Feldman. 2007. "Evolution Can Favor Antagonistic Epistasis." *Genetics* 177 (2): 1001–10.
- Dimmic, Matthew W., Melissa J. Hubisz, Carlos D. Bustamante, and Rasmus Nielsen. 2005. "Detecting Coevolving Amino Acid Sites Using Bayesian Mutational Mapping." *Bioinformatics* 21 Suppl 1 (June): i126–35.
- Diss, Guillaume, and Ben Lehner. 2018. "The Genetic Landscape of a Physical Interaction." *eLife* 7 (April).
- Dobzhansky, T. 1936. "Studies on Hybrid Sterility. II. Localization of Sterility Factors in *Drosophila Pseudoobscura* Hybrids." *Genetics* 21 (2): 113–35.
- . 1950. "Genetics of Natural Populations. XIX. Origin of Heterosis through Natural Selection in Populations of *Drosophila Pseudoobscura*." *Genetics* 35 (3): 288–302.
- Dobzhansky, Theodosius. 1937. *Genetics and the Origin of Species*. Columbia University Press.
- Dobzhansky, Theodosius, and Olga Pavlovsky. 1957. "An Experimental Study of Interaction between Genetic Drift and Natural Selection." *Evolution; International Journal of Organic Evolution* 11 (3): 311–19.
- Dobzhansky, T., and A. H. Sturtevant. 1938. "Inversions in the Chromosomes of *Drosophila Pseudoobscura*." *Genetics* 23 (1): 28–64.

- Dodson, M. M., and A. Hallam. 1977. "Allopatric Speciation and the Fold Catastrophe." *The American Naturalist* 111 (979): 415–33.
- Domingo, Júlia, Guillaume Diss, and Ben Lehner. 2018. "Pairwise and Higher-Order Genetic Interactions during the Evolution of a tRNA." *Nature* 558 (7708): 117–21.
- Doud, Michael B., Orr Ashenberg, and Jesse D. Bloom. 2015. "Site-Specific Amino Acid Preferences Are Mostly Conserved in Two Closely Related Protein Homologs." *Molecular Biology and Evolution* 32 (11): 2944–60.
- Draghi, Jeremy A., Todd L. Parsons, Günter P. Wagner, and Joshua B. Plotkin. 2010. "Mutational Robustness Can Facilitate Adaptation." *Nature* 463 (7279): 353–55.
- Ellegren, Hans, and Nicolas Galtier. 2016. "Determinants of Genetic Diversity." *Nature Reviews. Genetics* 17 (7): 422–33.
- Eyre-Walker, Adam, and Peter D. Keightley. 2009. "Estimating the Rate of Adaptive Molecular Evolution in the Presence of Slightly Deleterious Mutations and Population Size Change." *Molecular Biology and Evolution* 26 (9): 2097–2108.
- Falconer, R. S., and Douglas Scott Falconer. 1989. *Introduction to Quantitative Genetics*. Longman, Scientific & Technical.
- Faria, Rui, Kerstin Johannesson, Roger K. Butlin, and Anja M. Westram. 2019. "Evolving Inversions." *Trends in Ecology & Evolution* 34 (3): 239–48.
- Fenster, C. B., L. F. Galloway, and L. Chao. 1997. "Epistasis and Its Consequences for the Evolution of Natural Populations." *Trends in Ecology & Evolution* 12 (7): 282–86.
- Ferretti, Luca, B. Schmiegel, and D. Weinreich. 2016. "Measuring Epistasis in Fitness Landscapes: The Correlation of Fitness Effects of Mutations." *Journal of Theoretical Biology* 396: 132–43.
- Ferretti, Luca, Daniel Weinreich, Fumio Tajima, and Guillaume Achaz. 2018. "Evolutionary Constraints in Fitness Landscapes." *Heredity* 121 (5): 466–81.
- Ferrière, Régis, Ulf Dieckmann, and Denis Couvet. 2004. *Evolutionary Conservation Biology*. Cambridge University Press.
- Figliuzzi, Matteo, Hervé Jacquier, Alexander Schug, Oliver Tenailon, and Martin Weigt. 2016. "Coevolutionary Landscape Inference and the Context-Dependence of Mutations in Beta-Lactamase TEM-1." *Molecular Biology and Evolution* 33 (1): 268–80.
- Fisher, Ronald Aylmer. 1941. "THE THEORETICAL CONSEQUENCES OF POLYPLOID INHERITANCE FOR THE MID STYLE FORM OF LYTHRUM SALICARIA." *Annals of Eugenics*.  
 ———. 1930. *The Genetical Theory of Natural Selection*. The Clarendon Press.  
 ———. 1918. "The Correlation Between Relatives on the Supposition of Mendelian Inheritance." *Transactions of the Royal Society of Edinburgh* 52: 399–433.  
 ———. 1922. "On the Dominance Ratio." *Proceedings of the Royal Society of Edinburgh* 42: 321–41.
- Flynn, William F., Allan Haldane, Bruce E. Torbett, and Ronald M. Levy. 2017. "Inference of Epistatic Effects Leading to Entrenchment and Drug Resistance in HIV-1 Protease." *Molecular Biology and Evolution* 34 (6): 1291–1306.
- Fowler, Douglas M., and Stanley Fields. 2014. "Deep Mutational Scanning: A New Style of Protein Science." *Nature Methods* 11 (8): 801–7.
- Fragata, Inês, Alexandre Blanckaert, Marco António Dias Louro, David A. Liberles, and Claudia Bank. 2019. "Evolution in the Light of Fitness Landscape Theory." *Trends in Ecology & Evolution* 34 (1): 69–82.
- Fragata, Inês, Sebastian Matuszewski, Mark A. Schmitz, Thomas Bataillon, Jeffrey D. Jensen, and Claudia Bank. 2018. "The Fitness Landscape of the Codon Space across Environments." *Heredity* 121 (5): 422–37.
- Franke, Jasper, Alexander Klözer, J. Arjan G. M. de Visser, and Joachim Krug. 2011. "Evolutionary Accessibility of Mutational Pathways." *PLoS Computational Biology* 7 (8): e1002134.
- Franklin, I., and R. C. Lewontin. 1970. "Is the Gene the Unit of Selection?" *Genetics* 65 (4): 707–34.
- Friedlander, Tamar, Roshan Prizak, Nicholas H. Barton, and Gašper Tkačik. 2017. "Evolution of New Regulatory Functions on Biophysically Realistic Fitness Landscapes." *Nature*

- Communications* 8 (1): 216.
- Galtier, Nicolas, Laurent Duret, Sylvain Glémin, and Vincent Ranwez. 2009. "GC-Biased Gene Conversion Promotes the Fixation of Deleterious Amino Acid Changes in Primates." *Trends in Genetics: TIG* 25 (1): 1–5.
- Garcia, Jesse A., and Kirk E. Lohmueller. 2021. "Negative Linkage Disequilibrium between Amino Acid Changing Variants Reveals Interference among Deleterious Mutations in the Human Genome." *PLoS Genetics* 17 (7): 1–25.
- Gavrilets, S. 1997. "Evolution and Speciation on Holey Adaptive Landscapes." *Trends in Ecology & Evolution* 12 (8): 307–12.
- Gavrilets, Sergey. 1999. "Evolution and Speciation in a Hyperspace: The Roles of Neutrality, Selection, Mutation, and Random Drift." In *Towards a Comprehensive Dynamics of Evolution - Exploring the Interplay of Selection, Neutrality, Accident, and Function*, edited by Crutchfield J And Schuster, 135–62. Oxford University Press Oxford.
- . 2004. *Fitness Landscapes and the Origin of Species*. Princeton Univ. Press, Princeton, NJ.
- Gerrish, P. J., and R. E. Lenski. 1998. "The Fate of Competing Beneficial Mutations in an Asexual Population." *Genetica* 102-103 (1-6): 127–44.
- Gigord, L. D., M. R. Macnair, and A. Smithson. 2001. "Negative Frequency-Dependent Selection Maintains a Dramatic Flower Color Polymorphism in the Rewardless Orchid *Dactylorhiza Sambucina* (L.) Soo." *Proc. Natl. Acad. Sci. U. S. A.* 98 (11): 6253–55.
- Gilbert, Kimberly J., Fanny Pouyet, Laurent Excoffier, and Stephan Peischl. 2020. "Transition from Background Selection to Associative Overdominance Promotes Diversity in Regions of Low Recombination." *Current Biology: CB* 30 (1): 101–7.e3.
- Gillespie, J. H. 1984. "The Molecular Clock May Be an Episodic Clock." *Proceedings of the National Academy of Sciences of the United States of America* 81 (24): 8009–13.
- . 1973. "Polymorphism in Random Environments." *Theoretical Population Biology* 4 (2): 193–95.
- . 1991. *The Causes of Molecular Evolution*. Oxford University Press.
- . 1994. *The Causes of Molecular Evolution*. Oxford University Press.
- Gillespie, J. H., and C. H. Langley. 1979. "Are Evolutionary Rates Really Variable?" *Journal of Molecular Evolution* 13 (1): 27–34.
- Göbel, U., C. Sander, R. Schneider, and A. Valencia. 1994. "Correlated Mutations and Residue Contacts in Proteins." *Proteins* 18 (4): 309–17.
- Gokhale, Chaitanya S., Yoh Iwasa, Martin A. Nowak, and Arne Traulsen. 2009. "The Pace of Evolution across Fitness Valleys." *Journal of Theoretical Biology* 259 (3): 613–20.
- Goldstein, Richard A., Stephen T. Pollard, Seena D. Shah, and David D. Pollock. 2015. "Nonadaptive Amino Acid Convergence Rates Decrease over Time." *Molecular Biology and Evolution* 32 (6): 1373–81.
- Goldstein, Richard A., and David D. Pollock. 2017. "Sequence Entropy of Folding and the Absolute Rate of Amino Acid Substitutions." *Nature Ecology & Evolution* 1 (12): 1923–30.
- Gong, Lizhi Ian, Marc A. Suchard, and Jesse D. Bloom. 2013. "Stability-Mediated Epistasis Constrains the Evolution of an Influenza Protein." *eLife* 2 (May): e00631.
- Good, Benjamin H. 2020. "Linkage Disequilibrium between Rare Mutations." *bioRxiv*. <https://doi.org/10.1101/2020.12.10.420042>.
- Good, Benjamin H., and Michael M. Desai. 2015. "The Impact of Macroscopic Epistasis on Long-Term Evolutionary Dynamics." *Genetics* 199 (1): 177–90.
- Good, Benjamin H., Michael J. McDonald, Jeffrey E. Barrick, Richard E. Lenski, and Michael M. Desai. 2017. "The Dynamics of Molecular Evolution over 60,000 Generations." *Nature* 551 (7678): 45–50.
- Gould, S. J., and N. Eldredge. 1993. "Punctuated Equilibrium Comes of Age." *Nature* 366 (6452): 223–27.
- Graves, J. L., Jr, K. L. Hertweck, M. A. Phillips, M. V. Han, L. G. Cabral, T. T. Barter, L. F. Greer, M. K. Burke, L. D. Mueller, and M. R. Rose. 2017. "Genomics of Parallel Experimental Evolution in *Drosophila*." *Molecular Biology and Evolution* 34 (4): 831–42.
- Gros, Pierre-Alexis, Hervé Le Nagard, and Olivier Tenaillon. 2009. "The Evolution of Epistasis and

