

# Adaptive Partitioning of the tRNA Interaction Interface by Aminoacyl-tRNA-Synthetases

Andy Collins-Hed<sup>a</sup>, David H. Ardell<sup>a,b</sup>

<sup>a</sup>*Quantitative and Systems Biology Program, 5200 North Lake Road, University of California, Merced, CA, 95306, United States*

<sup>b</sup>*Molecular Cell Biology Unit, 5200 North Lake Road, University of California, Merced, CA, 95306, United States*

---

## Abstract

We introduce rugged fitness landscapes called match landscapes for the co-evolution of feature-based assortative interactions between  $P \geq 2$  cognate pairs of tRNAs and aminoacyl-tRNA synthetases (aaRSs) in aaRS-tRNA interaction networks. Our genotype-phenotype-fitness maps assume additive feature-matching energies, a macroscopic theory of aminoacylation kinetics including proofreading, and selection for translational accuracy in multiple, perfectly encoded site-types. We compute the stationary genotype distributions of finite panmictic, asexual populations of haploid aaRS-tRNA interaction networks evolving under mutation, genetic drift, and selection for cognate matching and non-cognate mismatching of aaRS-tRNA pairs. We compared expected genotype frequencies under different matching rules and fitness functions, both with and without linked site-specific modifiers of interaction. Under selection for translational accuracy alone, our model predicts no selection on modifiers to eliminate non-cognate interactions, so long as they are compensated by tighter cognate interactions. Only under combined selection for both translational accuracy and rate do modifiers adaptively eliminate cross-matching in non-cognate aaRS/tRNA pairs. We theorize that the encoding of macromolecular interaction networks is a genetic language that symbolically maps identifying structural and dynamic features of genes and gene-products to functions within cells. Our theory helps explain 1) the remarkable divergence in how aaRSs bind tRNAs, 2) why interaction-

---

*Email address:* [dardell@ucmerced.edu](mailto:dardell@ucmerced.edu) (David H. Ardell)

*URL:* <http://davidardell.org> (David H. Ardell)

informative features are phylogenetically informative, 3) why the Statistical Tree of Life became more tree-like after the Darwinian Transition, and 4) an approach towards computing the probability of the random origin of an interaction network.

*Keywords:* Rugged Landscapes, Darwinian Transition, Reciprocal Sign Epistasis, Modifier Model, Hamming Code, Karlin-Altschul Theory

---

## 1 Introduction

2 Carl Woese and his co-authors argue influentially that all Earth’s cells and  
3 organelles descend not from one universal ancestor cell, but rather a com-  
4 munal ancestral genetic code — the one operating in ribosomal protein  
5 synthesis [1–5]. Woese’s theory is that our ancestral genetic code evolved  
6 collectively in a community of cells that exchanged genes more frequently  
7 and translated them more ambiguously than we imagine most living cells  
8 would tolerate today (although increasing the accuracy of protein synthe-  
9 sis can be costly, for example in bacteria competing to grow [6–9]). Our  
10 ancestral genetic code evolved as an innovation-sharing protocol [4] in a  
11 “winner-takes-all” or big bang process [10] analogous to systems competi-  
12 tion in economics [11]. That is, the ancestral community of cells converged  
13 on one genetic code in parallel to exploit a convergently encoded pool of  
14 genes that they shared. Once enough genes came to depend on this code,  
15 and cellular fitness increasingly depended on interdependent coordination of  
16 the action of many gene products, an evolutionary phase transition occurred  
17 that “froze” the genetic code [12, 13]. In parallel, increasingly complex fitness  
18 interactions among genes, called generally *epistasis* [14, 15], cooled the rate  
19 of gene sharing, changing the evolution of cells from a genetically commu-  
20 nal to a more vertical mode of inheritance in a Statistical Tree of Life [16],  
21 in what Woese called the *Darwinian Transition*. Broadly consistent with  
22 this theory, it was found that complexity of gene interactions (the number  
23 of pairwise interactions a gene undertakes) constrains “informational” genes  
24 from transferring horizontally between cells relative to condition-dependent  
25 “operational” genes [17, 18] and increasing pairwise protein-protein interac-  
26 tions, as measured in yeast two-hybrid data, reduces substitution rates in  
27 proteins [19].

