

An extension of the Walsh-Hadamard transform to calculate and model epistasis in genetic landscapes of arbitrary shape and complexity

Andre J. Faure^{1*}, Ben Lehner^{1,2,3,4}, Verónica Miró Pina¹, Claudia Serrano Colome¹, and Donat Waghorn^{1,2}

¹Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona, Spain.

²Universitat Pompeu Fabra (UPF), Barcelona, Spain

³Institució Catalana de Recerca i estudis Avançats (ICREA), Barcelona, Spain

⁴Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK

*Corresponding author email: andre.faure@crg.eu (AJF)

Abstract

Accurate models describing the relationship between genotype and phenotype are necessary in order to understand and predict how mutations to biological sequences affect the fitness and evolution of living organisms. The apparent abundance of epistasis (genetic interactions), both between and within genes, complicates this task and how to build mechanistic models that incorporate epistatic coefficients (genetic interaction terms) is an open question. The Walsh-Hadamard transform represents a rigorous computational framework for calculating and modeling epistatic interactions at the level of individual genotypic values (known as genetical, biological or physiological epistasis), and can therefore be used to address fundamental questions related to sequence-to-function encodings. However, one of its main limitations is that it can only accommodate two alleles (amino acid or nucleotide states) per sequence position. In this paper we provide an extension of the Walsh-Hadamard transform that allows the calculation and modeling of background-averaged epistasis (also known as ensemble epistasis) in genetic landscapes with an arbitrary number of states per position (20 for amino acids, 4 for nucleotides, etc.). We also provide a recursive formula for the inverse matrix and then derive formulae to directly extract any element of either matrix without having to rely on the computationally intensive task of constructing or inverting large matrices. Finally, we demonstrate the utility of our theory by using it to model epistasis within a combinatorially complete multiallelic genetic landscape of a tRNA, revealing that both pairwise and higher-order genetic interactions are enriched between physically interacting positions.

30 **Author Summary**

31 An important question in genetics is how the effects of mutations combine to alter phenotypes. Genetic in-
32 teractions (epistasis) describe non-additive effects of pairs of mutations, but can also involve higher-order
33 (three- and four-way etc.) combinations. Quantifying higher-order interactions is experimentally very chal-
34 lenging requiring a large number of measurements. Techniques based on deep mutational scanning (DMS,
35 also known as MPRA and MAVEs) represent valuable sources of data to study epistasis. However, the
36 best way to extract the relevant pair-wise and higher-order epistatic coefficients (genetic interaction terms)
37 from this data for the task of phenotypic prediction remains an unresolved problem. The Walsh-Hadamard
38 transform represents a rigorous computational framework for calculating and modeling epistatic interactions
39 at the level of individual genotypic values. Critically, this formalism currently only allows for two alleles
40 (amino acid or nucleotide states) per sequence position, hampering applications in more biologically realis-
41 tic scenarios. Here we present an extension of the Walsh-Hadamard transform that overcomes this limitation
42 and demonstrate the utility of our theory by using it to model epistasis within a combinatorially complete
43 multiallelic genetic landscape of a tRNA.

44 **Introduction**

45 A fundamental challenge in biology is to understand and predict how changes (or mutations) to biologi-
46 cal sequences (DNA, RNA, proteins) affect their molecular function and ultimately the phenotype of living
47 organisms. The phenomenon of ‘epistasis’ (genetic interactions) – broadly defined as the dependence of mu-
48 tational effects on the genetic context in which they occur [1, 2, 3] – is widespread in biological systems, yet
49 knowledge of the underlying mechanisms remains limited. Defining the extent of epistasis and better under-
50 standing of its origins has relevance in fields ranging from genetic prediction, molecular evolution, infectious
51 disease and cancer drug development, to biomolecular structure determination and protein engineering [3].

52 Evolutionarily related sequences, natural genetic variation within populations, and more recently results of
53 techniques such as deep mutational scanning (DMS) [4] – also known as massively parallel reporter as-
54 says (MPRAs) and multiplex assays of variant effect (MAVEs) – represent valuable sources of data to study
55 epistasis [5, 1]. In particular, DMS enables the systematic measurement of mutational effects across entire
56 combinatorially complete genetic landscapes [5, 6, 7, 8, 9, 10, 11, 12, 13]. Importantly, the typical use of
57 engineered genotypes, haploid individuals and near-identical environmental (laboratory) conditions in these
58 experiments allows population genetic considerations – such as dominance, variable allele frequencies and
59 linkage disequilibrium – to be ignored [14]. In other words, measurements obtained from deep mutational
60 scanning and related methods permit the modeling of epistasis in the mechanistic sense (sequence-to-function
61 encoding) rather than in the evolutionary sense i.e. at the population genetic level. Nevertheless, precisely
62 how to extract the most biologically relevant pairwise and higher-order epistatic coefficients (genetic inter-
63 action terms) from this type of data is an unresolved problem.

64 Quantitative definitions of epistasis vary among fields, but it has been argued that one particular formula-

65 tion termed ‘background-averaged’ epistasis, also known as ‘ensemble’ epistasis [1, 12], may provide the
66 most useful information on the epistatic structure of biological systems [2]. The underlying rationale is that
67 by averaging the effects of mutations across many different genetic backgrounds (contexts), the method is
68 robust to local idiosyncrasies in the relationship between genotype and phenotype. It has been previously
69 pointed out that the definition of background-averaged epistasis is conceptually similar to that of ‘statisti-
70 cal epistasis’ attributed to Fisher, but instead of measuring the average effect of allele substitutions against
71 the population average genetic background i.e. averaging over all genotypes present in a given population
72 (taking into account their individual frequencies), the approach instead averages over all possible genotypes
73 (assuming equal genotype weights) [1, 2].

74 The current mathematical formalism of background-averaged epistasis is based on the Walsh-Hadamard
75 transform [2]. Interestingly, although widely used in physics and engineering, the Walsh-Hadamard trans-
76 form was first applied to non-biological fitness landscapes in the field of genetic algorithms (GA) [15], subse-
77 quently being proposed as the basis of a framework for the computation of higher-order epistasis in empirical
78 settings [16]. However, the Walsh-Hadamard transform can only accommodate two alleles (amino acid or
79 nucleotide states) per sequence position, with no extension to multiallelic landscapes (cardinality greater
80 than two) yet made, as confirmed by multiple recent reports [2, 17, 18, 19]. Alternative implementations for
81 multiallelic landscapes either rely on ‘one-hot encoding’ elements of larger alphabets as biallelic sequences –
82 requiring the manipulation of prohibitively large Walsh-Hadamard matrices – or constructing graph Fourier
83 bases [18], which is mathematically complex and provides no straightforward way to interpret epistatic co-
84 efficients. The result is that the application of background-averaged epistasis has been severely limited and
85 its properties remain largely unexplored in more biologically realistic scenarios.

86 In this work we provide an extension of the Walsh-Hadamard transform that allows the calculation and mod-
87 eling of background-averaged epistasis in genetic landscapes with an arbitrary number of states (20 for amino
88 acids, 4 for nucleotides, etc.). We also provide a recursive formula for the inverse matrix, which is required to
89 infer epistatic coefficients using regression. Furthermore, we derive convenient formulae to directly extract
90 any element of either matrix without having to rely on the computationally intensive task of constructing
91 or inverting large matrices. Lastly, we apply these formulae to the analysis of a multiallelic DMS dataset,
92 demonstrating that sparse models inferred from the background-averaged representation (embedding) of the
93 underlying genetic landscape more regularly include epistatic terms corresponding to direct physical inter-
94 actions.

95 **Results**

96 **Extension of the Walsh-Hadamard transform to multiallelic landscapes**

97 In this work, a genotype sequence is represented as a one-dimensional ordering of monomers, each of which
98 can take on s possible states (or alleles), for example $s = 4$ for nucleotide sequences or $s = 20$ for amino
99 acid sequences. Without loss of generality, the s states can be labelled $0, 1, 2, \dots, s - 1$, where 0 denotes the

100 wild-type allele. We are going to consider genotype sequences of length $n \in \mathbb{N}$, i.e. sequences taking values
 101 in S^n , where $S := \{0, 1, \dots, s - 1\}$.

102 Each genotype $\vec{i} \in S^n$ is associated with its phenotype $y_{\vec{i}}$. Note that here we use the term ‘phenotype’ as
 103 shorthand for ‘molecular phenotype score’ from a quantitative laboratory assay (DMS) reporting on a molec-
 104 ular function for each genotype of interest. In quantitative genetics terminology this might be referred to as
 105 ‘genotypic value’ because environmental deviation is negligible due to the controlled nature of the experi-
 106 ments, but our subject here is the macromolecule not an individual from a population [14]. In the context of
 107 empirical genotype-phenotype landscapes, the phenotypic effect of a genotype \vec{i} is typically measured with
 108 respect to the wild-type, i.e. it is given by $y_{\vec{i}} - y_{(0,\dots,0)}$.

109 It is important to emphasize that in what follows we implicitly restrict ourselves to the haploid reference
 110 base, because our primary goal is the modeling of sequence-to-function encodings for *individual* genotype
 111 sequences – for the ultimate purpose of understanding and engineering macromolecules – not the modeling
 112 of sequence evolution or quantification of sources of phenotypic variance in populations.