- Its Links with Genetic Robustness, Complexity and Drift in a Phenotypic Model of Adaptation." *Genetics* 182 (1): 277–93.
- Gutiérrez-Valencia, Juanita, P. William Hughes, Emma L. Berdan, and Tanja Slotte. 2021. "The Genomic Architecture and Evolutionary Fates of Supergenes." *Genome Biology and Evolution* 13 (5).
- Haddox, Hugh K., Adam S. Dingens, Sarah K. Hilton, Julie Overbaugh, and Jesse D. Bloom. 2018. "Mapping Mutational Effects along the Evolutionary Landscape of HIV Envelope." *eLife* 7 (March).
- Haller, Benjamin C., and Philipp W. Messer. 2019. "SLiM 3: Forward Genetic Simulations Beyond the Wright–Fisher Model." *Molecular Biology and Evolution* 36 (3): 632–37.
- Harmand, Noémie, Romain Gallet, Roula Jabbour-Zahab, Guillaume Martin, and Thomas Lenormand. 2017. "Fisher's Geometrical Model and the Mutational Patterns of Antibiotic Resistance across Dose Gradients." *Evolution; International Journal of Organic Evolution* 71 (1): 23–37.
- Hayman, B. I., and K. Mather. 1955. "The Description of Genic Interactions in Continuous Variation." *Biometrics* 11 (1): 69–82.
- Hedrick, Philip W. 2012. "What Is the Evidence for Heterozygote Advantage Selection?" *Trends in Ecology & Evolution* 27 (12): 698–704.
- Hegreness, Matthew, Noam Shores, Doris Damian, Daniel Hartl, and Roy Kishony. 2008. "Accelerated Evolution of Resistance in Multidrug Environments." *Proc. Natl. Acad. Sci. U. S. A.* 105 (37): 13977–81.
- Hellmann, Ines, Ingo Ebersberger, Susan E. Ptak, Svante Pääbo, and Molly Przeworski. 2003. "A Neutral Explanation for the Correlation of Diversity with Recombination Rates in Humans." *American Journal of Human Genetics* 72 (6): 1527–35.
- Hemani, Gibran, Joseph E. Powell, Huanwei Wang, Konstantin Shakhbazov, Harm-Jan Westra, Tonu Esko, Anjali K. Henders, et al. 2021. "Phantom Epistasis between Unlinked Loci." *Nature* 596 (7871): E1–3.
- Hermesen, Rutger, J. Barrett Deris, and Terence Hwa. 2012. "On the Rapidity of Antibiotic Resistance Evolution Facilitated by a Concentration Gradient." *Proc. Natl. Acad. Sci. U. S. A.* 109 (27): 10775–80.
- Hietpas, Ryan T., Claudia Bank, Jeffrey D. Jensen, and Daniel N. A. Bolon. 2013. "Shifting Fitness Landscapes in Response to Altered Environments." *Evolution; International Journal of Organic Evolution* 67 (12): 3512–22.
- Hill, W. G., and A. Robertson. 1966. "The Effect of Linkage on Limits to Artificial Selection." *Genetical Research* 8 (3): 269–94.
- Hill, William G., Michael E. Goddard, and Peter M. Visscher. 2008. "Data and Theory Point to Mainly Additive Genetic Variance for Complex Traits." *PLoS Genetics* 4 (2): e1000008.
- Hinkley, Trevor, João Martins, Colombe Chappay, Mojgan Haddad, Eric Stawiski, Jeannette M. Whitcomb, Christos J. Petropoulos, and Sebastian Bonhoeffer. 2011. "A Systems Analysis of Mutational Effects in HIV-1 Protease and Reverse Transcriptase." *Nature Genetics* 43 (5): 487–89.
- Hivert, V., J. Sidorenko, F. Rohart, M. E. Goddard, and J. Yang. 2021. "Estimation of Non-Additive Genetic Variance in Human Complex Traits from a Large Sample of Unrelated Individuals." *American Journal of Human Genetics* 108,5: 786–98.
- Hough, Josh, Wei Wang, Spencer C. H. Barrett, and Stephen I. Wright. 2017. "Hill-Robertson Interference Reduces Genetic Diversity on a Young Plant Y-Chromosome." *Genetics* 207 (2): 685–95.
- Ho, Wei-Chin, and Jianzhi Zhang. 2018. "Evolutionary Adaptations to New Environments Generally Reverse Plastic Phenotypic Changes." *Nature Communications* 9 (1): 350.
- Huerta-Sanchez, Emilia, Rick Durrett, and Carlos D. Bustamante. 2008. "Population Genetics of Polymorphism and Divergence under Fluctuating Selection." *Genetics* 178 (1): 325–37.
- Hunt, Gene. 2007. "The Relative Importance of Directional Change, Random Walks, and Stasis in the Evolution of Fossil Lineages." *Proceedings of the National Academy of Sciences of the United States of America* 104 (47): 18404–8.

- . 2008. “Gradual or Pulsed Evolution: When Should Punctuational Explanations Be Preferred?” *Paleobiology* 34 (03): 360–77.
- Hunt, Gene, Melanie J. Hopkins, and Scott Lidgard. 2015. “Simple versus Complex Models of Trait Evolution and Stasis as a Response to Environmental Change.” *Proceedings of the National Academy of Sciences of the United States of America* 112 (16): 4885–90.
- Ingram, Travis. 2011. “Speciation along a Depth Gradient in a Marine Adaptive Radiation.” *Proceedings. Biological Sciences / The Royal Society* 278 (1705): 613–18.
- Jacquier, Hervé, André Birgy, Hervé Le Nagard, Yves Mechulam, Emmanuelle Schmitt, Jérémy Glodt, Beatrice Bercot, et al. 2013. “Capturing the Mutational Landscape of the Beta-Lactamase TEM-1.” *Proc. Natl. Acad. Sci. U. S. A.* 110 (32): 13067–72.
- Jain, Kavita, and Joachim Krug. 2007. “Deterministic and Stochastic Regimes of Asexual Evolution on Rugged Fitness Landscapes.” *Genetics* 175 (3): 1275–88.
- Johnson, Milo S., Shreyas Gopalakrishnan, Juhee Goyal, Megan E. Dillingham, Christopher W. Bakerlee, Parris T. Humphrey, Tanush Jagdish, et al. 2021. “Phenotypic and Molecular Evolution across 10,000 Generations in Laboratory Budding Yeast Populations.” *eLife* 10 (January).
- Johnson, Milo S., Alena Martsul, Sergey Kryazhimskiy, and Michael M. Desai. 2019. “Higher-Fitness Yeast Genotypes Are Less Robust to Deleterious Mutations.” *Science* 366 (6464): 490–93.
- Jordan, Daniel M., Stephan G. Frangakis, Christelle Golzio, Christopher A. Cassa, Joanne Kurtzberg, Task Force for Neonatal Genomics, Erica E. Davis, Shamil R. Sunyaev, and Nicholas Katsanis. 2015. “Identification of Cis-Suppression of Human Disease Mutations by Comparative Genomics.” *Nature* 524 (7564): 225–29.
- Joron, Mathieu, Lise Frezal, Robert T. Jones, Nicola L. Chamberlain, Siu F. Lee, Christoph R. Haag, Annabel Whibley, et al. 2011. “Chromosomal Rearrangements Maintain a Polymorphic Supergene Controlling Butterfly Mimicry.” *Nature* 477 (7363): 203–6.
- Kamisetty, H., and S. Ovchinnikov. 2013. “Assessing the Utility of Coevolution-Based Residue–residue Contact Predictions in a Sequence-and Structure-Rich Era.” *Proc. Natl. Acad. Sci. U. S. A.* 110 (39): 15674–79.
- Katsnelson, Mikhail I., Yuri I. Wolf, and Eugene V. Koonin. 2019. “On the Feasibility of Saltational Evolution.” *Proc. Natl. Acad. Sci. U. S. A.* 116 (42): 21068–75.
- Kauffman, S. A., and E. D. Weinberger. 1989. “The NK Model of Rugged Fitness Landscapes and Its Application to Maturation of the Immune Response.” *Journal of Theoretical Biology* 141 (2): 211–45.
- Kauffman, S., and S. Levin. 1987. “Towards a General Theory of Adaptive Walks on Rugged Landscapes.” *Journal of Theoretical Biology* 128 (1): 11–45.
- Kern, Andrew D., and Fyodor A. Kondrashov. 2004. “Mechanisms and Convergence of Compensatory Evolution in Mammalian Mitochondrial tRNAs.” *Nature Genetics* 36 (11): 1207–12.
- Kimura, M. 1954. “Process Leading to Quasi-Fixation of Genes in Natural Populations Due to Random Fluctuation of Selection Intensities.” *Genetics* 39 (3): 280–95.
- . 1965. “Attainment of Quasi Linkage Equilibrium When Gene Frequencies Are Changing by Natural Selection.” *Genetics* 52 (5): 875–90.
- . 1977. “Preponderance of Synonymous Changes as Evidence for the Neutral Theory of Molecular Evolution.” *Nature* 267 (5608): 275–76.
- . 1983. *The Neutral Theory of Molecular Evolution*. Cambridge University Press.
- Kimura, M., and T. Maruyama. 1966. “The Mutational Load with Epistatic Gene Interactions in Fitness.” *Genetics* 54 (6): 1337–51.
- Kimura, M., and T. Ohta. 1974. “Probability of Gene Fixation in an Expanding Finite Population.” *Proceedings of the National Academy of Sciences of the United States of America* 71 (9): 3377–79.
- Kim, Yuseob, and H. Allen Orr. 2005. “Adaptation in Sexuals vs. Asexuals: Clonal Interference and the Fisher-Muller Model.” *Genetics* 171 (3): 1377–86.
- Kingman, J. F. C. 1978. “A Simple Model for the Balance between Selection and Mutation.” *Journal*

- of Applied Probability* 15 (1): 1–12.
- Klink, Galya V., and Georgii A. Bazykin. 2017. "Parallel Evolution of Metazoan Mitochondrial Proteins." *Genome Biology and Evolution* 9 (5): 1341–50.
- Klink, Galya V., Andrey V. Golovin, and Georgii A. Bazykin. 2017. "Substitutions into Amino Acids That Are Pathogenic in Human Mitochondrial Proteins Are More Frequent in Lineages Closely Related to Human than in Distant Lineages." *PeerJ* 5 (December): e4143.
- Koch, Evan, Mickey Ristroph, and Mark Kirkpatrick. 2013. "Long Range Linkage Disequilibrium across the Human Genome." *PLoS One* 8 (12): e80754.
- Kogenaru, Manjunatha, Marjon G. J. de Vos, and Sander J. Tans. 2009. "Revealing Evolutionary Pathways by Fitness Landscape Reconstruction." *Critical Reviews in Biochemistry and Molecular Biology* 44 (4): 169–74.
- Kondrashov, Alexey S. 2018. "Through Sex, Nature Is Telling Us Something Important." *Trends in Genetics: TIG* 34 (5): 352–61.
- Kondrashov, Alexey S., Inna S. Povolotskaya, Dmitry N. Ivankov, and Fyodor A. Kondrashov. 2010. "Rate of Sequence Divergence under Constant Selection." *Biology Direct* 5 (January): 5.
- Kondrashov, Alexey S., Shamil Sunyaev, and Fyodor A. Kondrashov. 2002. "Dobzhansky-Muller Incompatibilities in Protein Evolution." *Proceedings of the National Academy of Sciences of the United States of America* 99 (23): 14878–83.
- Kondrashov, Dmitry A., and Fyodor A. Kondrashov. 2015. "Topological Features of Rugged Fitness Landscapes in Sequence Space." *Trends in Genetics: TIG* 31 (1): 24–33.
- Kondrashov, F. A., and A. S. Kondrashov. 2001. "Multidimensional Epistasis and the Disadvantage of Sex." *Proceedings of the National Academy of Sciences* 98 (21): 12089–92.
- Kothe, E. 1999. "Mating Types and Pheromone Recognition in the Homobasidiomycete *Schizophyllum Commune*." *Fungal Genetics and Biology: FG & B* 27 (2-3): 146–52.
- Kouyos, Roger D., Sarah P. Otto, and Sebastian Bonhoeffer. 2006. "Effect of Varying Epistasis on the Evolution of Recombination." *Genetics* 173 (2): 589–97.
- Kouyos, Roger D., Olin K. Silander, and Sebastian Bonhoeffer. 2007. "Epistasis between Deleterious Mutations and the Evolution of Recombination." *Trends in Ecology & Evolution* 22 (6): 308–15.
- Krishna, K. Hari, Yallamandayya Vadlamudi, and Muthuvel Suresh Kumar. 2016. "Viral Evolved Inhibition Mechanism of the RNA Dependent Protein Kinase PKR's Kinase Domain, a Structural Perspective." *PLoS One* 11 (4): e0153680.
- Kryazhimskiy, Sergey, Jonathan Dushoff, Georgii A. Bazykin, and Joshua B. Plotkin. 2011. "Prevalence of Epistasis in the Evolution of Influenza A Surface Proteins." *PLoS Genetics* 7 (2): e1001301.
- Kryazhimskiy, Sergey, Daniel P. Rice, Elizabeth R. Jerison, and Michael M. Desai. 2014. "Microbial Evolution. Global Epistasis Makes Adaptation Predictable despite Sequence-Level Stochasticity." *Science* 344 (6191): 1519–22.
- Kryazhimskiy, Sergey, Gasper Tkacik, and Joshua B. Plotkin. 2009. "The Dynamics of Adaptation on Correlated Fitness Landscapes." *Proc. Natl. Acad. Sci. U. S. A.* 106 (44): 18638–43.
- Kulathinal, Rob J., Brian R. Bettencourt, and Daniel L. Hartl. 2004. "Compensated Deleterious Mutations in Insect Genomes." *Science* 306 (5701): 1553–54.
- Kumar, Sudhir, and S. Blair Hedges. 1998. "A Molecular Timescale for Vertebrate Evolution." *Nature* 392 (6679): 917–20.
- Kumar, Sujai, Martin Jones, Georgios Koutsovoulos, Michael Clarke, and Mark Blaxter. 2013. "Blobology: Exploring Raw Genome Data for Contaminants, Symbionts and Parasites Using Taxon-Annotated GC-Coverage Plots." *Frontiers in Genetics* 4 (November): 237.
- Kunte, K., W. Zhang, A. Tenger-Trolander, D. H. Palmer, A. Martin, R. D. Reed, S. P. Mullen, and M. R. Kronforst. 2014. "Doublesex Is a Mimicry Supergene." *Nature* 507 (7491): 229–32.
- Kvitek, Daniel J., and Gavin Sherlock. 2011. "Reciprocal Sign Epistasis between Frequently Experimentally Evolved Adaptive Mutations Causes a Rugged Fitness Landscape." *PLoS Genetics* 7 (4): e1002056.
- Lack, Justin B., Charis M. Cardeno, Marc W. Crepeau, William Taylor, Russell B. Corbett-Detig, Kristian A. Stevens, Charles H. Langley, and John E. Pool. 2015. "The *Drosophila* Genome