28 In protein biosynthesis, the translation of sense codons depends directly  
29 on the identity and distribution of amino acids attached or *aminoacylated*

30 to the 3' ends of tRNAs at their *acceptor stems*. Aminoacylation of amino  
31 acids to tRNA acceptor stems is catalyzed in an ATP-dependent two-step  
32 reaction [20] by amino-acid-specific catalytic core domains in tRNA-binding  
33 proteins called aminoacyl-tRNA synthetases (aaRSs) [21–23]. The conserved  
34 and modular domain structure of aaRSs and the ability of some aaRSs to  
35 specifically aminoacylate model acceptor stem hairpins led to the proposal  
36 that aaRS-tRNA interactions evolved through a primordial stage of an “oper-  
37 ational RNA code” depending on a small number of base-pairs in the acceptor  
38 stem [24, 25].

39 However, it is unclear how this theory fully accounts for diversity in  
40 tRNA-binding by aaRSs. As shown in Figure 1, aaRSs exhibit remarkable  
41 diversity in how they bind and interact with tRNAs. AaRSs come in two  
42 conserved and ancient superfamilies called Class I and Class II, with distinct  
43 folds, distinct mechanistic details of catalysis and — critical for our argument  
44 — distinct modes of binding to tRNAs, through opposing major or minor  
45 grooves of tRNA acceptor stems [26]. The two superfamilies may further be  
46 divided each into three subclasses [27], which pre-date the divergence of bac-  
47 teria, archaea and eukarotes [23] as exemplified by the consistency with which  
48 aaRSs can be used to root the statistical Tree of Life [28]. Striking examples  
49 of aaRS pairs of different classes were found that could be docked simulta-  
50 neously on tRNAs [29], which led to the hypothesis that aaRSs may have  
51 originally bound tRNAs in paired complexes to help protect tRNA acceptor  
52 stems, and subsequently diverged to single aaRS-binding with expansion of  
53 the code [30].

54 Because all tRNAs conform to a universal structure, tRNAs must distin-  
55 guish themselves to specific aaRSs through interaction-determining features  
56 called tRNA identity elements, which vary over the major domains of life [34].  
57 We say that the functional identity of a tRNA determines its assortative in-  
58 teraction with proteins as mediated by mutually compatible structural and  
59 dynamical features. Earlier, we applied an information theoretical approach  
60 to predict tRNA identity elements [35]; we call the features we predict *tRNA*  
61 *Class-Informative Features (CIFs)* (they could perhaps more specifically be  
62 called *Interaction-Informative Features (IIFs)*). Through comparative anal-  
63 ysis of tRNA CIFs and also through our tRNA functional classifier [36], we  
64 have shown that tRNA CIFs are variable and phylogenetically informative  
65 within the major domains of life [37–39].

66 There is a widely perceived need for genetically explicit models to inves-  
67 tigate theories about the origin and evolution of the aaRS-tRNA interaction