113 If the phenotypic effects of individual mutations were independent, they would be additive, meaning that
 114 the phenotypic effect of $\vec{i} = (i_1, \dots, i_n)$ would be the sum of the phenotypic effects of the single mutants
 115 $(i_1, 0, \dots, 0), \dots, (0, \dots, 0, i_n)$. The epistatic coefficient quantifies how much the observed phenotypic effect
 116 of \vec{i} deviates from this assumption. In the case of background-averaged epistasis, we quantify the interac-
 117 tions between a set of mutations by averaging over all possible genotypes for the remaining positions in the
 118 sequence. For example, if $n = 3$ and $s = 2$, the pairwise epistatic coefficient involving the mutations at posi-
 119 tions 2 and 3 is calculated by averaging over all states (backgrounds) for the remaining positions, in this case
 120 given by the two states of the first position (* denotes the positions at which the averaging is performed), i.e.

$$\begin{aligned} \epsilon_{(*,1,1)} &= \frac{1}{2} \left([(y_{(1,1,1)} - y_{(1,0,0)}) - (y_{(1,1,0)} - y_{(1,0,0)}) - (y_{(1,0,1)} - y_{(1,0,0)})] + \right. \\ &\quad \left. [(y_{(0,1,1)} - y_{(0,0,0)}) - (y_{(0,1,0)} - y_{(0,0,0)}) - (y_{(0,0,1)} - y_{(0,0,0)})] \right) \\ &= \frac{1}{2} \left([y_{(1,1,1)} - y_{(1,1,0)} - y_{(1,0,1)} + y_{(1,0,0)}] + [y_{(0,1,1)} - y_{(0,1,0)} - y_{(0,0,1)} + y_{(0,0,0)}] \right). \end{aligned}$$

121 More generally, in [2] it is shown that for $s = 2$ and any sequence length n , phenotypic effects can be
 122 decomposed into background-averaged epistatic coefficients with

$$\bar{\epsilon}_n = \hat{V}_n \cdot \hat{H}_n \cdot \bar{y}_n,$$

123 where \bar{y}_n is the vector $(y_{\vec{i}}, \vec{i} \in [0, 1]^n)$, $\bar{\epsilon}_n$ is the vector $(\epsilon_{\vec{j}}, j \in [*, 1]^n)$ and \hat{H}_n and \hat{V}_n are $2^n \times 2^n$ matrices
 124 defined recursively as follows:

$$\hat{H}_{n+1} = \begin{pmatrix} \hat{H}_n & \hat{H}_n \\ \hat{H}_n & -\hat{H}_n \end{pmatrix} \quad \hat{H}_0 = 1,$$

$$\hat{V}_{n+1} = \begin{pmatrix} \frac{1}{2}\hat{V}_n & 0 \\ 0 & -\hat{V}_n \end{pmatrix} \quad \hat{V}_0 = 1.$$

125 The matrix \hat{H} is known as the Walsh-Hadamard transform [20, 21] and \hat{V} is a diagonal weighting (or nor-
 126 malisation) matrix to correct the sign and account for averaging over different numbers of backgrounds as a
 127 function of epistatic order [2].

128 In this work, we provide an extension of this theory to describe background-averaged epistasis for sequences
 129 with an arbitrary number of states s . Before writing a general formula, we consider the simplest possible
 130 multi-state (multiallelic) landscape i.e. a sequence of length $n = 1$ with $s = 3$,

$$\begin{pmatrix} \varepsilon_{(*)} \\ \varepsilon_{(1)} \\ \varepsilon_{(2)} \end{pmatrix} = \begin{pmatrix} 1/3 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 1 \\ 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} y_{(0)} \\ y_{(1)} \\ y_{(2)} \end{pmatrix} := V_1 \cdot H_1 \cdot \bar{y}_1.$$

131 Consistent with the definition of background-averaged epistasis for biallelic landscapes [2], the zeroth-order
 132 epistatic coefficient $\varepsilon_{(*)}$ is the mean phenotypic value across all genotypes and the first-order epistatic coef-
 133 ficients $\varepsilon_{(1)}$ and $\varepsilon_{(2)}$ are simply the respective individual phenotypic effects of genotypes $y_{(1)}$ and $y_{(2)}$ with
 134 respect to the wild-type. However, the key feature of H_1 for multiallelic landscapes – and where it departs
 135 from the canonical Walsh-Hadamard transform – is the introduction of zero elements to exclude phenotypes
 136 that are irrelevant for the calculation of a given epistatic coefficient. In other words, these phenotypes are
 137 excluded because they correspond neither to relevant intermediate genotypes nor alternative genetic back-
 138 grounds. We remind the reader that as we are interested in phenotypes at the level of *individual* genotypes,
 139 i.e. the haploid reference base, additive effects of different alleles at the same position (locus) are irrelevant
 140 and can be ignored.

141 If we now consider a sequence of length $n = 2$ with $s = 3$, then the H_2 and V_2 matrices become 9×9 ($s^n \times s^n$)
 142 and can be constructed from recurring to the case $n = 1$ above, giving

$$\begin{pmatrix} \varepsilon_{(*,*)} \\ \varepsilon_{(*,1)} \\ \varepsilon_{(*,2)} \\ \varepsilon_{(1,*)} \\ \varepsilon_{(1,1)} \\ \varepsilon_{(1,2)} \\ \varepsilon_{(2,*)} \\ \varepsilon_{(2,1)} \\ \varepsilon_{(2,2)} \end{pmatrix} = \begin{pmatrix} \frac{1}{9} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -\frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -\frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -\frac{1}{3} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -\frac{1}{3} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & -1 & 0 & 1 & -1 & 0 & 1 & -1 & 0 \\ 1 & 0 & -1 & 1 & 0 & -1 & 1 & 0 & -1 \\ 1 & 1 & 1 & -1 & -1 & -1 & 0 & 0 & 0 \\ 1 & -1 & 0 & -1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & -1 & -1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & -1 & -1 & -1 \\ 1 & -1 & 0 & 0 & 0 & 0 & -1 & 1 & 0 \\ 1 & 0 & -1 & 0 & 0 & 0 & -1 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} y_{(0,0)} \\ y_{(0,1)} \\ y_{(0,2)} \\ y_{(1,0)} \\ y_{(1,1)} \\ y_{(1,2)} \\ y_{(2,0)} \\ y_{(2,1)} \\ y_{(2,2)} \end{pmatrix} \\ := V_2 \cdot H_2 \cdot \bar{y}_2,$$

143 where the colors highlight the block structure of the matrices. In V_2 , the red square corresponds to $\frac{1}{s}V_1$ and
 144 the light red squares to $-V_1$. In H_2 , the gray squares correspond to H_1 and the blue squares to $-H_1$. In Table
 145 1 we show the results of background-averaged epistatic coefficients calculated by applying the above formula
 146 to an empirical multiallelic landscape with $n = 2$ and $s = 3$ [6].

Nucleic acid sequence	Base $s = 3$ representation	Phenotypic effect \bar{y}_2	Epistatic term $\bar{\epsilon} = V_2 \cdot H_2 \cdot \bar{y}_2$
GC	(0,0)	0	-0.17
GA	(0,1)	-0.14	-0.21
GT	(0,2)	-0.07	0.02
AC	(1,0)	-0.13	-0.24
AA	(1,1)	-0.8	-0.53
AT	(1,2)	-0.01	0.19
TC	(2,0)	-0.19	-0.05
TA	(2,1)	0	0.33
TT	(2,2)	-0.18	0.08

Table 1: Interaction terms based on background-averaged epistasis ($\bar{\epsilon}$) for an empirical multiallelic genotype-phenotype landscape consisting of all combinations of two mutations each at positions 6 and 66 in the tRNA-Arg(CCU) [6], i.e. $n = 2$ and $s = 3$. The first two columns indicate nucleic acid sequences and their base 3 representations. Here the ‘GC’ reference (wild-type) genotype corresponds to that of *S. cerevisiae*, denoted by (0, 0). The second two columns show the measured phenotypic effects and corresponding background-averaged epistatic coefficients. See Results for a regression analysis of the entire dataset.

147 More generally, for any value of s , when $n = 1$,

$$\begin{pmatrix} \epsilon_{(*)} \\ \epsilon_{(1)} \\ \epsilon_{(2)} \\ \vdots \\ \epsilon_{(s-1)} \end{pmatrix} = \begin{pmatrix} 1/s & 0 & 0 & \dots & 0 \\ 0 & -1 & 0 & \dots & 0 \\ 0 & 0 & -1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & \dots & 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & -1 & 0 & \dots & 0 \\ 1 & 0 & -1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 1 & 0 & \dots & 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} y_{(0)} \\ y_{(1)} \\ y_{(2)} \\ \vdots \\ y_{(s-1)} \end{pmatrix} := V_1 \cdot H_1 \cdot \bar{y}_1,$$

148 where $\epsilon_{(*)}$ corresponds to averaging phenotypes over all possible genotypes and the remaining coefficients
 149 simply correspond to the phenotypic effects of each mutation.

150 For $n = 2$, we have to consider different combinations of mutations in both positions. In this case, the
 151 phenotypes can be written as

$$y_{(0,0)}, y_{(0,1)}, \dots, y_{(0,(s-1))}, y_{(1,0)}, \dots, y_{(1,(s-1))}, \dots, y_{((s-1),0)}, \dots, y_{((s-1),(s-1))}.$$

152 A natural ordering of the phenotypes is given by interpreting genotype \vec{i} as the base s representation of an
 153 integer (see Table 1). From this, we can see how the first s genotypes correspond to combining the wild-type
 154 allele at the first position with a state from the case $n = 1$, i.e. to genotypes that can be written $0 \frown \vec{i} := (0, \vec{i})$,
 155 with $\vec{i} \in \mathcal{S}^1$. The next s genotypes correspond to the first mutated allele at the first position combined with

156 all the genotypes of $n = 1$, i.e. $1 \curvearrowright \vec{i}, \vec{i} \in S^1$, and so on. Therefore, we can write the matrices H and V
 157 following a block structure. In the case $n = 2$ and any given s , we would then have

$$H_2 = \begin{pmatrix} H_1 & H_1 & H_1 & \dots & H_1 \\ H_1 & -H_1 & 0 & \dots & 0 \\ H_1 & 0 & -H_1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ H_1 & 0 & \dots & 0 & -H_1 \end{pmatrix},$$

158 where the number of H_1 blocks corresponds to the number of states of the first position, so s . Moreover,
 159 each of these blocks must be normalized to yield the corresponding background-averaged epistatic terms.
 160 Therefore V_2 can also be expressed as a function of V_1 as follows:

$$V_2 = \begin{pmatrix} \frac{1}{s}V_1 & 0 & \dots & 0 \\ 0 & -V_1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -V_1 \end{pmatrix}.$$

161 Given these two matrices, we can write the background-averaged epistatic coefficients for the case of $n = 2$
 162 and s different states per position as $\bar{\epsilon}_2 = V_2 \cdot H_2 \cdot \bar{y}_2$. More generally, the decomposition of phenotypic
 163 effects into background-averaged epistatic coefficients is given by

$$\bar{\epsilon}_n = V_n \cdot H_n \cdot \bar{y}_n, \quad (1)$$

164 where H_n and V_n can be defined recursively as

$$H_{n+1} = \begin{pmatrix} H_n & H_n & H_n & \dots & H_n \\ H_n & -H_n & 0 & \dots & 0 \\ H_n & 0 & -H_n & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ H_n & 0 & \dots & 0 & -H_n \end{pmatrix} \quad H_0 = 1 \quad \text{and} \quad H_n \text{ is } s^n \times s^n, \quad (2)$$

165

$$V_{n+1} = \begin{pmatrix} \frac{1}{s}V_n & 0 & 0 & \dots & 0 \\ 0 & -V_n & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & \dots & 0 & -V_n \end{pmatrix} \quad V_0 = 1 \quad \text{and} \quad V_n \text{ is } s^n \times s^n. \quad (3)$$

166 Recursive inverse matrix

167 Equation (1) defines the vector of epistatic coefficients, $\bar{\epsilon}_n$, as a function of the vector of phenotypes, \bar{y}_n , which
 168 in general is the quantity that is measured experimentally. However, usually phenotypic measurements are
 169 only available for a subset of genotypes. An alternative is therefore to estimate the epistatic coefficients $\bar{\epsilon}_n$
 170 by regression,

$$\bar{y}_n = H_n^{-1} \cdot V_n^{-1} \cdot \bar{\epsilon}_n, \quad (4)$$

171 where the product $H_n^{-1} \cdot V_n^{-1}$ represents a matrix of sequence features. This is analogous to the more widely
 172 used one-hot encoding strategy, which implicitly relies on a ‘background-relative’ (or ‘biochemical’) view of
 173 epistasis when regressing to full order [2]. We discuss other advantages of estimating background-averaged
 174 epistatic coefficients using regression at the end of this manuscript.