- Nexus: A Population Genomic Resource of 623 *Drosophila Melanogaster* Genomes, Including 197 from a Single Ancestral Range Population." *Genetics* 199 (4): 1229–41.
- Lande, Russell. 1976. "NATURAL SELECTION AND RANDOM GENETIC DRIFT IN PHENOTYPIC EVOLUTION." *Evolution; International Journal of Organic Evolution* 30 (2): 314–34.
- Lang, Gregory I., David Botstein, and Michael M. Desai. 2011. "Genetic Variation and the Fate of Beneficial Mutations in Asexual Populations." *Genetics* 188 (3): 647–61.
- Lang, Gregory I., and Michael M. Desai. 2014. "The Spectrum of Adaptive Mutations in Experimental Evolution." *Genomics* 104 (6 Pt A): 412–16.
- Lässig, Michael, Ville Mustonen, and Aleksandra M. Walczak. 2017. "Predicting Evolution." *Nature Ecology & Evolution* 1 (3): 77.
- Lawrie, David S., Dmitri A. Petrov, and Philipp W. Messer. 2011. "Faster than Neutral Evolution of Constrained Sequences: The Complex Interplay of Mutational Biases and Weak Selection." *Genome Biology and Evolution* 3 (April): 383–95.
- Lee, Juhye M., John Huddleston, Michael B. Doud, Kathryn A. Hooper, Nicholas C. Wu, Trevor Bedford, and Jesse D. Bloom. 2018. "Deep Mutational Scanning of Hemagglutinin Helps Predict Evolutionary Fates of Human H3N2 Influenza Variants." *Proc. Natl. Acad. Sci. U. S. A.* 115 (35): E8276–85.
- Leffler, Ellen M., Kevin Bullaughey, Daniel R. Matute, Wynn K. Meyer, Laure Ségurel, Aarti Venkat, Peter Andolfatto, and Molly Przeworski. 2012. "Revisiting an Old Riddle: What Determines Genetic Diversity Levels within Species?" *PLoS Biology* 10 (9): e1001388.
- Leffler, Ellen M., Ziyue Gao, Susanne Pfeifer, Laure Ségurel, Adam Auton, Oliver Venn, Rory Bowden, et al. 2013. "Multiple Instances of Ancient Balancing Selection Shared between Humans and Chimpanzees." *Science* 339 (6127): 1578–82.
- Lenski, R. E. 1998. "Bacterial Evolution and the Cost of Antibiotic Resistance." *International Microbiology: The Official Journal of the Spanish Society for Microbiology* 1 (4): 265–70.
- Lenski, Richard E., Michael R. Rose, Suzanne C. Simpson, and Scott C. Tadler. 1991. "Long-Term Experimental Evolution in *Escherichia Coli*. I. Adaptation and Divergence During 2,000 Generations." *The American Naturalist* 138 (6): 1315–41.
- Lewontin, R. C., and Ken-Ichi Kojima. 1960. "The Evolutionary Dynamics of Complex Polymorphisms." *Evolution; International Journal of Organic Evolution* 14 (4): 458–72.
- Li, Chuan, and Jianzhi Zhang. 2018. "Multi-Environment Fitness Landscapes of a tRNA Gene." *Nature Ecology & Evolution* 2 (6): 1025–32.
- Li, W. H., C. I. Wu, and C. C. Luo. 1985. "A New Method for Estimating Synonymous and Nonsynonymous Rates of Nucleotide Substitution Considering the Relative Likelihood of Nucleotide and Codon Changes." *Molecular Biology and Evolution* 2 (2): 150–74.
- Lunzer, Mark, G. Brian Golding, and Antony M. Dean. 2010. "Pervasive Cryptic Epistasis in Molecular Evolution." *PLoS Genetics* 6 (10): e1001162.
- Lynch, M. 1987. "The Consequences of Fluctuating Selection for Isozyme Polymorphisms in *Daphnia*." *Genetics* 115 (4): 657–69.
- Lynch, Michael. 2010. "Evolution of the Mutation Rate." *Trends in Genetics: TIG* 26 (8): 345–52.
- Lyons, Daniel M., Zhengting Zou, Haiqing Xu, and Jianzhi Zhang. 2020. "Idiosyncratic Epistasis Creates Universals in Mutational Effects and Evolutionary Trajectories." *Nature Ecology & Evolution* 4 (12): 1685–93.
- Mackay, Trudy Fc, and Jason H. Moore. 2014. "Why Epistasis Is Important for Tackling Complex Human Disease Genetics." *Genome Medicine* 6 (6): 124.
- Mahler, D. Luke, Liam J. Revell, Richard E. Glor, and Jonathan B. Losos. 2010. "Ecological Opportunity and the Rate of Morphological Evolution in the Diversification of Greater Antillean Anoles." *Evolution; International Journal of Organic Evolution* 64 (9): 2731–45.
- Mäki-Tanila, Asko, and William G. Hill. 2014. "Influence of Gene Interaction on Complex Trait Variation with Multilocus Models." *Genetics* 198 (1): 355–67.
- Malmberg, R. L. 1977. "The Evolution of Epistasis and the Advantage of Recombination in Populations of Bacteriophage T4." *Genetics* 86 (3): 607–21.
- Malmgren, Björn A., W. A. Berggren, and G. P. Lohmann. 1983. "Evidence for Punctuated Gradualism in the Late Neogene *Globorotalia Tumida* Lineage of Planktonic Foraminifera."

- Paleobiology* 9 (04): 377–89.
- Marks, Debora S., Lucy J. Colwell, Robert Sheridan, Thomas A. Hopf, Andrea Pagnani, Riccardo Zecchina, and Chris Sander. 2011. “Protein 3D Structure Computed from Evolutionary Sequence Variation.” *PloS One* 6 (12): e28766.
- Martin, A. P., and S. R. Palumbi. 1993. “Body Size, Metabolic Rate, Generation Time, and the Molecular Clock.” *Proc. Natl. Acad. Sci. U. S. A.* 90 (9): 4087–91.
- Martin, Guillaume, Santiago F. Elena, and Thomas Lenormand. 2007. “Distributions of Epistasis in Microbes Fit Predictions from a Fitness Landscape Model.” *Nature Genetics* 39 (4): 555–60.
- Masel, Joanna. 2006. “Cryptic Genetic Variation Is Enriched for Potential Adaptations.” *Genetics* 172 (3): 1985–91.
- Masel, Joanna, and Meredith V. Trotter. 2010. “Robustness and Evolvability.” *Trends in Genetics: TIG* 26 (9): 406–14.
- Mather, Kenneth. 1950. “The Genetical Architecture of Heterostyly in *Primula Sinensis*.” *Evolution; International Journal of Organic Evolution* 4 (4): 340–52.
- Mattila, Tiina M., and Folmer Bokma. 2008. “Extant Mammal Body Masses Suggest Punctuated Equilibrium.” *Proceedings. Biological Sciences / The Royal Society* 275 (1648): 2195–99.
- Maynard, John, and John Haigh. 1974. “The Hitch-Hiking Effect of a Favourable Gene.” *Genetics Research* 23: 23–35.
- Maynard Smith, J. 1970. “Natural Selection and the Concept of a Protein Space.” *Nature* 225 (5232): 563–64.
- McCandlish, David M. 2011. “Visualizing Fitness Landscapes.” *Evolution; International Journal of Organic Evolution* 65 (6): 1544–58.
- McCandlish, David M., Etienne Rajon, Premal Shah, Yang Ding, and Joshua B. Plotkin. 2013. “The Role of Epistasis in Protein Evolution.” *Nature*.
- McCandlish, David M., Premal Shah, and Joshua B. Plotkin. 2016. “Epistasis and the Dynamics of Reversion in Molecular Evolution.” *Genetics* 203 (3): 1335–51.
- Meer, Margarita V., Alexey S. Kondrashov, Yael Artzy-Randrup, and Fyodor A. Kondrashov. 2010. “Compensatory Evolution in Mitochondrial tRNAs Navigates Valleys of Low Fitness.” *Nature* 464 (7286): 279–82.
- Miller, Craig R. 2019. “The Treacheries of Adaptation.” *Science*.
- Mooers, Mooers, Vamosi, and Schluter. 1999. “Using Phylogenies to Test Macroevolutionary Hypotheses of Trait Evolution in Cranes (Gruinae).” *The American Naturalist* 154 (2): 249.
- Morcos, Faruck, Andrea Pagnani, Bryan Lunt, Arianna Bertolino, Debora S. Marks, Chris Sander, Riccardo Zecchina, José N. Onuchic, Terence Hwa, and Martin Weigt. 2011. “Direct-Coupling Analysis of Residue Coevolution Captures Native Contacts across Many Protein Families.” *Proc. Natl. Acad. Sci. U. S. A.* 108 (49): E1293–1301.
- Muller, H. 1942. “Isolating Mechanisms, Evolution, and Temperature.” *Biol. Symp.* 6: 71–125.
- Muller, H. J. 1932. “Some Genetic Aspects of Sex.” *The American Naturalist* 66 (703): 118–38.
- . 1964. “THE RELATION OF RECOMBINATION TO MUTATIONAL ADVANCE.” *Mutation Research* 106 (May): 2–9.
- Mustonen, Ville, and Michael Lässig. 2008. “Molecular Evolution under Fitness Fluctuations.” *Physical Review Letters* 100 (10).
- . 2009. “From Fitness Landscapes to Seascape: Non-Equilibrium Dynamics of Selection and Adaptation.” *Trends in Genetics: TIG* 25 (3): 111–19.
- Mustonen, V., and M. Lässig. 2007. “Adaptations to Fluctuating Selection in *Drosophila*.” *Proc. Natl. Acad. Sci. U. S. A.* 104 (7): 2277–82.
- . 2010. “Fitness Flux and Ubiquity of Adaptive Evolution.” *Proc. Natl. Acad. Sci. U. S. A.* 107 (9): 4248–53.
- Nabieva, Elena, and Georgii A. Bazykin. 2019. “SELVa: Simulator of Evolution with Landscape Variation.” *bioRxiv*. <https://doi.org/10.1101/647834>.
- Nasrallah, Chris A., and John P. Huelsenbeck. 2013. “A Phylogenetic Model for the Detection of Epistatic Interactions.” *Molecular Biology and Evolution* 30 (9): 2197–2208.
- Natarajan, Chandrasekhar, Federico G. Hoffmann, Roy E. Weber, Angela Fago, Christopher C. Witt,