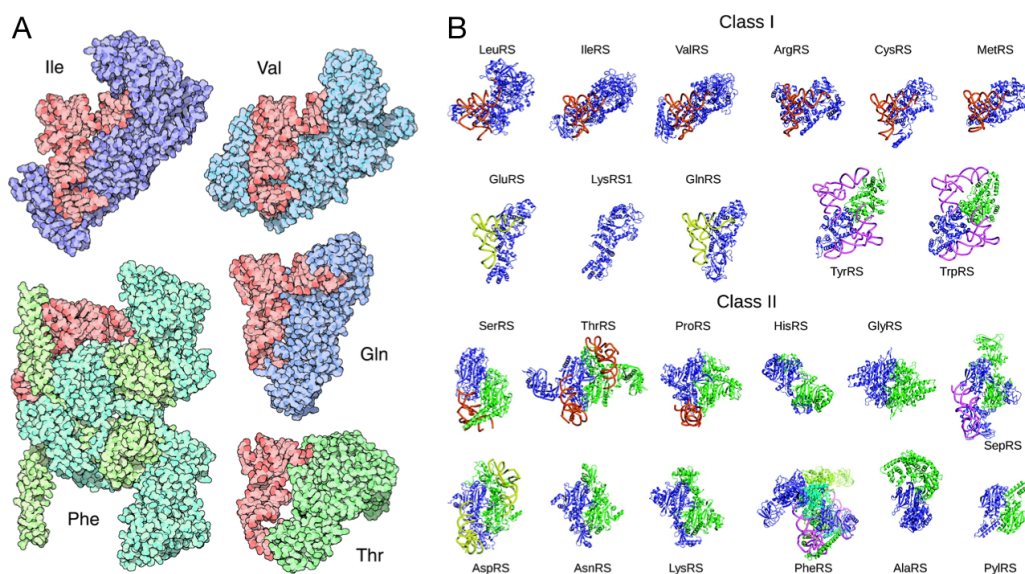


Figure 1: Diversity in tRNA-binding by Class I and Class II aaRSs and within aaRS subclasses. Panel A, reproduced from [31] without modification under License CC-BY 4.0.: Shown in red are different species of tRNAs, all oriented identically. Aminoacyl-tRNA synthetases are shown in purple and green. Class I aaRSs, such as IleRS, ValRS and GlnRS, and class II aaRSs, such as PheRS and ThrRS, bind tRNAs on opposite faces and catalyze aminoacylation on different carbons of the last tRNA base, A76 (in Sprinzl standard coordinates [32]). Panel B: Gallery of aaRS structures co-complexed with tRNAs when available, reprinted (adapted) with permission from [33]. Copyright (2008) American Chemical Society. Subclasses a, b and c of both Class I and Class II aaRS superfamilies are indicated by orange, yellow and pink tRNA colors respectively. aaRSs are visualized with their catalytic domains in the same orientation. (two-column figure)

68 network. For example, Vetsigian et al. [4] showed that horizontal gene trans-  
69 fer of protein-coding genes across a structured population of evolving codes  
70 improves the error-minimizing optimality of genetic codes, but they were  
71 unable to model the effect of horizontal transfer of components of the trans-  
72 lational apparatus itself. They write, “a fuller account of the evolution of  
73 the genetic code requires modeling physical components of the translational  
74 apparatus, including the dynamics of tRNAs and the aminoacyl-tRNA syn-  
75 thetases.” Similarly, Koonin and Novozhilov [40] write, “A real understanding  
76 of the code origin and evolution is likely to be attainable only in conjunction  
77 with a credible scenario for the evolution of the coding principle itself and  
78 the translation system.” Having code evolution models with explicit evolu-

79 tionary dynamics for tRNA and aaRS genes would help test other hypotheses  
80 including the roles of duplication and divergence of tRNA and aaRS genes in  
81 codon assignments [41], and even the dynamics of antisense-encoded aaRSs  
82 according to the Rodin-Ohno hypothesis [42–44].

83 In this work we introduce a theory for feature-based encoding of aaRS-  
84 tRNA interactions that helps answer the following questions:

- 85 1. Why are interaction-informative features phylogenetically informative?
- 86 2. How do interaction-determining features evolve and diverge while still  
87 strongly selected for function and fitness?
- 88 3. Why did more than one superfamily of aaRSs evolve with such different  
89 modes of binding tRNAs? Why is there such diversity in aaRS-binding  
90 of tRNAs even within subclasses (Fig. 1)?
- 91 4. What caused the Darwinian Transition to a more tree-like Statistical  
92 Tree of Life?
- 93 5. What is the probability of random origin of an aaRS-tRNA network of  
94 a given size?

95 At the outset, we considered that divergence in tRNA-binding by aaRSs  
96 could provide increased *robustness* [45] in translational accuracy to muta-  
97 tions in tRNAs and aaRSs (*i.e.* “survival of the flattest” [46]), or potentially  
98 could have evolved to increase the *evolvability* of new aaRS-tRNA interac-  
99 tions. Yet in the results we report here, we show that neither evolutionary  
100 robustness nor increased evolvability is necessary to positively select for di-  
101 vergence in tRNA-binding by aaRSs. Furthermore, selection on translational  
102 accuracy alone was insufficient to select for divergence in tRNA-binding. We  
103 found that combined selection on both accuracy and rate was necessary and  
104 sufficient for aaRS genes to evolve to adaptively partition the tRNA interac-  
105 tion interface. Our results depend on assumptions and modeling concepts as  
106 briefly introduced in the remainder of this section.

### 107 1.1. Additivity of macromolecular interaction energies

108 We assume that tRNAs and aaRSs interact through sets of paired fea-  
109 tures that contribute additively to their overall binding energy as manifested  
110 through dissociation rate constants. This assumption has long featured in  
111 models of DNA-protein interactions in transcription factor binding sites and  
112 their evolution [47–52] as well as on the structure and evolution of protein-  
113 protein interaction networks [53–55]. Such studies have also used the ab-  
114 straction of working with simplified binary genotypes as we do here.