175 Since V_n is a diagonal matrix, its inverse is also a diagonal matrix whose elements are the inverse of the
 176 elements of V_n .

177 The inverse of H_n is the matrix A_n which can be defined recursively as

$$A_{n+1} = \frac{1}{s} \begin{pmatrix} A_n & A_n & A_n & \dots & A_n \\ A_n & (1-s)A_n & A_n & \dots & A_n \\ A_n & A_n & (1-s)A_n & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & A_n \\ A_n & A_n & \dots & A_n & (1-s)A_n \end{pmatrix} \quad A_0 = 1 \quad \text{and} \quad A_n \text{ is } s^n \times s^n. \quad (5)$$

178 See Proposition 1 in S1 Text for a proof of this result. This is the most efficient method to determine the full
 179 matrix A_n (see Results) and, to the best of our knowledge, the first reported recursive definition of the inverse
 180 Walsh-Hadamard transform.

181 Formulae to obtain elements of the matrices

182 When regressing phenotypes on genotypes, a common goal is to determine whether epistatic coefficients up
 183 to the r^{th} order (where $r < n$) are sufficient to describe the complexity of the biological system. Furthermore,
 184 as mentioned above, some fraction of phenotype values within combinatorially complete genetic landscapes
 185 are typically unavailable, representing missing data. Restricting the epistatic order and missing phenotypes
 186 respectively correspond to omitting rows and columns from H_n (and vice versa from A_n). Formulae to directly
 187 obtain elements of the matrices in equations (1) and (4) would therefore be convenient.

188 In order to write the matrix element $(H_n)_{ij}$, we need to compare the genotype sequences $\vec{i}, \vec{j} \in S^n$,

$$\vec{i} = (i_1, i_2, \dots, i_n)$$

$$\vec{j} = (j_1, j_2, \dots, j_n),$$

189 where \vec{i} denotes the i^{th} element in S^n , $S = \{0, 1, \dots, s-1\}$, and the elements of S^n are ordered by the base
 190 s representation of integers. For instance, for any value of n , we will denote the wild-type state with index
 191 $i = 1$ and write $\vec{i} = \vec{1} = (0, \dots, 0)$. The element denoted with index $i = 2$ would be $\vec{i} = \vec{2} = (0, \dots, 0, 1)$ and
 192 so on.

193 The elements of H_n can be written as

$$(H_n)_{ij} = \begin{cases} (-1)^{(E_n)_{ij}} & \text{if } (M_n)_{ij} = n \\ 0 & \text{otherwise,} \end{cases}$$

194 where M and E are $s^n \times s^n$ matrices whose elements are

$$(E_n)_{ij} = \sum_{\substack{k=1 \\ i_k \cdot j_k > 0}}^n \delta_{i_k j_k} \quad (6)$$

$$(M_n)_{ij} = \sum_{\substack{k=1 \\ i_k \cdot j_k > 0}}^n \delta_{i_k j_k} + \sum_{\substack{k=1 \\ i_k \cdot j_k = 0}}^n 1 = (E_n)_{ij} + \sum_{\substack{k=1 \\ i_k \cdot j_k = 0}}^n 1,$$

195 where δ_{ij} denotes the Kronecker delta of i, j , which is equal to 1 when $i = j$ and 0 if $i \neq j$. In words, $(E_n)_{ij}$
 196 counts the number of positions at which the genotype sequences \vec{i} and \vec{j} carry the same mutated allele and
 197 $(M_n)_{ij}$ is equal to $(E_n)_{ij}$ plus the number of positions where \vec{i} or \vec{j} carry the wild-type allele. See Proposition
 198 2 in S1 Text for a proof of this result.

199 Furthermore, the elements of A_n can be written as

$$(A_n)_{ij} = \frac{1}{s^n} (1 - s)^{(E_n)_{ij}}, \quad (7)$$

200 where E_n is defined as in (6). See Proposition 3 in S1 Text for a proof of this result.

201 Finally, the matrices V_n and V_n^{-1} are diagonal matrices whose diagonal elements can be written as

$$(V_n)_{ii} = (-1)^{n - W_n(\vec{i})} \frac{1}{s^{W_n(\vec{i})}} \quad (8)$$

202 and

$$(V_n^{-1})_{ii} = (-1)^{n - W_n(\vec{i})} s^{W_n(\vec{i})}, \quad (9)$$

where

$$W_n(\vec{i}) := \sum_{k=1}^n w_k, \text{ with } w_k := \delta_{i_k 0}$$

203 and \vec{i} again denotes the i^{th} element in S^n when ordered by the base s representation of integers. In words,
 204 $w_k = 1$ if the genotype sequence \vec{i} carries the wild-type allele at position k and $W_n(\vec{i})$ counts the number of
 205 positions in \vec{i} carrying the wild-type allele. We prove this result in Proposition 4 in S1 Text.

206 Generalization to different numbers of states per position

207 We can generalize the formulae described in the previous subsection further by considering that each position
 208 can have different numbers of states. In this case, we can denote s_k the number of possible states at position
 209 k . For $n = 1$, this corresponds to exactly the same matrix as in the previous case but with $s = s_1$, which is
 210 the number of possible states in this position. For $n = 2$, the matrix changes because now the new position
 211 can have a different number of possible states, s_2 . Following the recursive definition of H_n , we can construct
 212 H_2 by repeating H_1 s_2 times, with the structure stated in (2). Therefore, we have

$$H_2 = \underbrace{\begin{pmatrix} \overbrace{H_1}^{s_1} & H_1 & H_1 & \dots & H_1 \\ H_1 & -H_1 & 0 & \dots & 0 \\ H_1 & 0 & -H_1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ H_1 & 0 & \dots & 0 & -H_1 \end{pmatrix}}_{s_2 \text{ blocks of size } s_1 \implies s_2 s_1}.$$

213 So the structure is exactly the same but the size of the matrix for each n varies according to the number of
 214 possible states of the new position. The definition of H_n is the same as in (2) but the dimensions of the matrix
 215 are $\prod_{k=1}^n s_k \times \prod_{k=1}^n s_k$. Similarly, the inverse matrix A_{n+1} can be written recursively as

$$A_{n+1} = \frac{1}{s_{n+1}} \begin{pmatrix} A_n & A_n & A_n & \dots & A_n \\ A_n & (1 - s_{n+1})A_n & A_n & \dots & A_n \\ A_n & A_n & (1 - s_{n+1})A_n & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & A_n \\ A_n & A_n & \dots & A_n & (1 - s_{n+1})A_n \end{pmatrix}, \quad A_0 = 1 \quad \text{and} \quad A_n \text{ is } \prod_{k=1}^n s_k \times \prod_{k=1}^n s_k. \quad (10)$$

216 The matrix A_n defined in (10) is the inverse of the matrix H_n in the general case where each position can
 217 have a different number of states.

218 In this general case, the elements of H_n and A_n can be written as

$$(H_n)_{ij} = \begin{cases} (-1)^{(E_n)_{ij}} & \text{if } (M_n)_{ij} = n \\ 0 & \text{otherwise} \end{cases}$$

$$(A_n)_{ij} = \frac{\prod_{k=1}^n (1 - s_k)^{e_k}}{\prod_{k=1}^n s_k},$$

219 where E_n and M_n are defined as in (6) and $e_k = \begin{cases} 1 & \text{if } i_k = j_k \neq 1 \\ 0 & \text{otherwise} \end{cases}$.

220 The matrices V_n and V_n^{-1} are diagonal matrices whose diagonal elements can be written as

$$(V_n)_{ii} = (-1)^{n - W_n(\vec{i})} \prod_{k=1}^n \left(\frac{1}{s_k} \right)^{w_k}$$

and

$$(V_n^{-1})_{ii} = (-1)^{n - W_n(\vec{i})} \prod_{k=1}^n s_k^{w_k},$$

where

$$W_n(\vec{i}) := \sum_{k=1}^n w_k, \text{ with } w_k := \delta_{i_k 0}.$$

221 We prove the results in this subsection in Propositions 5, 6 and 7 in S1 Text.

222 The above formulae permit the calculation and modeling of background-averaged epistasis in arbitrarily-
 223 shaped genetic landscapes, i.e. with any number of alleles (states) per position, as well as the direct construc-
 224 tion of sub-matrices for regression to any desired epistatic order and/or in the presence of missing data. In the
 225 following subsections we report benchmarking results comparing the performance of alternative methods to
 226 obtain H_n and A_n , as well as results from the application of our theory extension to an empirical multiallelic
 227 genotype-phenotype landscape.

228 Benchmarking

229 Fig 1a-d provides a visualization of the matrices H_n and A_n for different values of n and s , clearly showing
 230 a fractal pattern in all cases due to their recursive nature.

231 In this paper, we provide different methods to construct $A_n = H_n^{-1}$. First, H_n can be numerically inverted
 232 using standard matrix inversion algorithms (here we use the `linalg.inv` function from the SciPy library
 233 in Python), referred to as “Recursive H_n inverse” in Fig 1e,f. Alternatively, the recursive definition of the
 234 inverse given by equation (5) can be used, which we refer to as “Recursive A_n ”. As can be seen in Fig 1e,
 235 this method is faster than numerically inverting H_n .