- and Jay F. Storz. 2016. "Predictable Convergence in Hemoglobin Function Has Unpredictable Molecular Underpinnings." *Science* 354 (6310): 336–39.
- Naumenko, Sergey A., Alexey S. Kondrashov, and Georgii A. Bazykin. 2012. "Fitness Conferred by Replaced Amino Acids Declines with Time." *Biology Letters* 8 (5): 825–28.
- Naumenko, Sergey A., Maria D. Logacheva, Nina V. Popova, Anna V. Klepikova, Aleksey A. Penin, Georgii A. Bazykin, Anna E. Etingova, Nikolai S. Mugue, Alexey S. Kondrashov, and Lev Y. Yampolsky. 2017. "Transcriptome-Based Phylogeny of Endemic Lake Baikal Amphipod Species Flock: Fast Speciation Accompanied by Frequent Episodes of Positive Selection." *Molecular Ecology* 26 (2): 536–53.
- Neher, E. 1994. "How Frequent Are Correlated Changes in Families of Protein Sequences?" *Proc. Natl. Acad. Sci. U. S. A.* 91 (1): 98–102.
- Neher, Richard A., and Thomas Leitner. 2010. "Recombination Rate and Selection Strength in HIV Intra-Patient Evolution." *PLoS Computational Biology* 6 (1): e1000660.
- Neher, Richard A., and Boris I. Shraiman. 2009. "Competition between Recombination and Epistasis Can Cause a Transition from Allele to Genotype Selection." *Proc. Natl. Acad. Sci. U. S. A.* 106 (16): 6866–71.
- Neidhart, Johannes, Ivan G. Szendro, and Joachim Krug. 2014. "Adaptation in Tunably Rugged Fitness Landscapes: The Rough Mount Fuji Model." *Genetics* 198 (2): 699–721.
- Nei, M., and T. Gojobori. 1986. "Simple Methods for Estimating the Numbers of Synonymous and Nonsynonymous Nucleotide Substitutions." *Molecular Biology and Evolution* 3 (5): 418–26.
- Nei, M., and W. H. Li. 1973. "Linkage Disequilibrium in Subdivided Populations." *Genetics* 75 (1): 213–19.
- Nelson, Erik D., and Nick V. Grishin. 2019. "How Often Do Protein Genes Navigate Valleys of Low Fitness?" *Genes* 10 (4).
- Neverov, Alexey D., Sergey Kryazhimskiy, Joshua B. Plotkin, and Georgii A. Bazykin. 2014. "Coordinated Evolution of Influenza A Surface Proteins." *PLoS Genet* 11 (8): e1005404.
- Neverov, Alexey D., Anfisa V. Popova, Gennady G. Fedonin, Evgeny A. Cheremukhin, Galya V. Klink, and Georgii A. Bazykin. 2021. "Episodic Evolution of Coadapted Sets of Amino Acid Sites in Mitochondrial Proteins." *PLoS Genetics* 17 (1): e1008711.
- Nordborg, M., B. Charlesworth, and D. Charlesworth. 1996. "The Effect of Recombination on Background Selection." *Genetical Research* 67 (2): 159–74.
- Ochs, Ian E., and Michael M. Desai. 2015. "The Competition between Simple and Complex Evolutionary Trajectories in Asexual Populations." *BMC Evolutionary Biology* 15 (March): 55.
- Ohm, Robin A., Jan F. de Jong, Luis G. Lugones, Andrea Aerts, Erika Kothe, Jason E. Stajich, Ronald P. de Vries, et al. 2010. "Genome Sequence of the Model Mushroom *Schizophyllum commune*." *Nature Biotechnology* 28 (9): 957–63.
- Ohta, Tomoko. 1971. "Associative Overdominance Caused by Linked Detrimental Mutations." *Genetics Research* 18 (3): 277–86.
- Ohta, Tomoko, and Motoo Kimura. 1971. "On the Constancy of the Evolutionary Rate of Cistrons." *Journal of Molecular Evolution* 1 (1): 18–25.
- Olson, C. Anders, Nicholas C. Wu, and Ren Sun. 2014. "A Comprehensive Biophysical Description of Pairwise Epistasis throughout an Entire Protein Domain." *Current Biology: CB* 24 (22): 2643–51.
- O'Meara, Brian C., Cécile Ané, Michael J. Sanderson, and Peter C. Wainwright. 2006. "Testing for Different Rates of Continuous Trait Evolution Using Likelihood." *Evolution; International Journal of Organic Evolution* 60 (5): 922–33.
- Orr, H. Allen. 1995. "The Population Genetics of Speciation: The Evolution of Hybrid Incompatibilities." *Genetics* 139 (4): 1805–13.
- . 1998. "The Population Genetics of Adaptation: The Distribution of Factors Fixed during Adaptive Evolution." *Evolution* 52 (4): 935–49.
- . 2005. "The Genetic Theory of Adaptation: A Brief History." *Nature Reviews. Genetics* 6 (2): 119–27.
- Otto, Sarah P., Maria R. Servedio, and Scott L. Nuismer. 2008. "Frequency-Dependent Selection

- and the Evolution of Assortative Mating." *Genetics* 179 (4): 2091–2112.
- Otto, S. P., and M. W. Feldman. 1997. "Deleterious Mutations, Variable Epistatic Interactions, and the Evolution of Recombination." *Theoretical Population Biology* 51 (2): 134–47.
- Otto, S. P., and A. C. Gerstein. 2006. "Why Have Sex? The Population Genetics of Sex and Recombination." *Biochemical Society Transactions* 34 (Pt 4): 519–22.
- Otwinowski, Jakub, David M. McCandlish, and Joshua B. Plotkin. 2018. "Inferring the Shape of Global Epistasis." *Proc. Natl. Acad. Sci. U. S. A.* 115 (32): E7550–58.
- Ovchinnikov, Sergey, Hetunandan Kamisetty, and David Baker. 2014. "Robust and Accurate Prediction of Residue–residue Interactions across Protein Interfaces Using Evolutionary Information." *eLife* 3 (May): e02030.
- Ovchinnikov, Sergey, Lisa Kinch, Hahnbeom Park, Yuxing Liao, Jimin Pei, David E. Kim, Hetunandan Kamisetty, Nick V. Grishin, and David Baker. 2015. "Large-Scale Determination of Previously Unsolved Protein Structures Using Evolutionary Information." *eLife* 4 (September): e09248.
- Ovchinnikov, Sergey, Hahnbeom Park, Neha Varghese, Po-Ssu Huang, Georgios A. Pavlopoulos, David E. Kim, Hetunandan Kamisetty, Nikos C. Kyrpides, and David Baker. 2017. "Protein Structure Determination Using Metagenome Sequence Data." *Science* 355 (6322): 294–98.
- Pagel, Mark, Chris Venditti, and Andrew Meade. 2006. "Large Punctuational Contribution of Speciation to Evolutionary Divergence at the Molecular Level." *Science* 314 (5796): 119–21.
- Pál, Csaba, and Balázs Papp. 2017. "Evolution of Complex Adaptations in Molecular Systems." *Nature Ecology & Evolution* 1 (8): 1084–92.
- Paquin, Charlotte, and Julian Adams. 1983. "Frequency of Fixation of Adaptive Mutations Is Higher in Evolving Diploid than Haploid Yeast Populations." *Nature*.
- Park, Leeyoung. 2019. "Population-Specific Long-Range Linkage Disequilibrium in the Human Genome and Its Influence on Identifying Common Disease Variants." *Scientific Reports* 9 (1): 11380.
- Park, Su-Chan, and Joachim Krug. 2007. "Clonal Interference in Large Populations." *Proc. Natl. Acad. Sci. U. S. A.* 104 (46): 18135–40.
- Passagem-Santos, Diogo, Simone Zacarias, and Lilia Perfeito. 2018. "Power Law Fitness Landscapes and Their Ability to Predict Fitness." *Heredity* 121 (5): 482–98.
- Pedruzzi, Gabriele, Ayuna Barlukova, and Igor M. Rouzine. 2018. "Evolutionary Footprint of Epistasis." *PLoS Computational Biology* 14 (9): e1006426.
- Pedruzzi, Gabriele, and Igor M. Rouzine. 2019. "Epistasis Detectably Alters Correlations between Genomic Sites in a Narrow Parameter Window." *PloS One* 14 (5): e0214036.
- Pennell, Matthew W., Luke J. Harmon, and Josef C. Uyeda. 2014. "Is There Room for Punctuated Equilibrium in Macroevolution?" *Trends in Ecology & Evolution* 29 (1): 23–32.
- Pennings, Pleuni S., Sergey Kryazhimskiy, and John Wakeley. 2014. "Loss and Recovery of Genetic Diversity in Adapting Populations of HIV." *PLoS Genetics* 10 (1): e1004000.
- Perelman, Polina, Warren E. Johnson, Christian Roos, Hector N. Seuánez, Julie E. Horvath, Miguel A. M. Moreira, Bailey Kessing, et al. 2011. "A Molecular Phylogeny of Living Primates." *PLoS Genetics* 7 (3): e1001342.
- Phillips, Angela M., Katherine R. Lawrence, Alief Moulana, Thomas Dupic, Jeffrey Chang, Milo S. Johnson, Ivana Cvijovic, Thierry Mora, Aleksandra M. Walczak, and Michael M. Desai. 2021. "Binding Affinity Landscapes Constrain the Evolution of Broadly Neutralizing Anti-Influenza Antibodies." *eLife* 10 (September).
- Phillips, Patrick C. 2008. "Epistasis--the Essential Role of Gene Interactions in the Structure and Evolution of Genetic Systems." *Nature Reviews. Genetics* 9 (11): 855–67.
- Pillai, Arvind S., Shane A. Chandler, Yang Liu, Anthony V. Signore, Carlos R. Cortez-Romero, Justin L. P. Benesch, Arthur Laganowsky, Jay F. Storz, Georg K. A. Hochberg, and Joseph W. Thornton. 2020. "Origin of Complexity in Haemoglobin Evolution." *Nature* 581 (7809): 480–85.
- Podgornaia, Anna Igorevna. 2014. *Pervasive Degeneracy and Epistasis in a Protein-Protein Interface*. Massachusetts Institute of Technology, Computational and Systems Biology Program.

- Poelwijk, Frank J., Daniel J. Kiviet, Daniel M. Weinreich, and Sander J. Tans. 2007. "Empirical Fitness Landscapes Reveal Accessible Evolutionary Paths." *Nature* 445 (7126): 383–86.
- Poelwijk, Frank J., Michael Socolich, and Rama Ranganathan. n.d. "Learning the Pattern of Epistasis Linking Genotype and Phenotype in a Protein." <https://doi.org/10.1101/213835>.
- Poelwijk, Frank J., Sorin Tănase-Nicola, Daniel J. Kiviet, and Sander J. Tans. 2011. "Reciprocal Sign Epistasis Is a Necessary Condition for Multi-Peaked Fitness Landscapes." *Journal of Theoretical Biology* 272 (1): 141–44.
- Pokusaeva, Victoria O., Dinara R. Usmanova, Ekaterina V. Putintseva, Lorena Espinar, Karen S. Sarkisyan, Alexander S. Mishin, Natalya S. Bogatyreva, et al. 2019. "An Experimental Assay of the Interactions of Amino Acids from Orthologous Sequences Shaping a Complex Fitness Landscape." *PLoS Genetics* 15 (4): e1008079.
- Pollard, Katherine S., Sofie R. Salama, Bryan King, Andrew D. Kern, Tim Dreszer, Sol Katzman, Adam Siepel, et al. 2006. "Forces Shaping the Fastest Evolving Regions in the Human Genome." *PLoS Genetics* 2 (10): e168.
- Pollock, David D., Grant Thiltgen, and Richard A. Goldstein. 2012. "Amino Acid Coevolution Induces an Evolutionary Stokes Shift." *Proc. Natl. Acad. Sci. U. S. A.* 109 (21): E1352–59.
- Pollock, D. D., W. R. Taylor, and N. Goldman. 1999. "Coevolving Protein Residues: Maximum Likelihood Identification and Relationship to Structure." *Journal of Molecular Biology* 287 (1): 187–98.
- Popova, Anfisa V., Ksenia R. Safina, Vasily V. Ptushenko, Anastasia V. Stolyarova, Alexander V. Favorov, Alexey D. Neverov, and Georgii A. Bazykin. 2019. "Allele-Specific Nonstationarity in Evolution of Influenza A Virus Surface Proteins." *Proc. Natl. Acad. Sci. U. S. A.*, 201904246.
- Povolotskaya, Inna S., and Fyodor A. Kondrashov. 2010. "Sequence Space and the Ongoing Expansion of the Protein Universe." *Nature* 465 (7300): 922–26.
- Poynton, Helen C., Simone Hasenbein, Joshua B. Benoit, Maria S. Sepulveda, Monica F. Poelchau, Daniel S. T. Hughes, Shwetha C. Murali, et al. 2018. "The Toxicogenome of *Hyalella Azteca*: A Model for Sediment Ecotoxicology and Evolutionary Toxicology." *Environmental Science & Technology* 52 (10): 6009–22.
- Pritchard, J. K., M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman. 1999. "Population Growth of Human Y Chromosomes: A Study of Y Chromosome Microsatellites." *Molecular Biology and Evolution* 16 (12): 1791–98.
- Proschan, Frank. 1963. "Theoretical Explanation of Observed Decreasing Failure Rate." *Technometrics: A Journal of Statistics for the Physical, Chemical, and Engineering Sciences* 5 (3): 375–83.
- Provine, William B. 1989. *Sewall Wright and Evolutionary Biology*. University of Chicago Press.
- Puranen, Santeri, Maiju Pesonen, Johan Pensar, Ying Ying Xu, John A. Lees, Stephen D. Bentley, Nicholas J. Croucher, and Jukka Corander. 2018. "SuperDCA for Genome-Wide Epistasis Analysis." *Microbial Genomics* 4 (6).
- Ragsdale, Aaron P. 2021. "Can We Distinguish Modes of Selective Interactions Using Linkage Disequilibrium?" *bioRxiv*. <https://doi.org/10.1101/2021.03.25.437004>.
- Ranwez, Vincent, Sébastien Harispe, Frédéric Delsuc, and Emmanuel J. P. Douzery. 2011. "MACSE: Multiple Alignment of Coding SEquences Accounting for Frameshifts and Stop Codons." *PloS One* 6 (9): e22594.
- Rasmussen, Matthew D., Melissa J. Hubisz, Ilan Gronau, and Adam Siepel. 2014. "Genome-Wide Inference of Ancestral Recombination Graphs." *PLoS Genetics* 10 (5): e1004342.
- Reddy, Gautam, and Michael M. Desai. 2021. "Global Epistasis Emerges from a Generic Model of a Complex Trait." *eLife* 10 (March).
- Reis, Mario dos, Jun Inoue, Masami Hasegawa, Robert J. Asher, Philip C. J. Donoghue, and Ziheng Yang. 2012. "Phylogenomic Datasets Provide Both Precision and Accuracy in Estimating the Timescale of Placental Mammal Phylogeny." *Proceedings. Biological Sciences / The Royal Society* 279 (1742): 3491–3500.
- Risso, Valeria A., Fadia Manssour-Triedo, Asunción Delgado-Delgado, Rocio Arco, Alicia Barroso-delJesus, Alvaro Ingles-Prieto, Raquel Godoy-Ruiz, et al. 2015. "Mutational Studies on Resurrected Ancestral Proteins Reveal Conservation of Site-Specific Amino Acid