### 115 1.2. Kinetic proofreading

116 Hopfield [56] and Ninio [57] were the first to demonstrate the fundamen-  
117 tal mechanism of *kinetic proofreading* now shown to underlie the accuracy of  
118 information transduction in biopolymerization reactions such as aminoacyl-  
119 tRNA selection by the ribosome [58], tRNA selection in aminoacylation [59],  
120 and nucleotide selection in transcription [60], but also cellular signal trans-  
121 duction [61] (but see *e.g.* [62]), including T-cell activation [63] and recently,  
122 morphogenesis [64]. In kinetic prooreading, the dissipation of cellular free  
123 energy coupled to internal, allosteric non-reactive state transitions amplifies  
124 the kinetic discrimination of substrates at some combined expense of overall  
125 reaction rate, energy, and the stochastic discard of partially processed pre-  
126 ferred substrates [65, 66]. The discovery of proofreading was motivated in  
127 part by the observation that the amino acid selectivities of aaRSs are greater  
128 than can be explained by differences in free-energy of binding of different  
129 amino acids [56]. Ehrenberg and Blomberg [67] first derived the thermo-  
130 dynamic limits of kinetic proofreading in terms of the displacement from  
131 thermodynamic equilibrium of high energy cofactors such as ATP or GTP,  
132 as discussed by Kurland [68]. The kinetics of the two aaRSs classes is dif-  
133 ferent; In class I aaRSs, product release is rate-limiting, while in Class II  
134 aaRSs, aminoacyl transfer is rate-limiting [69]. However, a range of different  
135 regimes of kinetic rates and allosteric state transition networks can exhibit  
136 proofreading [65, 70]. In this work we use theoretical bounds for proofreading  
137 over all possible schemes to derive bounds on aminoacylation rates.

### 138 1.3. Rugged landscapes, epistatic gene interactions, and modifier models

139 Fitness landscapes, introduced by Sewall Wright [71], map genotypes to  
140 fitnesses either directly, or via phenotypes, as recently reviewed by Ahnert  
141 [72]. Interactions between genes can cause double, triple, *etc.* mutants to  
142 have greater or lesser fitness than expected from the isolated fitness effects  
143 of their component mutations, a phenomenon known as *epistasis*. Epistasis  
144 can take place across the genotype-phenotype map at multiple scales of bio-  
145 logical organization simultaneously [73]. Reciprocal sign epistasis (in which  
146 recombinants of haplotypes have lower fitness than non-recombinants) is a  
147 necessary (but not sufficient [74]) condition for fitness landscapes to become  
148 *rugged* [75], exhibiting potentially many separated local fitness maxima. Ab-  
149 stract genotype-fitness and genotype-phenotype-fitness models, such as the  
150 tunably rugged NK model [76, 77] or other regulatory or metabolic network  
151 evolution models [78–80] typically lack a concrete, mechanistic interpretation

152 for how epistasis actually manifests through the combined actions of genes  
153 on the basis of sequences.

154 In this work, we allow the availability of sites for matching or mis-  
155 matching between tRNA and aaRS gene products to evolve under direct  
156 genetic control, making epistasis evolveable at site resolution. As such, our  
157 work is related to population genetic models that study the genetic modi-  
158 fication of evolutionary forces such as mutation, recombination or epistasis,  
159 called *modifier models*. Modifier models encode evolutionary parameters at  
160 neutral loci that co-evolve under uniform genetic dynamics as other *major*  
161 *loci* that directly impact fitness. Original applications of modifier models  
162 were aimed at studying the evolution of recombination [81, 82]. Under very  
163 general conditions near an equilibrium under viability selection, modifier loci  
164 evolve to reduce rates of mutation, migration, or recombination [83]. With-  
165 out recombination, near mutation-selection balance, modifiers that increase  
166 positive or antagonistic epistasis will evolve, increasing the robustness of hap-  
167 loid asexual populations to mutations [14, 15], this robustness is an intrinsic  
168 property of the fitness landscape [84, 85]. An analysis of *fitness valley cross-*  
169 *ing* in asexual haploid populations with reciprocal sign epistasis [86] points  
170 to the critical role of the high-dimensional structure of fitness landscapes in  
171 determining evolutionary outcomes [87].