236 Finally, we also provide a convenient formula for extracting specific individual elements of A_n (Proposition
 237 3), referred to as “All elements A_n ” in Fig 1e,f. This method is more computationally intensive than the
 238 previously described methods, due to the formula relying on the computation of $(E_n)_{ij}$, which equates to
 239 counting the number of sequence positions that are identically mutated in vectors \vec{i} and \vec{j} , each of size n .
 240 However, in situations where subsets of elements (or sub-matrices) – rather than full matrices – are desired,

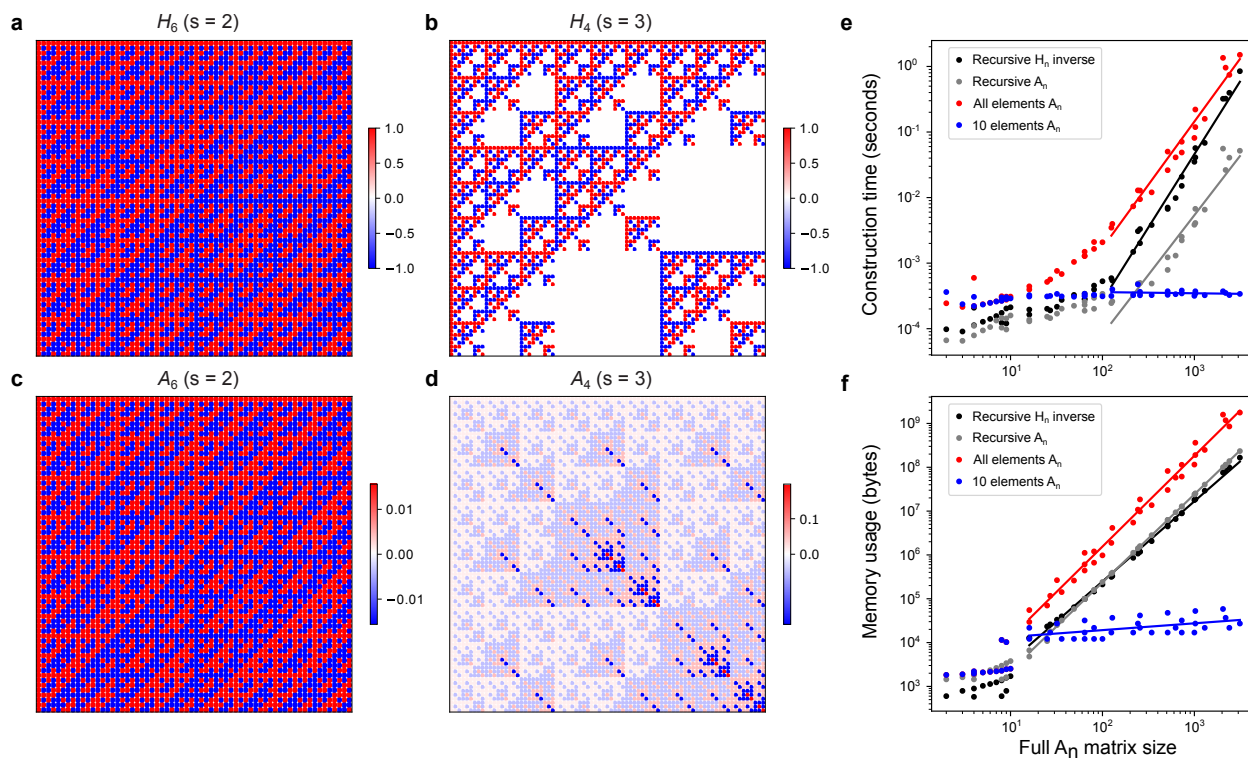


Figure 1: Benchmarking results and heat map representations of matrices corresponding to the binary (biallelic) and multi-state (multiallelic) extension of the Walsh-Hadamard transform, and their corresponding inverses. **a**, H_6 Walsh-Hadamard transform. **b**, H_4 multi-state extension of the Walsh-Hadamard transform for $s = 3$. **c**, A_6 Inverse Walsh-Hadamard transform. **d**, A_4 multi-state extension of the inverse Walsh-Hadamard transform for $s = 3$. **e**, Computational time on a MacBook Pro (13-inch, 2017, 2.3GHz dual-core Intel Core i5) for extracting elements of A_n matrices of various dimensions and numbers of states (alleles) per position ($s \in [2, 10]$). Comparisons are shown between numerically inverting the recursively constructed H_n (using `scipy.linalg.inv`), i.e. “Recursive H_n inverse”, using the recursive formula for A_n , using the formula to extract all elements of A_n and extracting 10 random elements of A_n (see legend). The mean across 10 replicates is depicted. Linear regression lines were fit to data from matrices with at least 100 elements. **f**, Similar to **e** but indicating memory usage. Linear regression lines were fit to data from matrices with at least 10 elements.

241 Proposition 3 provides a method that can be faster and more memory efficient (see “10 elements A_n ” in Fig
242 1e,f).

243 For example, in the case of a 10-mer DNA sequence, constructing the full inverse transform A_{10} with $s = 4$
244 would require $> 10^{23}$ bytes (100 million petabytes) of memory in the best-case scenario (“Recursive H_n
245 inverse” in Fig 1f, log-linear extrapolation). Similarly, the full inverse transform for a 4-mer amino acid se-
246 quence (A_4 with $s = 20$) would impose a memory footprint $> 10^{20}$ bytes. On the other hand, calculating
247 the subset of elements from these matrices required for the prediction of a single phenotype using epistatic
248 coefficients up to third order (three-way genetic interaction terms) is feasible in both situations using Propo-
249 sition 3 (3,675 and 29,678 elements; 2.5 GB and 192 GB of memory; 1.8 and 99 seconds, respectively). This
250 memory footprint can easily be diminished further using data chunking, which is a unique benefit of this
251 method.

252 **Application to a multiallelic genotype-phenotype landscape**

253 In order to demonstrate the utility of our theory, we used it to model epistasis within a combinatorially com-
254 plete multiallelic genetic landscape of a tRNA. Fig 2a-c summarises the model system and DMS experimental
255 strategy employed in [6]. Briefly, a budding yeast strain was used in which the single-copy arginine-CCU
256 tRNA (tRNA-Arg(CCU)) gene is conditionally required for growth. A library of variants of this gene was
257 designed to cover all 5,184 ($2^6 \times 3^4$) combinations of the 14 nucleotide substitutions observed in ten positions
258 in post-whole-genome duplication yeast species (Fig 2a,b). The library was transformed into *S. cerevisiae*,
259 expressed under restrictive conditions and the enrichment of each genotype in the culture was quantified by
260 deep sequencing before and after selection (Fig 2c). After reprocessing of the raw data, we retained high
261 quality fitness estimates for 3,847 variants (74.2%).

262 Although the findings in [6] were based on the application of background-averaged epistasis theory, the prior
263 limitation of the Walsh-Hadamard transform to only two alleles per sequence position required the authors
264 to adopt an *ad hoc* strategy that involved performing separate analyses on combinatorially complete biallelic
265 sub-landscapes.

266 However, with the extensions provided in this work, we were able model background-averaged epistasis in
267 this multiallelic landscape using all available data simultaneously. We trained Lasso regression models of
268 the form in equation (4) to predict variant fitness from nucleotide sequences, where the inferred model pa-
269 rameters correspond to background-averaged epistatic coefficients up to eighth order (Fig 2d; see Methods).
270 To determine the effect of data sparsity on the results, we sub-sampled the original data to obtain training
271 dataset sizes ranging from 64% to 1% of all variants with high quality fitness estimates. The resulting mod-
272 els incorporate many higher-order epistatic coefficients (Fig 2e, ‘Background-averaged models’) yet exhibit
273 extreme sparsity, with the median number of non-zero coefficients of any order ranging from 19 to 60 i.e.
274 approximately 1% of all possible coefficients of eighth order or less (Fig 2f, S1 Fig). Fig 2e indicates model
275 performance on held out test data, with all models except those fit using the most severely subsampled data
276 (1%) tending to explain more than 50% of the total explainable variance.