- Preferences throughout Evolutionary History." *Molecular Biology and Evolution* 32 (2): 440–55.
- Rogozin, Igor B., Karen Thomson, Miklós Csürös, Liran Carmel, and Eugene V. Koonin. 2008. "Homoplasy in Genome-Wide Analysis of Rare Amino Acid Replacements: The Molecular-Evolutionary Basis for Vavilov's Law of Homologous Series." *Biology Direct* 3 (March): 7.
- Rohlf, Rori V., Willie J. Swanson, and Bruce S. Weir. 2010. "Detecting Coevolution through Allelic Association between Physically Unlinked Loci." *American Journal of Human Genetics* 86 (5): 674–85.
- Rollins, Nathan J., Kelly P. Brock, Frank J. Poelwijk, Michael A. Stiffler, Nicholas P. Gauthier, Chris Sander, and Debora S. Marks. 2019. "Inferring Protein 3D Structure from Deep Mutation Scans." *Nature Genetics* 51 (7): 1170–76.
- Romanova, Elena V., Vladimir V. Aleoshin, Ravil M. Kamaltynov, Kirill V. Mikhailov, Maria D. Logacheva, Elena A. Sirotnina, Alexander Yu Gornov, Anton S. Anikin, and Dmitry Yu Sherbakov. 2016. "Evolution of Mitochondrial Genomes in Baikalian Amphipods." *BMC Genomics* 17 (Suppl 14): 1016.
- Rosenbloom, Kate R., Joel Armstrong, Galt P. Barber, Jonathan Casper, Hiram Clawson, Mark Diekhans, Timothy R. Dreszer, et al. 2015. "The UCSC Genome Browser Database: 2015 Update." *Nucleic Acids Research* 43 (Database issue): D670–81.
- Rothenburg, Stefan, Eun Joo Seo, James S. Gibbs, Thomas E. Dever, and Katharina Dittmar. 2009. "Rapid Evolution of Protein Kinase PKR Alters Sensitivity to Viral Inhibitors." *Nature Structural & Molecular Biology* 16 (1): 63–70.
- Roze, Denis, and Nick H. Barton. 2006. "The Hill–Robertson Effect and the Evolution of Recombination." *Genetics* 173 (3): 1793–1811.
- Sackton, Timothy B., and Daniel L. Hartl. 2016. "Genotypic Context and Epistasis in Individuals and Populations." *Cell* 166 (2): 279–87.
- Sailer, Zachary R., and Michael J. Harms. 2017a. "Detecting High-Order Epistasis in Nonlinear Genotype-Phenotype Maps." *Genetics* 205 (3): 1079–88.
- . 2017b. "High-Order Epistasis Shapes Evolutionary Trajectories." *PLoS Computational Biology* 13 (5): e1005541.
- Salinas, Victor H., and Rama Ranganathan. 2018. "Coevolution-Based Inference of Amino Acid Interactions Underlying Protein Function." *eLife* 7 (July).
- Salverda, Merijn L. M., Eynat Dellus, Florian A. Gorter, Alfons J. M. Debets, John van der Oost, Rolf F. Hoekstra, Dan S. Tawfik, and J. Arjan G. M. de Visser. 2011. "Initial Mutations Direct Alternative Pathways of Protein Evolution." *PLoS Genetics* 7 (3): e1001321.
- Sandler, George, Stephen I. Wright, and Aneil F. Agrawal. 2021. "Patterns and Causes of Signed Linkage Disequilibria in Flies and Plants." *Molecular Biology and Evolution*, June. <https://www.ncbi.nlm.nih.gov/pubmed/34097067>.
- Sarkisyan, Karen S., Dmitry A. Bolotin, Margarita V. Meer, Dinara R. Usmanova, Alexander S. Mishin, George V. Sharonov, Dmitry N. Ivankov, et al. 2016. "Local Fitness Landscape of the Green Fluorescent Protein." *Nature* 533 (7603): 397–401.
- Schaper, E., A. Eriksson, M. Rafajlovic, S. Sagitov, and B. Mehlig. 2012. "Linkage Disequilibrium under Recurrent Bottlenecks." *Genetics* 190 (1): 217–29.
- Schlosser, Gerhard, and Günter P. Wagner. 2008. "A Simple Model of Co-Evolutionary Dynamics Caused by Epistatic Selection." *Journal of Theoretical Biology* 250 (1): 48–65.
- Schoustra, Sijmen, Sungmin Hwang, Joachim Krug, and J. Arjan G. M. de Visser. 2016. "Diminishing-Returns Epistasis among Random Beneficial Mutations in a Multicellular Fungus." *Proc. Biol. Sci.* 283 (1837).
- Schrag, S. J., V. Perrot, and B. R. Levin. 1997. "Adaptation to the Fitness Costs of Antibiotic Resistance in *Escherichia Coli*." *Proc. Biol. Sci.* 264 (1386): 1287–91.
- Sella, Guy, Dmitri A. Petrov, Molly Przeworski, and Peter Andolfatto. 2009. "Pervasive Natural Selection in the *Drosophila* Genome?" *PLoS Genetics* 5 (6): e1000495.
- Seplyarskiy, Vladimir B., Maria D. Logacheva, Aleksey A. Penin, Maria A. Baranova, Evgeny V. Leushkin, Natalia V. Demidenko, Anna V. Klepikova, Fyodor A. Kondrashov, Alexey S.

- Kondrashov, and Timothy Y. James. 2014. "Crossing-over in a Hypervariable Species Preferentially Occurs in Regions of High Local Similarity." *Molecular Biology and Evolution* 31 (11): 3016–25.
- Shah, Premal, David M. McCandlish, and Joshua B. Plotkin. 2015. "Contingency and Entrenchment in Protein Evolution under Purifying Selection." *Proceedings of the National Academy of Sciences of the United States of America* 112 (25): E3226–35.
- Shapiro, Beth, Andrew Rambaut, Oliver G. Pybus, and Edward C. Holmes. 2006. "A Phylogenetic Method for Detecting Positive Epistasis in Gene Sequences and Its Application to RNA Virus Evolution." *Molecular Biology and Evolution* 23 (9): 1724–30.
- Shnol, E. E., and A. S. Kondrashov. 1993. "The Effect of Selection on the Phenotypic Variance." *Genetics* 134 (3): 995–96.
- Singh, Bashisth N. 2008. "Chromosome Inversions and Linkage Disequilibrium in *Drosophila*." *Current Science* 94 (4): 459–64.
- Sjodt, Megan, Kelly Brock, Genevieve Dobihal, Patricia D. A. Rohs, Anna G. Green, Thomas A. Hopf, Alexander J. Meeske, et al. 2018. "Structure of the Peptidoglycan Polymerase RodA Resolved by Evolutionary Coupling Analysis." *Nature* 556 (7699): 118–21.
- Sohail, Mashaal, Olga A. Vakhrusheva, Jae Hoon Sul, Sara L. Pult, Laurent C. Francioli, Genome of the Netherlands Consortium, Alzheimer's Disease Neuroimaging Initiative, et al. 2017. "Negative Selection in Humans and Fruit Flies Involves Synergistic Epistasis." *Science* 356 (6337): 539–42.
- Soylemez, Onuralp, and Fyodor A. Kondrashov. 2012. "Estimating the Rate of Irreversibility in Protein Evolution." *Genome Biology and Evolution* 4 (12): 1213–22.
- Spencer, Chris C. A., Panos Deloukas, Sarah Hunt, Jim Mullikin, Simon Myers, Bernard Silverman, Peter Donnelly, David Bentley, and Gil McVean. 2006. "The Influence of Recombination on Human Genetic Diversity." *PLoS Genetics* 2 (9): e148.
- Stamatakis, Alexandros. 2014. "RAxML Version 8: A Tool for Phylogenetic Analysis and Post-Analysis of Large Phylogenies." *Bioinformatics* 30 (9): 1312–13.
- Stanley, Steven M. 1998. *Macroevolution : Pattern and Process*. Johns Hopkins University Press.
- Starr, Tyler N., Julia M. Flynn, Parul Mishra, Daniel N. A. Bolon, and Joseph W. Thornton. 2018. "Pervasive Contingency and Entrenchment in a Billion Years of Hsp90 Evolution." *Proc. Natl. Acad. Sci. U. S. A.* 115 (17): 4453–58.
- Starr, Tyler N., and Joseph W. Thornton. 2016. "Epistasis in Protein Evolution." *Protein Science: A Publication of the Protein Society* 25 (7): 1204–18.
- Steinberg, Barrett, and Marc Ostermeier. 2016. "Environmental Changes Bridge Evolutionary Valleys." *Science Advances* 2 (1): e1500921.
- Stiffler, Michael A., Frank J. Poelwijk, Kelly P. Brock, Richard R. Stein, Adam Riesselman, Joan Teyra, Sachdev S. Sidhu, Debora S. Marks, Nicholas P. Gauthier, and Chris Sander. 2020. "Protein Structure from Experimental Evolution." *Cell Systems* 10 (1): 15–24.e5.
- Storz, Jay F. 2016. "Causes of Molecular Convergence and Parallelism in Protein Evolution." *Nature Reviews. Genetics* 17 (4): 239–50.
- Strotz, Luke C., and Andrew P. Allen. 2013. "Assessing the Role of Cladogenesis in Macroevolution by Integrating Fossil and Molecular Evidence." *Proc. Natl. Acad. Sci. U. S. A.* 110 (8): 2904–9.
- Sturtevant, A. H., and K. Mather. 1938. "The Interrelations of Inversions, Heterosis and Recombination." *The American Naturalist* 72 (742): 447–52.
- Szendro, Ivan G., Jasper Franke, J. Arjan G. M. de Visser, and Joachim Krug. 2013. "Predictability of Evolution Depends Nonmonotonically on Population Size." *Proc. Natl. Acad. Sci. U. S. A.* 110 (2): 571–76.
- Szendro, Ivan G., Martijn F. Schenk, Jasper Franke, Joachim Krug, and J. Arjan G. M. de Visser. 2013. "Quantitative Analyses of Empirical Fitness Landscapes." *Journal of Statistical Mechanics* 2013 (01): P01005.
- Takahasi, K. Ryo, and Fumio Tajima. 2005. "Evolution of Coadaptation in a Two-Locus Epistatic System." *Evolution; International Journal of Organic Evolution* 59 (11): 2324–32.
- Takahata, N., K. Ishii, and H. Matsuda. 1975. "Effect of Temporal Fluctuation of Selection Coefficient on Gene Frequency in a Population." *Proc. Natl. Acad. Sci. U. S. A.* 72 (11):

4541–45.