#### 172 1.4. *Origin-fixation formalism for evolutionary genetics*

173 We model evolutionary dynamics in finite, haploid asexual populations of  
174 aaRS-tRNA networks using the *statistical mechanical* or *sequential fixation*  
175 Markov chain [49, 88], a variety of *origin-fixation model* [89] that assumes a  
176 maximum of two genotypes segregating in a population at a given time. Thus,  
177 it is assumed that the mutation rate is much smaller than the reciprocal of the  
178 square of the population size [89]. These assumptions yield an exact solution  
179 of the stationary distribution of fixed genotypes in finite populations of con-  
180 stant size experiencing selection, mutation and genetic drift [88, 90]. Models  
181 of this kind have been used to highlight the role of compensatory evolution  
182 on the complex genotype-phenotype-fitness landscapes of transcription-factor  
183 binding sites [50] and proteins [91]. In an appendix, Sella [90] shows results  
184 for the stationary genotype distribution of a population of haploid binary  
185 genomes selected to maximize their weight (in the coding theory sense), that  
186 is, to become “all ones.” In the Discussion, we return to this model as a  
187 natural modeling complement to the binary match landscape models that  
188 we introduce here.

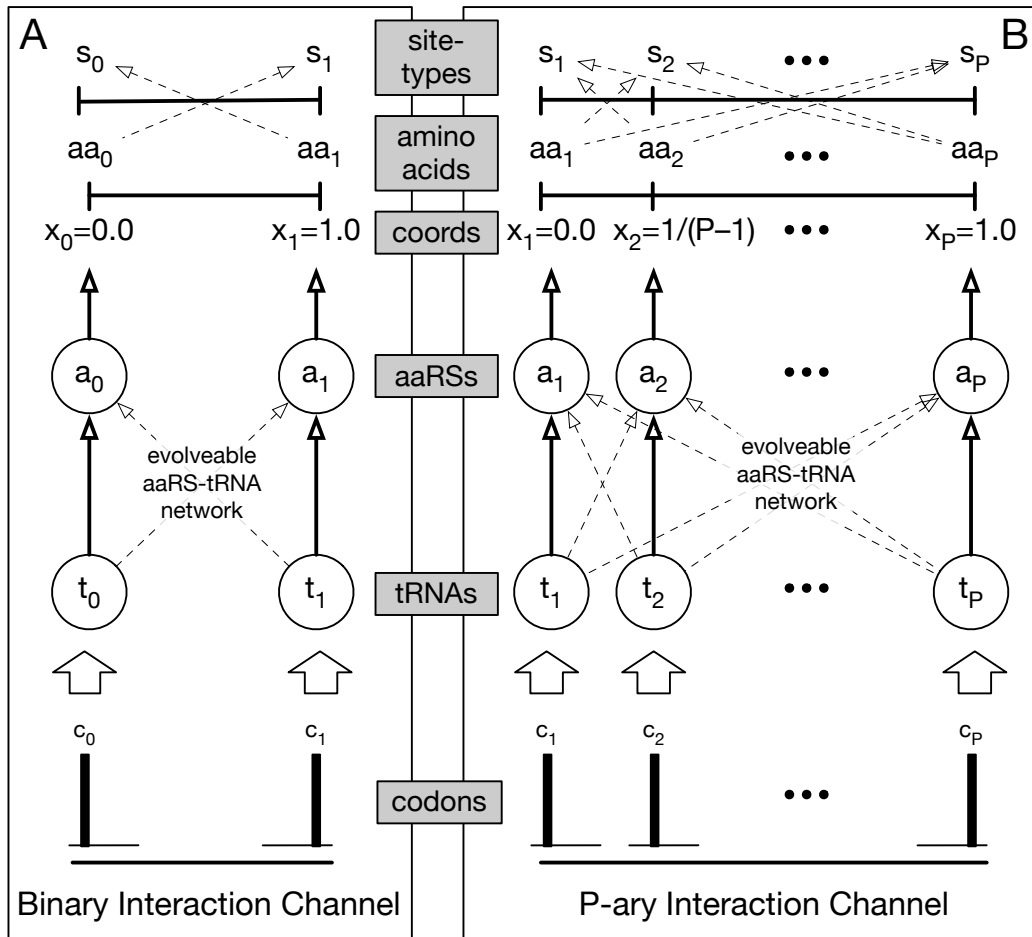


Figure 2: Set-up for models comparing fitness landscapes with different aaRS-tRNA networks and network encodings. Except in section 2.7, there are always a fixed and equal number  $P \geq 2$  species of tRNA,  $P$  species of aaRS,  $P$  codons,  $P$  available amino acids, and  $P$  site-types, the latter two of which are uniformly and maximally distributed within a one-dimensional amino-acid/site-type space representing differential selection on amino acid side chain properties such as hydrophobicity [92] (labelled as "coords"). To each site-type corresponds a unique codon that encodes it perfectly and a unique amino acid that fits it perfectly. To each codon corresponds a unique tRNA that reads it perfectly. To each amino acid corresponds a unique aaRS that charges it perfectly. Panel A. The Binary Interaction Network Channel ( $P = 2$ ) studied in subsection 2.3. Panel B. The P-ary Interaction Network Channel studied in subsections 2.8 and 2.9. (one-column figure)

## 189 2. Match landscapes: Models and Results

### 190 2.1. General Assumptions of the Current Work

191 Unless otherwise noted, genotypes  $g \in \mathbb{B}^L$  are haploid binary strings of  
 192 length  $L$  that undergo point mutation, selection and genetic drift in panmictic



193 Moran [93] populations of constant size — but experiencing neither recombina-  