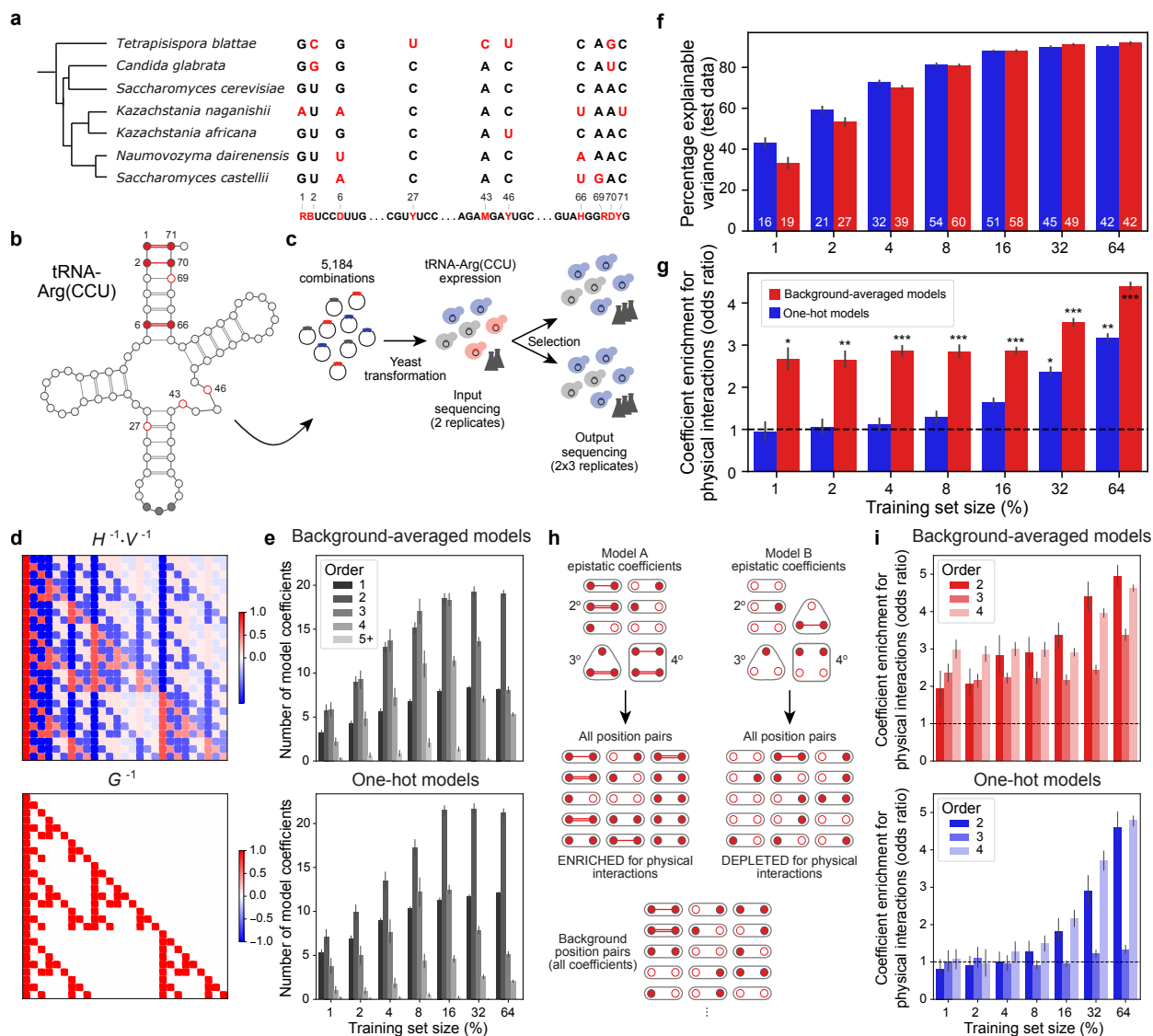


Figure 2: Learning sparse models from the near combinatorially complete fitness landscape of a tRNA. **a**, Species phylogenetic tree and multiple sequence alignment of the tRNA-Arg(CCU) orthologues indicating variable positions across the seven yeast species and the synthesized library below: R (A or G); B (C, G or T); D (A, G or T); Y (C or T); M (A or C); H (A, C or T). **b**, Secondary structure of *S. cerevisiae* tRNA-Arg(CCU) indicating variable positions (open and closed red circles) and three Watson–Crick base pairing (WCBP) interactions between pairs of variable positions i.e. [1,71], [2,70] and [6,66] (red lines and closed red circles). **c**, DMS experiment to quantify the phenotypic effects of all variants in the combinatorially complete genetic landscape. See [6] for details. **d**, Cartoon depiction of alternative feature matrices for inferring epistatic coefficients by linear regression. G^{-1} in the lower panel indicates the matrix of one-hot encoded sequence features – or embeddings – typically used when fitting models of genotype-phenotype landscapes [2]. The upper panel represents the matrix of sequence features used to infer background-averaged epistatic coefficients, as in equation (4). **e**, Numbers of non-zero epistatic coefficients of different orders in Lasso regression models inferred using different random fractions of the DMS data indicated in panels a-c. **f**, Performance of Lasso regression models. The median number of model coefficients is indicated. Colour scale as in panel g. **g**, Enrichment of direct physical interactions (red lines in panel b) in non-zero epistatic coefficients (see panel h). *, $P < 0.05$; **, $P < 0.01$; ***, $P < 0.001$. **h**, Strategy for testing enrichment of direct physical interactions in Lasso regression model coefficients. **i**, Same as panel g, except enrichment test results are shown separately for epistatic coefficients of different orders. Error bars indicate nonparametric bootstrap 95% confidence intervals of the mean in all panels.

277 For comparison, we used the same procedure to fit Lasso regression models of the form $\bar{y} = G^{-1} \cdot \bar{\epsilon}$, where
278 G^{-1} represents a matrix of one-hot encoded sequence features i.e. the presence or absence of a given mutation
279 – or mutation combination (interaction) – with respect to the reference (wild-type) genotype is denoted by a
280 ‘1’ or ‘0’ respectively (Fig 2d-f, ‘One-hot models’). The definition of G and its relationship to the biochemical
281 (or background-relative) view of epistasis is explained in [2]. The sparsity of one-hot models is similar to
282 that of background-averaged models regardless of training set size (Fig 2e,f). However, the latter tend to
283 incorporate greater numbers of higher-order epistatic terms, particularly with larger training set sizes (Fig
284 2e, orders 3,4,5+), whereas the former tend to perform slightly better with very small training set sizes (Fig
285 2f).

286 To evaluate whether the inferred models report on biologically relevant features of the underlying genetic
287 landscapes, we tested whether sparse model coefficients were more likely to comprise genetic interactions (or
288 modulators thereof) involving known physically interacting positions in the wild-type tRNA secondary struc-
289 ture (Fig 2b,h). Regardless of data sparsity, background-averaged model coefficients tend to be significantly
290 enriched for physical interactions (Fig 2g, S1 Fig). On the other hand, in the case of even moderate sub-
291 sampling of training data (16%), one-hot model coefficients show no such enrichment (Fig 2g). Importantly,
292 repeating a similar enrichment analysis using randomly selected model coefficients of identical number and
293 distribution over coefficient orders speaks to the validity of the Fisher’s Exact Test null hypothesis with only
294 minor inflation of the corresponding test statistic (S1 Fig). Restricting the enrichment analysis to epistatic
295 coefficients of specific orders shows qualitatively similar results, with background-averaged model coeffi-
296 cients up to fourth order significantly enriched for physical interacting position pairs, even at the most severe
297 sub-sampling fractions (Fig 2i, S1 Fig).

298 Discussion

299 We have provided an extension to the most rigorous computational framework available for describing and
300 modeling empirical genotype-phenotype mappings. Beyond the study of background-averaged epistasis with
301 respect to mutations in the primary sequence, this also permits the inclusion of ‘epimutations’ (changes in
302 the epigenetic state of DNA), amino acid post-translational modifications or even particular environmen-
303 tal/experimental conditions.

304 In the simplest application, background-averaged epistatic coefficients (genetic interaction terms) can be di-
305 rectly computed from phenotypic measurements via the decomposition in equation (1). However, estimating
306 epistatic coefficients by regression – as in equation (4) – is a more natural choice in the presence of missing
307 data, when data for multiple related phenotypes is available [22] and/or in the presence of global epistasis
308 [23, 24]. Our mathematical results provide three alternative methods to compute the multi-state (multiallelic)
309 extension of the inverse Walsh-Hadamard transform A_n , one of which allows the direct extraction of specific
310 elements or sub-matrices. In which situations might this capability be desirable?

311 First, constructing full A_n matrices – particularly by numerical inversion – is impractical for large genetic
312 landscapes. Second, the result of the product $H_n^{-1} \cdot V_n^{-1}$ represents a matrix of sequence features when

313 setting up the inference of epistatic (model) coefficients $\bar{\epsilon}_n$ from phenotypic measurements \bar{y}_n as a regression
314 task [22, 23, 25, 26]. The ability to construct this feature matrix in batches (of rows) allows computational
315 resource-efficient iteration over large datasets when using frameworks such as TensorFlow or PyTorch.

316 Third, there are currently no methods to comprehensively map empirical genotype-phenotype landscapes
317 with size greater than the low millions of genotypes. Therefore, assaying landscapes of this size or larger
318 will typically involve experimental measurement of a (random) sub-sample of genotypes, corresponding to
319 distinct rows in A_n . In other words, it is usually unnecessary to construct full A_n matrices when modeling real
320 experimental data. Finally, there is evidence of extreme sparsity in the epistatic architecture of biomolecules
321 where only a small fraction of theoretically possible genetic interactions are non-zero [7]. The feasibility of
322 sampling very large background-averaged epistatic coefficient spaces may improve methods to infer accurate
323 genotype-phenotype models.

324 Using results from the analysis of a near combinatorially complete multiallelic fitness landscape of a tRNA,
325 we have shown that sparse regression models relying on a background-averaged definition of epistasis can
326 efficiently capture salient features of the underlying biological system – namely direct physical interactions
327 – even in situations of sparse sampling of phenotypes. This behaviour, which we speculate is due to a richer
328 representation of the sequence feature space compared to one-hot models (i.e. higher level of constraint
329 during model fitting; Fig 2d), is particularly desirable in the case of very large genetic landscapes where
330 comprehensive phenotyping is infeasible. However, more work is needed to determine whether this result
331 holds more generally. One difficulty in such comparisons between approaches is the requirement for a set
332 of interactions or landscape features that are known to be critical for biomolecular function. Here we rely
333 on Watson–Crick base pairing interactions whose importance for RNA secondary structure and function is
334 well-established.

335 More broadly, this work opens the door to investigations of the biological properties of background-averaged
336 epistasis in empirical genetic landscapes of arbitrary shape and complexity. Beyond applications within the
337 field of DMS, we believe our theory extensions have the potential to influence research in evolutionary and
338 synthetic biology including protein engineering. In future it will be important to compare the performance
339 and properties of models relying on this definition of epistasis to those of other recently proposed models
340 that incorporate higher-order genetic interactions for phenotypic prediction [27, 28].

341 **Methods**

342 Raw sequencing (FASTQ) files obtained from the tRNA-Arg(CCU) DMS experiment in [6] were re-processed
343 with DiMSum v1.3 [29] using default parameters with minor adjustments. We obtained fitness estimates for
344 5,059 out of a total of 5,184 possible variants (97.6%) in the combinatorially complete genetic landscape.
345 We restricted the data to a high quality subset by requiring fitness estimates in all six biological replicates as
346 well as at least 10 input read counts in all input samples. This resulted in a total of 3,847 retained variants
347 (74.2%) for downstream analysis.

348 We trained Lasso regression models to predict variant fitness estimates from nucleotide sequences using
349 the ‘scikit-learn’ Python package. Training data comprised random subsets of 1, 2, 4, 8, 16, 32 and 64% of
350 retained variants of all mutation orders. All remaining held out variants comprised the ‘test’ data which was
351 unseen during model training in each case.

352 To train models inferring background-averaged epistatic coefficients we used feature matrices of the form
353 $H_n^{-1} \cdot V_n^{-1}$ (see equation (4)). For comparison, one-hot encoded matrices of sequence features were used.
354 Linear regression was performed using 10-fold cross validation to determine the optimal value of the L1
355 regularization parameter λ in the range [0.005, 0.25] (‘LassoCV’ and ‘RepeatedKFold’ functions). Final
356 models were fit to all training data. In order to estimate model-related statistics and performance results we
357 fit 100 models to different random subsets of the training data for each model type and training data fraction.
358 In Fig 2 and S1 Fig we plot the mean or median of the indicated measures over all models, where 95%
359 confidence intervals were obtained using a nonparametric bootstrap approach and 1000 bootstrap samples.
360 For performance estimates in Fig 2f we estimated the maximum explainable variance by taking the square
361 of the mean Pearson correlation between replicate fitness estimates over all 15 pairwise combinations.