- Takahata, N., and M. Nei. 1990. "Allelic Genealogy under Overdominant and Frequency-Dependent Selection and Polymorphism of Major Histocompatibility Complex Loci." *Genetics* 124 (4): 967–78.
- Tamura, Koichiro, Fabia Ursula Battistuzzi, Paul Billing-Ross, Oscar Murillo, Alan Filipski, and Sudhir Kumar. 2012. "Estimating Divergence Times in Large Molecular Phylogenies." *Proceedings of the National Academy of Sciences of the United States of America* 109 (47): 19333–38.
- Tellier, Aurélien, and James K. M. Brown. 2011. "Spatial Heterogeneity, Frequency-Dependent Selection and Polymorphism in Host-Parasite Interactions." *BMC Evolutionary Biology* 11 (November): 319.
- Tenaillon, O. 2014. "The Utility of Fisher's Geometric Model in Evolutionary Genetics." *Annual Review of Ecology, Evolution, and Systematics* 45 (November): 179–201.
- Tenaillon, Olivier, Alejandra Rodríguez-Verdugo, Rebecca L. Gaut, Pamela McDonald, Albert F. Bennett, Anthony D. Long, and Brandon S. Gaut. 2012. "The Molecular Diversity of Adaptive Convergence." *Science* 335 (6067): 457–61.
- Thompson, M. J., and C. D. Jiggins. 2014. "Supergenes and Their Role in Evolution." *Heredity* 113 (1): 1–8.
- Trubenová, Barbora, Martin S. Krejca, Per Kristian Lehre, and Timo Kötzing. 2019. "Surfing on the Seascape: Adaptation in a Changing Environment." *Evolution; International Journal of Organic Evolution* 73 (7): 1356–74.
- Usmanova, Dinara R., Luca Ferretti, Inna S. Povolotskaya, Peter K. Vlasov, and Fyodor A. Kondrashov. 2015. "A Model of Substitution Trajectories in Sequence Space and Long-Term Protein Evolution." *Molecular Biology and Evolution* 32 (2): 542–54.
- Van Cleve, Jeremy, and Daniel B. Weissman. 2015. "Measuring Ruggedness in Fitness Landscapes." *Proc. Natl. Acad. Sci. U. S. A.*
- Vaupel, J. W., K. G. Manton, and E. Stallard. 1979. "The Impact of Heterogeneity in Individual Frailty on the Dynamics of Mortality." *Demography* 16 (3): 439–54.
- Venditti, Chris, and Mark Pagel. 2008. "Speciation and Bursts of Evolution." *Evolution: Education and Outreach* 1 (3): 274–80.
- Visscher, Peter M., Stuart Macgregor, Beben Benyamin, Gu Zhu, Scott Gordon, Sarah Medland, William G. Hill, et al. 2007. "Genome Partitioning of Genetic Variation for Height from 11,214 Sibling Pairs." *American Journal of Human Genetics* 81 (5): 1104–10.
- Visser, J. Arjan G. M. de, Tim F. Cooper, and Santiago F. Elena. 2011. "The Causes of Epistasis." *Proc. Royal Soc. B* 278 (1725): 3617–24.
- Visser, J. Arjan G. M. de, Santiago F. Elena, Inês Fragata, and Sebastian Matuszewski. 2018. "The Utility of Fitness Landscapes and Big Data for Predicting Evolution." *Heredity* 121 (5): 401–5.
- Visser, J. Arjan G. M. de, Joachim Hermisson, Günter P. Wagner, Lauren Ancel Meyers, Homayoun Bagheri-Chaichian, Jeffrey L. Blanchard, Lin Chao, et al. 2003. "Perspective: Evolution and Detection of Genetic Robustness." *Evolution; International Journal of Organic Evolution* 57 (9): 1959–72.
- Visser, J. Arjan G. M. de, Rolf F. Hoekstra, and Herman van den Ende. 1997. "Test of Interaction between Genetic Markers That Affect Fitness in *Aspergillus Niger*." *Evolution* 51 (5): 1499–1505.
- Visser, J. Arjan G. M. de, and Joachim Krug. 2014. "Empirical Fitness Landscapes and the Predictability of Evolution." *Nature Reviews. Genetics* 15 (7): 480–90.
- Voje, Kjetil Lysne. 2016. "Tempo Does Not Correlate with Mode in the Fossil Record." *Evolution; International Journal of Organic Evolution* 70 (12): 2678–89.
- Vos, Marjon G. J. de, Frank J. Poelwijk, Nico Battich, Joseph D. T. Ndika, and Sander J. Tans. 2013. "Environmental Dependence of Genetic Constraint." *PLoS Genetics* 9 (6): e1003580.
- Wagner, Andreas. 2008. "Robustness and Evolvability: A Paradox Resolved." *Proc. Biol. Sci.* 275 (1630): 91–100.
- Wang, Ming-Chih, Feng-Chi Chen, Yen-Zho Chen, Yao-Ting Huang, and Trees-Juen Chuang. 2012.

- “LDGIdb: A Database of Gene Interactions Inferred from Long-Range Strong Linkage Disequilibrium between Pairs of SNPs.” *BMC Research Notes* 5 (May): 212.
- Wang, Shenshen, and Lei Dai. 2019. “Evolving Generalists in Switching Rugged Landscapes.” *PLoS Computational Biology* 15 (10): e1007320.
- Weigt, Martin, Robert A. White, Hendrik Szurmant, James A. Hoch, and Terence Hwa. 2009. “Identification of Direct Residue Contacts in Protein–protein Interaction by Message Passing.” *Proc. Natl. Acad. Sci. U. S. A.* 106 (1): 67–72.
- Weinreich, Daniel M., and Lin Chao. 2005. “Rapid Evolutionary Escape by Large Populations from Local Fitness Peaks Is Likely in Nature.” *Evolution* 59 (6): 1175–82.
- Weinreich, Daniel M., Nigel F. Delaney, Mark A. Depristo, and Daniel L. Hartl. 2006. “Darwinian Evolution Can Follow Only Very Few Mutational Paths to Fitter Proteins.” *Science* 312 (5770): 111–14.
- Weinreich, Daniel M., Yinghong Lan, C. Scott Wylie, and Robert B. Heckendorn. 2013. “Should Evolutionary Geneticists Worry about Higher-Order Epistasis?” *Current Opinion in Genetics & Development* 23 (6): 700–707.
- Weinreich, Daniel M., Richard A. Watson, and Lin Chao. 2005. “Perspective: Sign Epistasis and Genetic Constraint on Evolutionary Trajectories.” *Evolution; International Journal of Organic Evolution* 59 (6): 1165–74.
- Weissman, Daniel B., Michael M. Desai, Daniel S. Fisher, and Marcus W. Feldman. 2009. “The Rate at Which Asexual Populations Cross Fitness Valleys.” *Theoretical Population Biology* 75 (4): 286–300.
- Weissman, Daniel B., Marcus W. Feldman, and Daniel S. Fisher. 2010. “The Rate of Fitness-Valley Crossing in Sexual Populations.” *Genetics* 186 (4): 1389–1410.
- Wei, Xinzhu, and Jianzhi Zhang. 2019. “Patterns and Mechanisms of Diminishing Returns from Beneficial Mutations.” *Molecular Biology and Evolution* 36 (5): 1008–21.
- Wiser, Michael J., Noah Ribeck, and Richard E. Lenski. 2013. “Long-Term Dynamics of Adaptation in Asexual Populations.” *Science* 342 (6164): 1364–67.
- Wolf, Yuri I., Cecile Viboud, Edward C. Holmes, Eugene V. Koonin, and David J. Lipman. 2006. “Long Intervals of Stasis Punctuated by Bursts of Positive Selection in the Seasonal Evolution of Influenza A Virus.” *Biology Direct* 1 (October): 34.
- Wood, Andrew R., Marcus A. Tuke, Mike A. Nalls, Dena G. Hernandez, Stefania Bandinelli, Andrew B. Singleton, David Melzer, Luigi Ferrucci, Timothy M. Frayling, and Michael N. Weedon. 2014. “Another Explanation for Apparent Epistasis.” *Nature*.
- Woodcock, G., and P. G. Higgs. 1996. “Population Evolution on a Multiplicative Single-Peak Fitness Landscape.” *Journal of Theoretical Biology* 179 (1): 61–73.
- Woods, Robert, Dominique Schneider, Cynthia L. Winkworth, Margaret A. Riley, and Richard E. Lenski. 2006. “Tests of Parallel Molecular Evolution in a Long-Term Experiment with *Escherichia Coli*.” *Proc. Natl. Acad. Sci. U. S. A.* 103 (24): 9107–12.
- Wright, S. 1931. “Evolution in Mendelian Populations.” *Genetics* 16 (2): 97–159.
- . 1932. “The Roles of Mutation, Inbreeding, Crossbreeding and Selection in Evolution.” *Proc 6th Int Cong Genet.* 1: 356–66.
- Wu, Nicholas C., Lei Dai, C. Anders Olson, James O. Lloyd-Smith, and Ren Sun. 2016. “Adaptation in Protein Fitness Landscapes Is Facilitated by Indirect Paths.” *eLife* 5 (July).
- Wu, Nicholas C., Jakub Otwinowski, Andrew J. Thompson, Corwin M. Nycholat, Armita Nourmohammad, and Ian A. Wilson. 2020. “Major Antigenic Site B of Human Influenza H3N2 Viruses Has an Evolving Local Fitness Landscape.” *Nature Communications* 11 (1): 1233.
- Wünsche, Andrea, Duy M. Dinh, Rebecca S. Satterwhite, Carolina Diaz Arenas, Daniel M. Stoebel, and Tim F. Cooper. 2017. “Diminishing-Returns Epistasis Decreases Adaptability along an Evolutionary Trajectory.” *Nature Ecology & Evolution* 1 (4): 61.
- Yang, Gloria, Dave W. Anderson, Florian Baier, Elias Dohmen, Nansook Hong, Paul D. Carr, Shina Caroline Lynn Kamerlin, Colin J. Jackson, Erich Bornberg-Bauer, and Nobuhiko Tokuriki. 2019. “Higher-Order Epistasis Shapes the Fitness Landscape of a Xenobiotic-Degrading Enzyme.” *Nature Chemical Biology* 15 (11): 1120–28.

- Yang, Z., and J. P. Bielawski. 2000. "Statistical Methods for Detecting Molecular Adaptation." *Trends in Ecology & Evolution* 15 (12): 496–503.
- Yang, Ziheng. 2007. "PAML 4: Phylogenetic Analysis by Maximum Likelihood." *Molecular Biology and Evolution* 24 (8): 1586–91.
- Yang, Ziheng, Wendy S. W. Wong, and Rasmus Nielsen. 2005. "Bayes Empirical Bayes Inference of Amino Acid Sites under Positive Selection." *Molecular Biology and Evolution* 22 (4): 1107–18.
- Yubero, Pablo, Susanna Manrubia, and Jacobo Aguirre. 2017. "The Space of Genotypes Is a Network of Networks: Implications for Evolutionary and Extinction Dynamics." *Scientific Reports* 7 (1): 13813.
- Zanini, Fabio, and Richard A. Neher. 2012. "FFPopSim: An Efficient Forward Simulation Package for the Evolution of Large Populations." *Bioinformatics* 28 (24): 3332–33.
- Zan, Yanjun, Simon K. G. Forsberg, and Örjan Carlborg. 2018. "On the Relationship Between High-Order Linkage Disequilibrium and Epistasis." *G3* 8 (8): 2817–24.
- Zhao, Lei, and Brian Charlesworth. 2016. "Resolving the Conflict Between Associative Overdominance and Background Selection." *Genetics* 203 (3): 1315–34.
- Zheng, Jia, Joshua L. Payne, and Andreas Wagner. 2019. "Cryptic Genetic Variation Accelerates Evolution by Opening Access to Diverse Adaptive Peaks." *Science* 365 (6451): 347–53.
- Zhou, Juannan, and David M. McCandlish. 2020. "Minimum Epistasis Interpolation for Sequence-Function Relationships." *Nature Communications* 11 (1): 1782.
- Zou, Zhengting, and Jianzhi Zhang. 2015. "Are Convergent and Parallel Amino Acid Substitutions in Protein Evolution More Prevalent Than Neutral Expectations?" *Molecular Biology and Evolution* 32 (8): 2085–96.
- Zuckermandl, Emile, and Linus Pauling. 1965. "Evolutionary Divergence and Convergence in Proteins." In *Evolving Genes and Proteins*, 97–166.
- Zuk, Or, Eliana Hechter, Shamil R. Sunyaev, and Eric S. Lander. 2012. "The Mystery of Missing Heritability: Genetic Interactions Create Phantom Heritability." *Proc. Natl. Acad. Sci. U. S. A.* 109 (4): 1193–98.