194 tion, duplication, deletion, insertion, inversion, conversion, nor drive (nor  
195 implicitly, horizontal transfer) of genes or genomes. We define here various  
196 *match fitness landscapes* that map genotypes  $g \in \mathbb{B}^L$  to real-valued fitnesses  
197  $w(g)$  with  $0 \leq w(g) \leq 1$ . The fitness functions we define all have at their roots  
198 a *matching function* that maps genotypes to *match matrices*, which unam-  
199 biguously predict the intensities at which pairs of tRNA and aaRS species  
200 interact in a model cell or cytoplasmic volume. Except in subsection 2.7,  
201 each genotype  $g \in \mathbb{B}^L$  expresses an equal number  $P$  species of tRNA and  $P$   
202 species of aaRS.

203 Any species of tRNA can potentially match any species of aaRS through  
204 an interaction interface shared by all. Each species of tRNA or aaRS contains  
205 the same number of sites in this shared interaction interface. A correspon-  
206 dence exists that partitions sites in the same way across all species, and  
207 thereby limits the way in which matches of species can occur. We call the  
208 union of single sites over all species that can potentially match or mismatch  
209 within any possible aaRS-tRNA species pair a *site-block*. Matching occurs  
210 exclusively within site-blocks, and matching is additive over site-blocks. We  
211 denote the number of site-blocks  $n$  and call it the *width* of the interaction  
212 interface.

## 213 2.2. Overview of Models and Results

214 A list of symbols and parameter values is given in Table 1. In sub-  
215 section 2.3 we define a model we call the binary interaction channel with  
216 one site-block and compute its average fitness, load and epistasis under two  
217 different matching rules. In section 2.4, we define the P-ary interaction chan-  
218 nel with multiple site-blocks, while in section 2.5 we present a result about  
219 its stationary genotype frequency distribution when fitness is multiplicative  
220 over site-blocks. In section 2.6 we develop an additive interaction model  
221 for aaRSs and tRNAs. In section 2.7 we re-derive a macroscopic model of  
222 aminoacylation kinetics in an interaction network with  $N$  tRNA species and  
223  $M$  aaRS species. In section 2.8 we present results on the dependency of  
224 fitness maxima and fixed drift load on the number of cognate pairs encoded  
225 in an aaRS-tRNA network. In section 2.9 we compare fitnesses and the sta-  
226 tionary expected frequency of masking in networks selected for translational  
227 accuracy alone versus networks selected for both accuracy and rate.