362 To test enrichment of physical interactions in Lasso model coefficients we used the strategy illustrated in Fig
363 2h. For each model, all position pairs represented in non-zero epistatic coefficients of at least second order
364 were determined. The number of position pairs corresponding to direct physical interactions was counted and
365 an associated enrichment score (odds ratio) and P-value calculated using Fisher’s Exact Test. The background
366 set consisted of all position pairs in all possible epistatic coefficients. To test the appropriateness of the
367 null hypothesis we also repeated enrichment analyses using random models i.e. randomly chosen sets of
368 epistatic coefficients matching the numbers of non-zero coefficients in Lasso models and their distribution
369 over different epistatic orders.

370 **Acknowledgments**

371 We thank all members of the Lehner and Weghorn Labs for helpful discussions and suggestions.

372 **References**

- 373 [1] Phillips PC. Epistasis—the essential role of gene interactions in the structure and evolution of genetic
374 systems. *Nature Reviews Genetics*. 2008;9(11):855-67.
- 375 [2] Poelwijk FJ, Krishna V, Ranganathan R. The Context-Dependence of Mutations: A Linkage of For-
376 malisms. *PLoS Computational Biology*. 2016 Jun;12(6):e1004771.
- 377 [3] Domingo J, Baeza-Centurion P, Lehner B. The Causes and Consequences of Genetic Interactions
378 (Epistasis). *Annu Rev Genomics Hum Genet*. 2019 Aug;20:433-60.

- 379 [4] Fowler DM, Fields S. Deep mutational scanning: a new style of protein science. *Nat Methods*. 2014
380 Aug;11(8):801-7.
- 381 [5] de Visser JAG, Krug J. Empirical fitness landscapes and the predictability of evolution. *Nature Reviews*
382 *Genetics*. 2014;15(7):480-90.
- 383 [6] Domingo J, Diss G, Lehner B. Pairwise and higher-order genetic interactions during the evolution of a
384 tRNA. *Nature*. 2018 Jun;558(7708):117-21.
- 385 [7] Poelwijk FJ, Socolich M, Ranganathan R. Learning the pattern of epistasis linking genotype and phe-
386 notype in a protein. *Nature communications*. 2019 Sep;10(1):1-11.
- 387 [8] Baeza-Centurion P, Miñana B, Schmiedel JM, Valcárcel J, Lehner B. Combinatorial genetics reveals a
388 scaling law for the effects of mutations on splicing. *Cell*. 2019;176(3):549-63.
- 389 [9] Pokusaeva VO, Usmanova DR, Putintseva EV, Espinar L, Sarkisyan KS, Mishin AS, et al. An experi-
390 mental assay of the interactions of amino acids from orthologous sequences shaping a complex fitness
391 landscape. *PLoS genetics*. 2019;15(4):e1008079.
- 392 [10] Bendixsen DP, Collet J, Østman B, Hayden EJ. Genotype network intersections promote evolutionary
393 innovation. *PLoS biology*. 2019;17(5):e3000300.
- 394 [11] Soo VW, Swadling JB, Faure AJ, Warnecke T. Fitness landscape of a dynamic RNA structure. *PLoS*
395 *genetics*. 2021;17(2):e1009353.
- 396 [12] Moulana A, Dupic T, Phillips AM, Chang J, Nieves S, Roffler AA, et al. Compensatory epistasis
397 maintains ACE2 affinity in SARS-CoV-2 Omicron BA. 1. *Nature Communications*. 2022;13(1):1-11.
- 398 [13] Rotrattanadumrong R, Yokobayashi Y. Experimental exploration of a ribozyme neutral network using
399 evolutionary algorithm and deep learning. *Nature communications*. 2022;13(1):1-14.
- 400 [14] Lynch M, Walsh B, et al. *Genetics and analysis of quantitative traits*. vol. 1. Sinauer Sunderland, MA;
401 1998.
- 402 [15] Goldberg DE. *Genetic Algorithms and Walsh Functions: Part I, A Genetic Introduction*. *Complex*
403 *systems*. 1989;3:129-52.
- 404 [16] Weinreich DM, Lan Y, Wylie CS, Heckendorn RB. Should evolutionary geneticists worry about higher-
405 order epistasis? *Current opinion in genetics & development*. 2013;23(6):700-7.
- 406 [17] Poelwijk FJ, Ranganathan R. The relation between alignment covariance and background-averaged
407 epistasis. *arXiv*. 2017;10.48550/ARXIV.1703.10996.
- 408 [18] Brookes DH, Aghazadeh A, Listgarten J. On the sparsity of fitness functions and implications for
409 learning. *Proc Natl Acad Sci U S A*. 2022 Jan;119(1).
- 410 [19] Ogbunugafor CB. The mutation effect reaction norm (μ -rn) highlights environmentally dependent
411 mutation effects and epistatic interactions. *Evolution*. 2022 02;76(s1):37-48.

- 412 [20] Beer T. Walsh transforms. *American Journal of Physics*. 1981;49(5):466-72.
- 413 [21] Stoffer DS. Walsh-Fourier Analysis and its Statistical Applications. *Journal of the American Statistical*
414 *Association*. 1991;86(414):461-79.
- 415 [22] Faure AJ, Domingo J, Schmiedel JM, Hidalgo-Carcedo C, Diss G, Lehner B. Mapping the energetic
416 and allosteric landscapes of protein binding domains. *Nature*. 2022;604(7904):175-83.
- 417 [23] Tareen A, Kooshkbaghi M, Posfai A, Ireland WT, McCandlish DM, Kinney JB. MAVE-NN: learning
418 genotype-phenotype maps from multiplex assays of variant effect. *Genome biology*. 2022;23(1):1-27.
- 419 [24] Otwinowski J, McCandlish DM, Plotkin JB. Inferring the shape of global epistasis. *Proceedings of the*
420 *National Academy of Sciences*. 2018;115(32):E7550-8.
- 421 [25] Forcier TL, Ayaz A, Gill MS, Jones D, Phillips R, Kinney JB. Measuring cis-regulatory energetics in
422 living cells using allelic manifolds. *Elife*. 2018;7:e40618.
- 423 [26] Kinney JB, Murugan A, Callan Jr CG, Cox EC. Using deep sequencing to characterize the biophysical
424 mechanism of a transcriptional regulatory sequence. *Proceedings of the National Academy of Sciences*.
425 2010;107(20):9158-63.
- 426 [27] Zhou J, Wong MS, Chen WC, Krainer AR, Kinney JB, McCandlish DM. Higher-order epistasis and
427 phenotypic prediction. *Proceedings of the National Academy of Sciences*. 2022;119(39):e2204233119.
- 428 [28] Zhou J, McCandlish DM. Minimum epistasis interpolation for sequence-function relationships. *Nature*
429 *communications*. 2020;11(1):1-14.
- 430 [29] Faure AJ, Schmiedel JM, Baeza-Centurion P, Lehner B. DiMSum: an error model and pipeline for
431 analyzing deep mutational scanning data and diagnosing common experimental pathologies. *Genome*
432 *Biology*. 2020;21(1):1-23.

433 **Supporting information captions**

434 **S1 Fig. Supplementary figure related to Fig 2.**

435 **S1 Text. Supplementary Methods.**

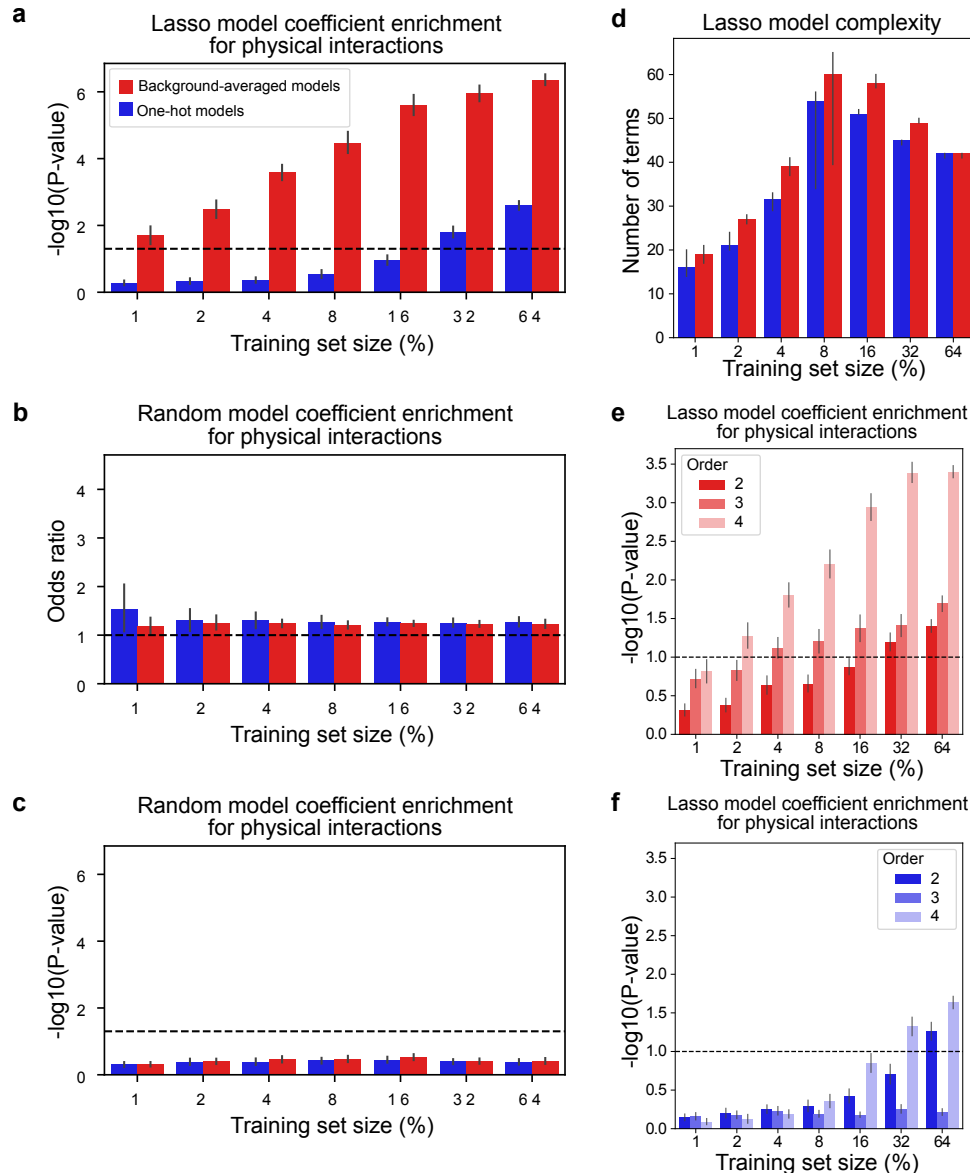


Figure S1: Supplementary figure related to Fig 2. **a**, P-value from Fisher's Exact Test for enrichment of direct physical interactions in non-zero epistatic coefficients (related to Fig 2g). **b**, Enrichment of direct physical interactions in non-zero epistatic coefficients of random models with matching numbers of epistatic coefficients of different orders (related to Fig 2g). P-value from Fisher's Exact Test for enrichment of direct physical interactions in non-zero epistatic coefficients of random models with matching numbers of epistatic coefficients of different order (related to panel b). **d**, Median number of epistatic terms in Lasso models (related to Fig 2f). **f**, P-value from Fisher's Exact Test for enrichment of direct physical interactions in non-zero epistatic shown separately for epistatic coefficients of different orders (related to Fig 2i). Error bars indicate nonparametric bootstrap 95% confidence intervals of the mean in all panels, except in panel d where these correspond to the median.