# Appendix A

**Table A1. *S. commune* genomes assembly statistics.**

| <b>Sample id</b> | <b>Specimen voucher</b> | <b>Origin</b>  | <b># contigs</b> | <b>total length (bp)</b> | <b>GC %</b> | <b>N50</b> | <b>coverage</b> |
|------------------|-------------------------|----------------|------------------|--------------------------|-------------|------------|-----------------|
| <b>s1401</b>     | WS-M161                 | USA; Ann Arbor | 2,462            | 37,201,238               | 57.5        | 153,743    | 113.2           |
| <b>s14101</b>    | WS-M180                 | USA; Ann Arbor | 2,161            | 36,629,295               | 57.6        | 208,079    | 74.2            |
| <b>s14102</b>    | WS-M181                 | USA; Ann Arbor | 2,895            | 37,669,799               | 57.6        | 146,599    | 67.8            |
| <b>s14104</b>    | WS-M183                 | USA; Ann Arbor | 2,750            | 37,829,033               | 57.6        | 151,136    | 75.8            |
| <b>s14112</b>    | WS-M191                 | USA; Ann Arbor | 2,577            | 37,171,981               | 57.6        | 173,824    | 124.5           |
| <b>s1411</b>     | WS-M188                 | USA; Ann Arbor | 2,634            | 37,679,657               | 57.6        | 158,664    | 64.0            |
| <b>s1425</b>     | WS-M206                 | USA; Ann Arbor | 2,762            | 38,042,099               | 57.6        | 160,044    | 93.8            |
| <b>s1429</b>     | WS-M210                 | USA; Ann Arbor | 2,665            | 37,691,449               | 57.5        | 158,777    | 95.9            |
| <b>s1431</b>     | WS-M212                 | USA; Ann Arbor | 2,453            | 37,348,833               | 57.5        | 161,384    | 100.0           |
| <b>s1432</b>     | WS-M213                 | USA; Ann Arbor | 2,923            | 37,685,895               | 57.6        | 145,350    | 62.9            |
| <b>s1434</b>     | WS-M215                 | USA; Ann Arbor | 2,455            | 37,403,482               | 57.5        | 185,879    | 89.2            |
| <b>s1467</b>     | WS-M247                 | USA; Ann Arbor | 2,900            | 37,778,589               | 57.6        | 195,995    | 70.8            |
| <b>s1470</b>     | WS-M247                 | USA; Ann Arbor | 2,809            | 37,546,616               | 57.6        | 154,362    | 91.4            |
| <b>s1485</b>     | WS-M265                 | USA; Ann Arbor | 2,347            | 37,174,196               | 57.6        | 176,501    | 45.3            |
| <b>s1489</b>     | WS-M269                 | USA; Ann Arbor | 2,352            | 37,218,933               | 57.6        | 177,695    | 111.6           |
| <b>s1490</b>     | WS-M270                 | USA; Ann Arbor | 2,957            | 37,322,559               | 57.6        | 139,455    | 110.3           |
| <b>s1514</b>     | WS-M292                 | USA; Florida   | 2,460            | 37,328,560               | 57.6        | 157,189    | 71.5            |
| <b>X12</b>       | WS-M12                  | Russia; Moscow | 3,879            | 38,221,043               | 57.6        | 75,624     | 117.4           |
| <b>X17</b>       | WS-M18                  | Russia; Moscow | 3,738            | 37,604,751               | 57.6        | 71,000     | 105.9           |
| <b>X21</b>       | WS-M22                  | Russia; Moscow | 5,012            | 39,204,396               | 57.6        | 63,280     | 77.3            |
| <b>X27</b>       | WS-M28                  | Russia; Moscow | 3,571            | 37,399,774               | 57.6        | 71,903     | 78.2            |
| <b>X30</b>       | WS-M31                  | Russia; Moscow | 4,487            | 38,310,778               | 57.6        | 66,442     | 84.7            |
| <b>X69</b>       | WS-M70                  | Russia; Moscow | 3,965            | 38,348,248               | 57.6        | 70,802     | 76.7            |
| <b>X9</b>        | WS-M9                   | Russia; Moscow | 4,590            | 38,741,959               | 57.6        | 67,770     | 74.8            |

**Table A2. List of genes with pairs of physically adjacent protein sites being under higher LD than pairs of distant sites in *S. commune*.** P-values are calculated using chi-square test and adjusted using Benjamini-Hochberg multiple testing correction.

| rna            | # distant & high LD | # close & high LD | # distant & low LD | # close & low LD | OR    | p-value | q-value | aligned PDB ID |
|----------------|---------------------|-------------------|--------------------|------------------|-------|---------|---------|----------------|
| RUS population |                     |                   |                    |                  |       |         |         |                |
| 10789          | 6                   | 6                 | 85                 | 8                | 10.63 | 4.3E-04 | 8.4E-03 | 4N6Q A         |
| 7636           | 40                  | 16                | 263                | 11               | 9.56  | 5.2E-09 | 5.5E-07 | 4FQG A         |
| 11223          | 7                   | 9                 | 87                 | 12               | 9.32  | 1.0E-04 | 2.8E-03 | 1S3S G         |
| 9853           | 21                  | 22                | 261                | 32               | 8.54  | 8.7E-11 | 1.7E-08 | 4X00 A         |
| 17085          | 16                  | 6                 | 183                | 11               | 6.24  | 1.6E-03 | 2.1E-02 | 1NLT A         |
| 12357          | 56                  | 30                | 245                | 24               | 5.47  | 1.5E-08 | 1.4E-06 | 2GUY A         |
| 1037           | 26                  | 13                | 113                | 11               | 5.14  | 4.6E-04 | 8.8E-03 | 1K8F A         |
| 6153           | 39                  | 14                | 126                | 9                | 5.03  | 5.2E-04 | 9.7E-03 | 5GVH A         |
| 14273          | 22                  | 9                 | 244                | 21               | 4.75  | 7.5E-04 | 1.3E-02 | 3LCC A         |
| 5725           | 26                  | 15                | 312                | 38               | 4.74  | 1.6E-05 | 6.3E-04 | 1TA3 B         |
| 18561          | 69                  | 32                | 558                | 55               | 4.71  | 3.0E-10 | 4.8E-08 | 1KSG A         |
| 3052           | 91                  | 25                | 222                | 13               | 4.69  | 1.3E-05 | 5.4E-04 | 4U9V B         |
| 3876           | 38                  | 22                | 373                | 47               | 4.59  | 4.1E-07 | 3.0E-05 | 1W63 A         |
| 4779           | 80                  | 26                | 873                | 63               | 4.50  | 1.6E-09 | 2.1E-07 | 2GJL A         |
| 16912          | 172                 | 54                | 1383               | 99               | 4.39  | 9.1E-17 | 8.8E-14 | 1WKR A         |
| 14670          | 25                  | 8                 | 273                | 20               | 4.37  | 2.2E-03 | 2.6E-02 | 6C6N A         |
| 14338          | 110                 | 28                | 150                | 9                | 4.24  | 2.8E-04 | 6.6E-03 | 6J3E A         |
| 8942           | 35                  | 10                | 279                | 19               | 4.20  | 1.1E-03 | 1.6E-02 | 3DH1 A         |
| 3214           | 37                  | 9                 | 189                | 11               | 4.18  | 4.4E-03 | 4.2E-02 | 1SZN A         |
| 1413           | 78                  | 24                | 253                | 19               | 4.10  | 1.8E-05 | 6.8E-04 | 5L3Q B         |
| 7650           | 75                  | 29                | 201                | 19               | 4.09  | 1.2E-05 | 5.1E-04 | 5EBE B         |
| 1071           | 178                 | 54                | 1002               | 75               | 4.05  | 1.0E-13 | 2.4E-11 | 1SXJ D         |
| 18096          | 42                  | 16                | 462                | 45               | 3.91  | 3.7E-05 | 1.2E-03 | 1WPX A         |
| 13142          | 57                  | 9                 | 562                | 23               | 3.86  | 1.6E-03 | 2.1E-02 | 5GHE A         |
| 16593          | 118                 | 34                | 954                | 72               | 3.82  | 1.7E-09 | 2.1E-07 | 3WDO A         |
| 14325          | 133                 | 28                | 621                | 36               | 3.63  | 1.1E-06 | 6.8E-05 | 5U03 A         |
| 10827          | 25                  | 11                | 286                | 35               | 3.60  | 2.1E-03 | 2.5E-02 | 2IHO A         |
| 10077          | 38                  | 11                | 397                | 32               | 3.59  | 1.3E-03 | 2.0E-02 | 3AKF A         |

|       |     |     |      |     |      |         |         |        |
|-------|-----|-----|------|-----|------|---------|---------|--------|
| 9626  | 47  | 9   | 468  | 25  | 3.58 | 3.2E-03 | 3.4E-02 | 2CVF A |
| 10648 | 59  | 20  | 587  | 56  | 3.55 | 1.4E-05 | 5.6E-04 | 3I83 A |
| 8095  | 192 | 60  | 1686 | 151 | 3.49 | 3.2E-14 | 1.0E-11 | 2YMU A |
| 5372  | 88  | 33  | 914  | 99  | 3.46 | 3.3E-08 | 2.9E-06 | 1RGI G |
| 14269 | 105 | 23  | 426  | 27  | 3.46 | 4.1E-05 | 1.3E-03 | 5UJ8 E |
| 14404 | 54  | 16  | 374  | 33  | 3.36 | 4.0E-04 | 8.0E-03 | 4IDA A |
| 17420 | 32  | 12  | 340  | 39  | 3.27 | 2.4E-03 | 2.8E-02 | 1W9P A |
| 6507  | 63  | 13  | 620  | 40  | 3.20 | 9.9E-04 | 1.6E-02 | 1AUA A |
| 9423  | 60  | 11  | 401  | 23  | 3.20 | 4.4E-03 | 4.2E-02 | 4Y42 A |
| 7878  | 41  | 17  | 446  | 58  | 3.19 | 3.5E-04 | 8.0E-03 | 4TYW A |
| 3307  | 87  | 17  | 720  | 45  | 3.13 | 2.3E-04 | 5.8E-03 | 6DVH A |
| 6285  | 55  | 13  | 565  | 43  | 3.11 | 1.4E-03 | 2.1E-02 | 2VWS A |
| 10049 | 68  | 16  | 668  | 51  | 3.08 | 4.0E-04 | 8.0E-03 | 1KH4 A |
| 6148  | 628 | 131 | 1361 | 93  | 3.05 | 1.6E-15 | 7.7E-13 | 4CHT A |
| 14511 | 42  | 13  | 414  | 44  | 2.91 | 3.7E-03 | 3.7E-02 | 3HG7 A |
| 2522  | 46  | 16  | 492  | 60  | 2.85 | 1.5E-03 | 2.1E-02 | 5L0R A |
| 5375  | 45  | 18  | 448  | 63  | 2.84 | 9.6E-04 | 1.5E-02 | 2WZO A |
| 73    | 68  | 20  | 715  | 74  | 2.84 | 2.5E-04 | 6.2E-03 | 1WPX A |
| 12131 | 210 | 35  | 816  | 48  | 2.83 | 8.7E-06 | 4.0E-04 | 6GKV A |
| 1097  | 63  | 17  | 534  | 51  | 2.83 | 1.1E-03 | 1.6E-02 | 2IW0 A |
| 18360 | 174 | 35  | 924  | 66  | 2.82 | 3.6E-06 | 2.0E-04 | 1ULT A |
| 5930  | 57  | 16  | 594  | 60  | 2.78 | 1.5E-03 | 2.1E-02 | 2PXX A |
| 8261  | 112 | 42  | 930  | 129 | 2.70 | 9.4E-07 | 6.5E-05 | 4QNW A |
| 18092 | 83  | 22  | 794  | 78  | 2.70 | 2.5E-04 | 6.1E-03 | 4QJY A |
| 1060  | 65  | 16  | 533  | 50  | 2.62 | 3.2E-03 | 3.3E-02 | 3WXB A |
| 17037 | 95  | 19  | 918  | 70  | 2.62 | 7.4E-04 | 1.3E-02 | 5YHP A |
| 15353 | 109 | 26  | 944  | 86  | 2.62 | 1.0E-04 | 2.8E-03 | 3WNV A |
| 7784  | 120 | 15  | 929  | 45  | 2.58 | 3.5E-03 | 3.5E-02 | 3L4G B |
| 10236 | 182 | 35  | 1810 | 135 | 2.58 | 3.5E-06 | 2.0E-04 | 1Q6X A |
| 2011  | 390 | 17  | 889  | 67  | 2.51 | 2.4E-03 | 2.7E-02 | 3A1K A |
| 8572  | 167 | 32  | 976  | 76  | 2.46 | 8.1E-05 | 2.4E-03 | 4AH6 A |
| 14282 | 153 | 18  | 1251 | 60  | 2.45 | 2.0E-03 | 2.4E-02 | 3QM4 A |
| 7836  | 83  | 17  | 685  | 58  | 2.42 | 4.4E-03 | 4.2E-02 | 1DQW A |
| 3610  | 196 | 45  | 1946 | 185 | 2.42 | 1.2E-06 | 7.3E-05 | 4BKX B |
| 8725  | 71  | 18  | 669  | 71  | 2.39 | 4.0E-03 | 4.0E-02 | 6G6M A |
| 11096 | 64  | 20  | 657  | 87  | 2.36 | 3.0E-03 | 3.2E-02 | 3AKF A |