228 *2.3. The Binary Interaction Channel with an Interface of One Site-Block*

229 Suppose that exactly one binary site in a gene for one tRNA species,  $t_0$ ,  
230 and another site in a gene for one aaRS species,  $a_0$ , are selected to match each  
231 other, so that genotypes 11 and 00 have equal and maximal viabilities greater  
232 than those of genotypes 10 and 01,  $w_{00} = w_{11} > w_{01} = w_{10}$ . This landscape is  
233 an example of “reciprocal sign epistasis.” [74, 86, 87]. In another landscape,  
234 one genotype, say 11, has higher viability than the other three, with  $w_{11} >$   
235  $w_{10} = w_{01} = w_{00}$ . This landscape is an example of positive or antagonistic  
236 epistasis [14], in which the fitness cost of the double mutant is less than either  
237 the sum or product of the costs of single mutants. The evolution of two-locus,  
238 two-allele models has been studied under very general settings, in haploid  
239 and diploid populations with and without recombination and modifiers of  
240 epistasis, most recently in the haploid setting by Liberman and Feldman [15].  
241 The minimal setting for a binary feed-forward interaction channel, encoding  
242 up to two amino acids, is only slightly more complex than the two-locus, two-  
243 allele model. It is a four-locus, two-allele model representing genes for two  
244 tRNA species  $t_0$  and  $t_1$  and two aaRS species  $a_0$  and  $a_1$ , in which either tRNA  
245 can potentially match either aaRS through a single site-block. Depending on  
246 the matching rule and the specific genotype, either of the two tRNA species  
247 may match zero, one or both aaRS species.

248 We define two different matching rules in our model through logical op-  
249 erations on bits. The first we call the *XNOR rule* and indicate it in Table 2  
250 and elsewhere with the  $\Leftrightarrow$  symbol. Using the XNOR rule, the match score  
251  $m_{i,j}^{\text{XNOR}}$  of  $t_i$  and  $a_j$ , with  $i, j \in \{0, 1\}$  is:

$$m_{i,j}^{\text{XNOR}} = t_i \Leftrightarrow a_j, \quad (1)$$

252 where  $(a \Leftrightarrow b) \equiv (a \odot b) \equiv \neg(a \oplus b)$  is the logical XNOR of  $a$  and  $b$ .

253 The second we call the *AND rule* and indicate it in Table 2 and elsewhere  
254 with the  $\wedge$  symbol. Using the AND rule, the match score  $m_{i,j}^{\text{AND}}$  of  $t_i$  and  $a_j$ ,  
255 with  $i, j \in 0, 1$  is:

$$m_{i,j}^{\text{AND}} = t_i \wedge a_j, \quad (2)$$

256 where  $(a \wedge b)$  is the logical AND of  $a$  and  $b$ .

257 According to the set-up in Panel A of Fig. 2, we suppose that all sources  
258 of ambiguity are collected into the network. The interaction of these four  
259 species of gene products occurs through a single site for each of them. Both  
260 aaRS species have equal concentration and efficiency, both tRNA species

261 have equal concentration and both amino acids have equal concentration.  
262 There are two equally frequent site-types  $s_0$  and  $s_1$  using the terminology  
263 and assumptions of [37, 92], one at coordinate  $x_0 = 0$  and the other at  
264 coordinate  $x_1 = 1$ . Amino acids  $aa_0$  and  $aa_1$  obtain maximal viability 1 in  
265 their respective site-types  $s_0$  and  $s_1$ . Amino acid  $aa_0$  obtains viability  $\phi$   
266 in site-type  $s_1$  and *vice versa*, while the viability of an unencoded amino acid  
267 (corresponding to when an aaRS species has no tRNA species that matches  
268 it) is  $\psi$ , with  $0 < \psi < \phi < 1$ . Only codons of type  $c_0$ , which are exclusively and  
269 perfectly read by tRNA species  $t_0$ , exist in sites of type  $s_0$ , while only codons  
270 of type  $c_1$ , which are exclusively and perfectly read by tRNA species  $t_1$ , exist  
271 in sites of type  $s_1$ . Amino acid  $aa_0$  is charged exclusively and perfectly by  
272 aaRS  $a_0$  and amino acid  $aa_1$  is charged exclusively and perfectly by aaRS  
273  $a_1$ . If a tRNA matches both aaRSs, the codons it reads achieve a fitness  
274  $\delta = (\phi + 1)/2$ , which is the arithmetic average of its translations. Thus,  
275 ambiguity is more fit than pure missense,  $\delta > \phi$ . The fitness of a genotype is  
276 the product of its fitness in the two site-types. With these assumptions, we  
277 write the fitnesses of the 16 possible genotypes under two different matching  
278 rules in Table 2.