436 **Supporting Information 1 - Supplementary Methods**

437 **Here we provide the proofs of the mathematical results shown in the main text.**

438 **Proposition 1.** *Let us define the matrices A_n recursively as*

$$A_{n+1} = \frac{1}{s} \begin{pmatrix} A_n & A_n & A_n & \dots & A_n \\ A_n & (1-s)A_n & A_n & \dots & A_n \\ A_n & A_n & (1-s)A_n & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & A_n \\ A_n & A_n & \dots & A_n & (1-s)A_n \end{pmatrix} \quad A_0 = 1 \quad \text{and} \quad A_n \text{ is } s^n \times s^n. \quad (5)$$

439 *For $n \in \mathbb{N}$, A_n is the inverse of the matrix H_n defined in equation (2).*

440 *Proof.* Let us prove this by induction. For $n = 1$ we have

$$H_1 = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & -1 & 0 & \dots & 0 \\ 1 & 0 & -1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 1 & 0 & \dots & 0 & -1 \end{pmatrix} \quad (11)$$

441

$$A_1 = \frac{1}{s} \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 1-s & 1 & \dots & 1 \\ 1 & 1 & 1-s & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 1 \\ 1 & 1 & \dots & 1 & 1-s \end{pmatrix}. \quad (12)$$

442 The rows and columns of these two matrices can be described as follows:

$$(H_1)_i = \begin{cases} (1, 1, \dots, 1) & \text{if } i = 1 \\ (1, 0, \dots, 0, h_i, 0, \dots, 0) & \text{if } i \neq 1, \end{cases}$$

$$(A_1)_{.j} = \frac{1}{s} \begin{cases} (1, 1, \dots, 1)^T & \text{if } j = 1 \\ (1, 1, \dots, 1, a_j, 1, \dots, 1)^T & \text{if } j \neq 1, \end{cases}$$

443 where $h_i := (H_1)_{ii} = -1 \forall i > 1$ and $a_j := (A_1)_{jj} = 1 - s \forall j > 1$.

Therefore,

$$(H_1 \cdot A_1)_{ij} = (H_1)_{i \cdot} \cdot (A_1)_{\cdot j} = \frac{1}{s} \begin{cases} s & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

$$H_1 \cdot A_1 = \frac{1}{s} \begin{pmatrix} s & 0 & \dots & 0 \\ 0 & s & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & s \end{pmatrix} = I_{s \times s},$$

444 where $I_{s \times s}$ is the identity matrix of size $s \times s$. Since both H_1 and A_1 are symmetric, it is also true that
 445 $A_1 \cdot H_1 = I_{s \times s}$. Therefore, A_1 is the inverse of H_1 .

446 Assume that the hypothesis is true for a fixed $n \in \mathbb{N}$. Let us now prove that it is also true for $n + 1$. Following
 447 the recursive definitions of H_{n+1} and A_{n+1} in equations (2) and (5), we can write the blocks of these matrices
 448 as follows:

$$(H_{n+1})_{[i][\cdot]} = \begin{cases} (H_n, H_n, \dots, H_n) & \text{if } i = 1 \\ (H_n, 0, \dots, 0, \tilde{h}_i, 0, \dots, 0) & \text{if } i \neq 1, \end{cases} \quad (13)$$

449 where $(H_{n+1})_{[i][j]}$ denotes the block at position i, j in H_{n+1} and $\tilde{h}_i := (H_{n+1})_{[i][i]} = -H_n \forall i > 1$;

$$(A_{n+1})_{[\cdot][j]} = \frac{1}{s} \begin{cases} (A_n, A_n, \dots, A_n)^T & \text{if } j = 1 \\ (A_n, A_n, \dots, A_n, \tilde{a}_j, A_n, \dots, A_n)^T & \text{if } j \neq 1, \end{cases} \quad (14)$$

where $(A_{n+1})_{[i][j]}$ denotes the block at position i, j in A_{n+1} and $\tilde{a}_j := s(A_{n+1})_{[j][j]} = (1 - s)A_n \forall j > 1$. We can therefore write the block at position i, j of the product of these matrices as follows:

$$(H_{n+1} \cdot A_{n+1})_{[i][j]} = (H_{n+1})_{[i][\cdot]} \cdot (A_{n+1})_{[\cdot][j]} = \frac{1}{s} \begin{cases} sH_n \cdot A_n & \text{if } i = j \\ 0 & \text{if } i \neq j. \end{cases}$$

450 According to the induction hypothesis we know that H_n and A_n are inverse matrices i.e. $H_n \cdot A_n = I_{s^n \times s^n}$.
 451 Therefore, the blocks on the diagonal are identity matrices and the blocks outside the diagonal are zeros.
 452 This means that $H_{n+1} \cdot A_{n+1} = I_{s^{n+1} \times s^{n+1}}$. Similarly, due to the symmetry of the matrices we can also prove
 453 that $A_{n+1} \cdot H_{n+1} = I_{s^{n+1} \times s^{n+1}}$.

454 We can then conclude that $A_n = H_n^{-1}$ for every $n \in \mathbb{N}$. □

455 **Proposition 2.** *The elements of H_n can be written as*

$$(H_n)_{ij} = \begin{cases} (-1)^{(E_n)_{ij}} & \text{if } (M_n)_{ij} = n \\ 0 & \text{otherwise,} \end{cases}$$

456 where M and E are $s^n \times s^n$ matrices whose elements are

$$(E_n)_{ij} = \sum_{\substack{k=1 \\ i_k \cdot j_k > 0}}^n \delta_{i_k j_k} \quad (6)$$

$$(M_n)_{ij} = \sum_{\substack{k=1 \\ i_k \cdot j_k > 0}}^n \delta_{i_k j_k} + \sum_{\substack{k=1 \\ i_k \cdot j_k = 0}}^n 1 = (E_n)_{ij} + \sum_{\substack{k=1 \\ i_k \cdot j_k = 0}}^n 1,$$

457 where δ_{ij} denotes the Kronecker delta of i, j .

458 *Proof.* Let us prove the formula by induction. For $n = 1$ and any given s , H_n is given by equation (11).

459 Therefore, we can write $(E_1)_{ij}$ and $(M_1)_{ij}$ as follows:

$$(M_1)_{ij} = \begin{cases} 1 & \text{if } i = 1 \text{ or } j = 1 \\ 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \neq 1 \end{cases}$$

$$(E_1)_{ij} = \begin{cases} 1 & \text{if } i = j \neq 1 \\ 0 & \text{otherwise.} \end{cases} \quad (15)$$

460 Therefore, since $n = 1$, $(M_1)_{ij} = n = 1$ only when either $i = 1, j = 1$ or $i = j$. In the rest of the cases
 461 $(M_1)_{ij} \neq n = 1$ and, according to the formula, the elements of the matrix will be 0. Now, for the cases where
 462 $(M_1)_{ij} = n = 1$, we need to check the value of $(E_1)_{ij}$. We can see how $(E_1)_{ij} = 1$ only when $i = j$ and
 463 they are different from 1. This means that all the elements of the diagonal of H_1 , except the first one, will be
 464 $(-1)^1 = -1$ and the first row and first columns will have $(-1)^0 = 1$. The rest of the elements correspond to
 465 $(M_1)_{ij} \neq n = 1$ so they will be filled with zeros. Putting all this together, we find the expression as H_1 from
 466 equation (11).

467 Assume now that the expression is true for a fixed $n \in \mathbb{N}$ and let us prove it for $n + 1$. In this case, the matrix
 468 of H_{n+1} is defined by blocks (see equation (2)). We first define the indices $P \in \{1, \dots, s\}$ and $Q \in \{1, \dots, s\}$
 469 for row and column blocks, respectively. The first matrix block of H_{n+1} corresponds to $P = 1$ and $Q = 1$. The
 470 corresponding blocks of the matrices M_{n+1} and E_{n+1} , which are necessary for the derivation of the block in
 471 H_{n+1} , are computed by comparing the genotype sequences $\vec{p}, \vec{q} \in S^{n+1}$ for which $\vec{p} = (0, i_1, \dots, i_n) := 0 \curvearrowright \vec{i}$

472 for $\vec{i} \in S^n$ and $\vec{q} = 0 \frown \vec{j}, \vec{j} \in S^n$. More generally, the block in the P^{th} position with respect to the rows and
 473 Q^{th} position with respect to the columns can be obtained by comparing the genotype sequences $\vec{p}, \vec{q} \in S^{n+1}$
 474 for which $\vec{p} = (P - 1) \frown \vec{i}, \vec{i} \in S^n$ and $\vec{q} = (Q - 1) \frown \vec{j}, \vec{j} \in S^n$. See below for a visual description of the
 475 notation.