|                |     |     |      |     |      |         |         |        |
|----------------|-----|-----|------|-----|------|---------|---------|--------|
| 6520           | 86  | 28  | 599  | 84  | 2.32 | 8.3E-04 | 1.4E-02 | 4K3A A |
| 12399          | 610 | 107 | 994  | 76  | 2.29 | 1.4E-07 | 1.1E-05 | 1JZQ A |
| 8945           | 147 | 32  | 558  | 53  | 2.29 | 7.9E-04 | 1.3E-02 | 3E5M A |
| 1744           | 148 | 45  | 616  | 83  | 2.26 | 9.7E-05 | 2.8E-03 | 1C7J A |
| 16360          | 134 | 34  | 1297 | 149 | 2.21 | 2.0E-04 | 5.3E-03 | 3WTC A |
| 12853          | 168 | 31  | 1634 | 139 | 2.17 | 3.8E-04 | 8.0E-03 | 6H7D A |
| 14137          | 118 | 27  | 1175 | 127 | 2.12 | 1.7E-03 | 2.1E-02 | 6C5B A |
| 7106           | 124 | 22  | 1167 | 98  | 2.11 | 4.4E-03 | 4.2E-02 | 5K8E A |
| 1523           | 77  | 29  | 402  | 72  | 2.10 | 4.4E-03 | 4.2E-02 | 5Y1B A |
| 2779           | 127 | 27  | 1155 | 120 | 2.05 | 2.8E-03 | 3.0E-02 | 1SXJ B |
| 4275           | 555 | 87  | 957  | 74  | 2.03 | 2.5E-05 | 8.4E-04 | 5VC7 A |
| 17782          | 184 | 53  | 581  | 83  | 2.02 | 4.1E-04 | 8.0E-03 | 4QNW A |
| 4829           | 350 | 62  | 504  | 46  | 1.94 | 1.7E-03 | 2.1E-02 | 5MXC A |
| 9827           | 237 | 42  | 2273 | 209 | 1.93 | 3.9E-04 | 8.0E-03 | 2VJY A |
| 1520           | 153 | 32  | 1498 | 163 | 1.92 | 2.6E-03 | 2.9E-02 | 3PQV A |
| 4468           | 238 | 48  | 1409 | 148 | 1.92 | 3.6E-04 | 8.0E-03 | 1V9L A |
| 8360           | 852 | 74  | 884  | 40  | 1.92 | 1.5E-03 | 2.1E-02 | 5YLW A |
| 16987          | 154 | 34  | 1482 | 174 | 1.88 | 2.8E-03 | 3.0E-02 | 3LWT X |
| 935            | 104 | 38  | 1036 | 203 | 1.86 | 3.0E-03 | 3.2E-02 | 3FGA A |
| 11732          | 118 | 40  | 1076 | 196 | 1.86 | 2.3E-03 | 2.7E-02 | 4C2L A |
| 15295          | 152 | 41  | 1326 | 193 | 1.85 | 1.7E-03 | 2.1E-02 | 4A69 A |
| 6753           | 246 | 36  | 1868 | 151 | 1.81 | 3.4E-03 | 3.5E-02 | 1SXJ C |
| 13863          | 319 | 86  | 414  | 62  | 1.80 | 1.6E-03 | 2.1E-02 | 6F43 A |
| USA population |     |     |      |     |      |         |         |        |
| 14970          | 13  | 6   | 160  | 8   | 9.23 | 1.8E-04 | 1.3E-02 | 2VFR A |
| 1536           | 12  | 13  | 184  | 23  | 8.67 | 4.6E-07 | 2.0E-04 | 6AHR E |
| 3618           | 9   | 10  | 139  | 24  | 6.44 | 2.1E-04 | 1.4E-02 | 5LCL B |
| 18366          | 44  | 15  | 486  | 41  | 4.04 | 3.5E-05 | 4.5E-03 | 1UPU D |
| 8253           | 44  | 12  | 467  | 35  | 3.64 | 5.8E-04 | 3.4E-02 | 6F87 A |
| 9241           | 56  | 23  | 624  | 81  | 3.16 | 2.6E-05 | 4.2E-03 | 1YCD A |
| 1743           | 49  | 19  | 510  | 64  | 3.09 | 2.1E-04 | 1.4E-02 | 4PEH A |
| 64             | 85  | 27  | 905  | 101 | 2.85 | 1.9E-05 | 3.6E-03 | 2B4Q A |
| 14128          | 77  | 28  | 804  | 103 | 2.84 | 1.9E-05 | 3.6E-03 | 2X8R A |
| 10841          | 120 | 20  | 1166 | 69  | 2.82 | 1.5E-04 | 1.2E-02 | 3TIK A |
| 17174          | 96  | 27  | 679  | 73  | 2.62 | 1.4E-04 | 1.2E-02 | 5EY6 A |

|       |     |    |      |     |      |         |         |        |
|-------|-----|----|------|-----|------|---------|---------|--------|
| 5725  | 73  | 30 | 799  | 126 | 2.61 | 5.9E-05 | 6.9E-03 | 1TA3 B |
| 10834 | 90  | 31 | 936  | 124 | 2.60 | 3.3E-05 | 4.5E-03 | 2QB6 A |
| 1267  | 117 | 29 | 1150 | 116 | 2.46 | 1.0E-04 | 1.0E-02 | 4CPD A |
| 9614  | 149 | 44 | 1550 | 187 | 2.45 | 1.9E-06 | 6.0E-04 | 2VGL B |
| 6148  | 487 | 75 | 4438 | 284 | 2.41 | 1.2E-10 | 1.5E-07 | 4CHT A |
| 161   | 227 | 52 | 2206 | 210 | 2.41 | 2.0E-07 | 1.3E-04 | 5DNC A |
| 621   | 105 | 35 | 1060 | 152 | 2.32 | 9.1E-05 | 9.7E-03 | 3WG6 A |
| 14368 | 124 | 26 | 1008 | 92  | 2.30 | 7.4E-04 | 4.1E-02 | 2X1C A |
| 9215  | 161 | 56 | 1546 | 243 | 2.21 | 3.0E-06 | 7.6E-04 | 4QNW A |
| 13117 | 140 | 44 | 1401 | 226 | 1.95 | 4.5E-04 | 2.8E-02 | 1W9P A |
| 3876  | 232 | 51 | 2244 | 259 | 1.90 | 1.5E-04 | 1.2E-02 | 1W63 A |

**Table A3. Senescence and entrenchment in mitochondrial genes of Metazoa.** The table shows mitochondrial sites which show statistically significant entrenchment (regression coefficient < 0) or senescence (regression coefficient > 0). Regression coefficients are calculated using binomial logistic regression (for more details, see main text).

| AC (human) | protein name | position | p-value  | regression coefficient | q-value  |
|------------|--------------|----------|----------|------------------------|----------|
| P00846.1   | ATP6         | 5        | 9.25E-07 | -29.13                 | 5.91E-04 |
| P00846.1   | ATP6         | 10       | 8.40E-07 | 11.30                  | 5.75E-04 |
| P00846.1   | ATP6         | 26       | 3.86E-07 | -10.14                 | 3.08E-04 |
| P00846.1   | ATP6         | 69       | 1.25E-04 | -37.11                 | 3.14E-02 |
| P00846.1   | ATP6         | 70       | 6.59E-05 | -16.23                 | 1.97E-02 |
| P00846.1   | ATP6         | 71       | 4.45E-14 | -9.02                  | 2.13E-10 |
| P00846.1   | ATP6         | 80       | 2.44E-08 | 7.34                   | 3.34E-05 |
| P00846.1   | ATP6         | 105      | 5.21E-06 | 9.09                   | 2.27E-03 |
| P00846.1   | ATP6         | 111      | 4.72E-08 | 7.73                   | 5.65E-05 |
| P00846.1   | ATP6         | 116      | 2.09E-04 | -8.38                  | 4.17E-02 |
| P00846.1   | ATP6         | 125      | 1.46E-08 | 6.34                   | 2.81E-05 |
| P00846.1   | ATP6         | 141      | 1.86E-08 | 5.95                   | 2.98E-05 |

|          |      |     |          |        |          |
|----------|------|-----|----------|--------|----------|
| P00846.1 | ATP6 | 178 | 1.47E-04 | -7.60  | 3.61E-02 |
| P00846.1 | ATP6 | 179 | 1.72E-09 | 4.23   | 5.49E-06 |
| P00846.1 | ATP6 | 224 | 1.89E-04 | -10.28 | 3.97E-02 |
| P00395.1 | COX1 | 52  | 6.69E-06 | 6.37   | 2.71E-03 |
| P00395.1 | COX1 | 132 | 2.30E-06 | -8.60  | 1.22E-03 |
| P00395.1 | COX1 | 185 | 1.88E-05 | -16.80 | 6.44E-03 |
| P00395.1 | COX1 | 256 | 2.10E-05 | -30.61 | 6.93E-03 |
| P00395.1 | COX1 | 336 | 6.79E-06 | -13.88 | 2.71E-03 |
| P00395.1 | COX1 | 389 | 1.07E-04 | -12.40 | 2.86E-02 |
| P00395.1 | COX1 | 399 | 1.71E-06 | -30.26 | 9.61E-04 |
| P00395.1 | COX1 | 414 | 3.58E-06 | -8.05  | 1.63E-03 |
| P00395.1 | COX1 | 415 | 2.58E-06 | -9.15  | 1.30E-03 |
| P00395.1 | COX1 | 448 | 7.98E-05 | -14.43 | 2.25E-02 |
| P00395.1 | COX1 | 491 | 1.92E-04 | -37.47 | 3.97E-02 |
| P00403.1 | COX2 | 21  | 1.75E-04 | 41.65  | 3.97E-02 |
| P00403.1 | COX2 | 37  | 1.03E-07 | 6.56   | 1.10E-04 |
| P00403.1 | COX2 | 51  | 2.74E-05 | 30.72  | 8.48E-03 |
| P00403.1 | COX2 | 56  | 3.91E-18 | 12.63  | 3.74E-14 |
| P00403.1 | COX2 | 57  | 1.85E-05 | -28.14 | 6.44E-03 |
| P00403.1 | COX2 | 112 | 2.28E-04 | 35.40  | 4.47E-02 |
| P00403.1 | COX2 | 217 | 2.40E-09 | 6.28   | 5.75E-06 |
| P00414.2 | COX3 | 47  | 1.52E-04 | 10.36  | 3.64E-02 |
| P00414.2 | COX3 | 60  | 1.18E-04 | 3.37   | 3.05E-02 |
| P00414.2 | COX3 | 166 | 7.85E-07 | -9.71  | 5.75E-04 |
| P00156.2 | CYTB | 4   | 3.42E-06 | -4.17  | 1.63E-03 |
| P00156.2 | CYTB | 10  | 1.93E-04 | -51.60 | 3.97E-02 |
| P00156.2 | CYTB | 38  | 1.95E-04 | 16.61  | 3.97E-02 |
| P00156.2 | CYTB | 81  | 1.10E-05 | -59.64 | 4.21E-03 |
| P00156.2 | CYTB | 93  | 1.90E-04 | -17.48 | 3.97E-02 |
| P00156.2 | CYTB | 121 | 2.17E-07 | 28.84  | 1.89E-04 |
| P00156.2 | CYTB | 180 | 1.35E-06 | 8.57   | 8.07E-04 |
| P00156.2 | CYTB | 183 | 7.39E-05 | 14.84  | 2.15E-02 |
| P00156.2 | CYTB | 218 | 1.83E-05 | -23.65 | 6.44E-03 |
| P00156.2 | CYTB | 225 | 8.90E-05 | 9.37   | 2.44E-02 |
| P00156.2 | CYTB | 235 | 1.70E-07 | -42.63 | 1.63E-04 |
| P00156.2 | CYTB | 337 | 1.89E-04 | -20.53 | 3.97E-02 |

|          |      |     |          |       |          |
|----------|------|-----|----------|-------|----------|
| P00156.2 | CYTB | 346 | 2.65E-05 | -4.01 | 8.47E-03 |
|----------|------|-----|----------|-------|----------|