476 From these observations, it can easily be deduced that for any $\vec{p} = (P - 1) \frown \vec{i}, \vec{q} = (Q - 1) \frown \vec{j}$ with
 477 $\vec{i}, \vec{j} \in S^n, P - 1, Q - 1 \in S$,

$$(M_{n+1})_{pq} = \begin{cases} (M_n)_{ij} + 1 & \text{if } P = 1 \text{ or } Q = 1 \\ (M_n)_{ij} + 1 & \text{if } P = Q \\ (M_n)_{ij} & \text{if } P \neq Q \neq 1 \end{cases} \quad (16)$$

$$(E_{n+1})_{pq} = \begin{cases} (E_n)_{ij} + 1 & \text{if } P = Q \neq 1 \\ (E_n)_{ij} & \text{otherwise,} \end{cases} \quad (17)$$

478 where $i = p - s^n(P - 1), j = q - s^n(Q - 1), P = \lceil p/s^n \rceil$ and $Q = \lceil q/s^n \rceil$, with $\lceil \cdot \rceil$ denoting the ceiling
 479 function. A visual description of the notation of the block structure of the matrices is given by

$$X_{n+1} = \text{row block } P \left\{ \begin{array}{cccc} X_n^{11} & \dots & X_n^{1Q} & \dots & X_n^{1s} \\ \vdots & \ddots & \vdots & & \vdots \\ X_n^{P1} & \dots & X_n^{PQ} & \dots & X_n^{Ps} \\ \vdots & & \vdots & \ddots & \vdots \\ X_n^{s1} & \dots & X_n^{sQ} & \dots & X_n^{ss} \end{array} \right\} \leftarrow \text{row } p$$

column block Q
↑
column q

480 Here X_{n+1} denotes any generic matrix following the structure of the matrices in equations (2), (3), (5), (16)
 481 and (17).

482 Now, similar to the case $n = 1$, we have that $(M_{n+1})_{pq} = n + 1$ only when $(M_n)_{ij} = n$ and either $P = 1$,
 483 $Q = 1$ or $P = Q$, which corresponds to the newly added state being $p_1 = P - 1 = 0, q_1 = Q - 1 = 0$ or
 484 $p_1 = q_1$. In the rest of the cases $(M_{n+1})_{pq} \neq n + 1$ and the elements of the matrix H_{n+1} will be 0. Now, for
 485 the cases where $(M_{n+1})_{pq} = n + 1$, we will have that $(E_{n+1})_{pq}$ has either the same value of the corresponding
 486 entry in E_n or it will be increased by 1. This means, that when $P = Q \neq 1$, the sign of the entry in H_{n+1} will
 487 be inverted. Otherwise, the sign of the element of the matrix stays the same. With this we prove the formula
 488 for H_{n+1} , since we have shown how we can find the same block structure as in equation (2). By induction,
 489 we can conclude that the formula holds for every $n \in \mathbb{N}$. □

490 **Proposition 3.** *The elements of A_n can be written as*

$$(A_n)_{ij} = \frac{1}{s^n} (1-s)^{(E_n)_{ij}}, \quad (7)$$

491 where E_n is defined as in equation (6).

492 *Proof.* Let us prove the formula by induction.

493 For any given s , A_1 is defined in equation 12. Its diagonal elements are equal to $(1-s)/s$ for $i > 1$, $1/s$ for
 494 $i = 1$ and its off-diagonal elements are equal to $1/s$. It can easily be observed from equation (15) that the
 495 RHS of equation (7) is equal to $(1-s)/s$ for $i = j \neq 1$ and to $1/s$ otherwise, so equation (7) is true for $n = 1$.

496 Assume now that equation (7) is true for a fixed $n \in \mathbb{N}$ and let us prove it for $n + 1$. We use the recursive
 497 definition of A_{n+1} in equation (5) and the block representation of the genotype sequences as in the proof of
 498 (2).

499 Let us start with the first block in the diagonal of A_{n+1} , i.e. $P = 1$ and $Q = 1$, where the entries of the
 500 corresponding block in E_{n+1} are derived from comparisons of pairs of genotype sequences of the form $\vec{p} =$
 501 $0 \frown \vec{i}, \vec{q} = 0 \frown \vec{j}$ for $\vec{i}, \vec{j} \in S^n$. From equation (5), this block is equal to $\frac{1}{s} A_n$, so writing $i = p \bmod s^n$ and
 502 $j = q \bmod s^n$, we have

$$(A_{n+1})_{pq} = \frac{1}{s} (A_n)_{ij} = \frac{1}{s^{n+1}} (1-s)^{(E_n)_{ij}}.$$

503 From equation (17), $(E_{n+1})_{pq} = (E_n)_{ij}$, which yields the desired result.

504 Now let us consider the elements in the other diagonal blocks of A_{n+1} , where the entries correspond to pairs
 505 of genotype sequences of the form $\vec{p} = (P-1) \frown \vec{i}, \vec{q} = (Q-1) \frown \vec{j}$ with $P = Q \neq 1$ and $\vec{i}, \vec{j} \in S^n$. From
 506 equation (5), this block is equal to $\frac{1}{s} (1-s) A_n$, i.e.

$$(A_{n+1})_{pq} = \frac{1}{s} (1-s) (A_n)_{ij} = \frac{1}{s^{n+1}} (1-s)^{(E_n)_{ij}+1}.$$

507 From equation (17), $(E_{n+1})_{pq} = (E_n)_{ij} + 1$, which yields the desired result.

508 Finally, let us consider the elements in the off-diagonal blocks of A_{n+1} , where the entries correspond to pairs
 509 of genotype sequences of the form $\vec{p} = (P-1) \frown \vec{i}, \vec{q} = (Q-1) \frown \vec{j}$ with $P \neq Q$ and $\vec{i}, \vec{j} \in S^n$. From
 510 equation (5), this block is equal to $\frac{1}{s} A_n$, so we have

$$(A_{n+1})_{pq} = \frac{1}{s} (A_n)_{ij} = \frac{1}{s^{n+1}} (1-s)^{(E_n)_{ij}}.$$

511 From equation (17), $(E_{n+1})_{pq} = (E_n)_{ij}$, which completes the proof. \square

512 **Proposition 4.** The matrices V_n and V_n^{-1} are diagonal matrices whose diagonal elements can be written as

$$(V_n)_{ii} = (-1)^{n-W_n(\vec{i})} \frac{1}{s^{W_n(\vec{i})}} \quad (8)$$

513 and

$$(V_n^{-1})_{ii} = (-1)^{n-W_n(\vec{i})} s^{W_n(\vec{i})}, \quad (9)$$

where

$$W_n(\vec{i}) := \sum_{k=1}^n w_k, \text{ with } w_k := \delta_{i_k 0}$$

514 and \vec{i} again denotes the i^{th} element in S^n when ordered by the base s representation of integers.

515 *Proof.* Let us prove equation (8) by induction, equation (9) follows directly.

516 One can easily check from equation (3) that the formula holds for $n = 1$. Let us now assume that equation
 517 (8) is true for a fixed $n \in \mathbb{N}$ and let us prove it for $n + 1$. We use the recursive definition of V_{n+1} in equation
 518 (3). Let us consider the element $(V_{n+1})_{pp}$. If $\vec{p} = 0 \frown \vec{i}$, this corresponds to the first block of V_{n+1} , i.e. $P = 1$,
 519 where the elements are multiplied by $1/s$ and $W_{n+1}(\vec{p}) = W_n(\vec{i}) + 1$, so

$$(V_{n+1})_{pp} = \frac{1}{s}(V_n)_{ii} = (-1)^{n-W_n(\vec{i})} \frac{1}{s^{W_n(\vec{i})+1}} = (-1)^{n+1-W_{n+1}(\vec{p})} \frac{1}{s^{W_{n+1}(\vec{p})}}.$$

520 Similarly, if $\vec{p} = (P - 1) \frown \vec{i}$ and $P > 1$, i.e. for the other diagonal blocks, from the recursive formula
 521 the elements are multiplied by -1 and $W_{n+1}(\vec{p}) = W_n(\vec{i})$, so by writing again $i = p - s^n(P - 1)$, where
 522 $P = \lceil p/s^n \rceil$, we have

$$(V_{n+1})_{pp} = -(V_n)_{ii} = (-1)^{1+n-W_n(\vec{i})} \frac{1}{s^{W_n(\vec{i})}} = (-1)^{n+1-W_{n+1}(\vec{p})} \frac{1}{s^{W_{n+1}(\vec{p})}},$$

523 which completes the proof. □

524 **Proposition 5.** *The matrix A_n defined in equation (10) is the inverse of the matrix H_n in the general case*
 525 *where each position can have a different number of states.*

526 *Proof.* Let us prove by induction that $H_n \cdot A_n = I$ where I is the identity matrix of the corresponding size.
 527 Since H_n and A_n are symmetric, this would imply that $A_n \cdot H_n = I$ as well, and therefore, $A_n = H_n^{-1}$.

528 The case $n = 1$ corresponds exactly to the case $n = 1$ of the proof of Proposition 1 by setting $s = s_1$.
 529 Therefore, $H_1 A_1 = I_{s_1 \times s_1}$, and A_1 is the inverse of H_1 .

530 Now, assume the hypothesis is true for a fixed $n \in \mathbb{N}$ and let us prove that this is also true for $n + 1$. We can
 531 write the rows and columns of the matrices H_{n+1} and A_{n+1} as equation (13) and equation (14), respectively.
 532 The only difference is that we need to replace s by s_{n+1} and the size of the matrices is different. Following
 533 exactly the same derivation as in Proposition 1 we can conclude that $H_{n+1} A_{n+1} = I$ and this proves by
 534 induction that $A_n = H_n^{-1}$. □

535 **Proposition 6.** *In this general case, the elements of H_n and A_n can be written as*

$$(H_n)_{ij} = \begin{cases} (-1)^{(E_n)_{ij}} & \text{if } (M_n)_{ij} = n \\ 0 & \text{otherwise} \end{cases}$$

$$(A_n)_{ij} = \frac{\prod_{k=1}^n (1 - s_k)^{e_k}}{\prod_{k=1}^n s_k},$$

536 where E_n and M_n are defined as in equation (6) and $e_k = \begin{cases} 1 & \text{if } i_k = j_k \neq 1 \\ 0 & \text{otherwise} \end{cases}$.

537 *Proof.* The proof follows directly from the proofs of Propositions 2 and 3. The only difference in the induction
538 step is that s is replaced by s_{n+1} . □

539 **Proposition 7.** *The matrices V_n and V_n^{-1} are diagonal matrices whose diagonal elements can be written as*

$$(V_n)_{ii} = (-1)^{n - W_n(\vec{i})} \prod_{k=1}^n \left(\frac{1}{s_k} \right)^{w_k}$$

and

$$(V_n^{-1})_{ii} = (-1)^{n - W_n(\vec{i})} \prod_{k=1}^n s_k^{w_k},$$

where

$$W_n(\vec{i}) := \sum_{k=1}^n w_k, \text{ with } w_k := \delta_{i_k 0}.$$

540 *Proof.* The proof follows the same steps as the proof of Proposition 4. The only difference in the induction
541 step is that s is replaced by s_{n+1} . □