279 Table 2 gives all genotype viabilities for the binary interaction channel  
280 with one site-block under the two different matching rules, XNOR and AND.  
281 The channel achieves greater maximum fitness using the XNOR rule because  
282 it can encode two interactions simultaneously with it, but only one with  
283 the AND rule. Inspecting the fitnesses of genotypes in consideration of the  
284 assumed inequality  $0 < \psi < \phi < \delta = (\phi + 1)/2 < 1$ , one finds that the  
285 fitness of every genotype with the XNOR rule is greater than or equal to  
286 its fitness with the AND rule. From eq. 9 in [90], one may infer directly  
287 that with these fitnesses under the stationary genotype distribution of the  
288 “sequential fixations” origin-fixation process [88, 90], the binary interaction  
289 channel has both a higher average fitness and a smaller fixed-drift load with  
290 the XNOR rule than it does with the AND rule, for all values of population  
291 size parameter  $\beta$  and for all  $0 < \psi < \phi < 1$ .

292 Liberman and Feldman [15] define multiplicative epistasis for the two-  
293 locus, two-allele model analogously to:

$$\epsilon_{2,2} = w_{11}w_{00} - w_{10}w_{01}. \quad (3)$$

294 A generalization of this expression to four loci and two alleles is:

$$\epsilon_{(4,2)} = (w_{1100}w_{0000} - w_{1000}w_{0100})(w_{1111}w_{0011} - w_{1011}w_{0111}) - (w_{1110}w_{0010} - w_{1010}w_{0110})(w_{1101}w_{0001} - w_{1001}w_{0101}). \quad (4)$$

295 After substituting fitnesses from Table 2 and simplification, we find that  
 296 the multiplicative epistasis  $\epsilon_{(4,2)}^{\text{AND}}$  of the AND rule is always positive:

$$\epsilon_{(4,2)}^{\text{AND}} = (1 - \phi)/2 > 0, \quad (5)$$

297 and that the multiplicative epistasis  $\epsilon_{(4,2)}^{\text{XNOR}}$  of the XNOR rule is also always  
 298 positive:

$$\epsilon_{(4,2)}^{\text{XNOR}} = (\delta - \phi\psi)^3(\delta + \phi\psi) > 0. \quad (6)$$

#### 299 2.4. The $P$ -ary Interaction Channel over an Interface of Multiple Site-Blocks

300 We now extend the model of section 2.3 by assuming that the interaction  
 301 intensities of  $P > 2$  tRNA species, labeled  $t_i$  with  $1 \leq i \leq P$ , and  $P$  aaRS  
 302 species, labeled  $a_j$  with  $1 \leq j \leq P$ , depend directly on their *match scores*  
 303  $m_{i,j}^R$  with matching rule  $R$ , which are additive over an interaction interface  
 304 of width  $n > 1$  site-blocks. To do so, we introduce two different combina-  
 305 tions of genotype spaces and matching rules to be used in the sequel. We  
 306 first define a genotype space  $G^{(P,P,n,1)}$  of dimension  $2Pn$  and explain how we  
 307 apply an XNOR matching function  $m_{i,j}^{\text{XNOR}}$  to genotypes from that space to  
 308 obtain the results of section 2.8. We then define a second larger genotype  
 309 space  $G^{(P,P,n,2)}$  of dimension  $4Pn$  and explain how we apply a more complex  
 310 matching function  $m_{i,j}^{\text{AND-XNOR}}$  on genotypes from that space to obtain the  
 311 results of section 2.9.

312 Assuming every species of tRNA or aaRS is produced by only one gene,  
 313 we assign  $n$  *state-bits* to each of the  $2P$  tRNA and aaRS genes and write  
 314 them as follows:  $t_i \equiv t_{i1}t_{i2} \dots t_{ir} \dots t_{in}$ , and  $a_j \equiv a_{j1} \dots a_{jr} \dots a_{jn}$  respectively,  
 315 where multiplication in this case implies string concatenation,  $1 \leq i, j \leq P$ ,  
 316  $1 \leq r \leq n$ , and  $t_{ir}, a_{jr} \in \mathbb{B}$ . We then order and concatenate genes into  
 317 genotypes as follows:  $g \equiv t_1a_1t_2a_2 \dots t_pa_p$ . Denote by  $G^{(P,P,n,1)}$  the set of all  
 318 possible binary genotypes with  $P$  tRNA genes and  $P$  aaRS genes of width  $n$   
 319 site-blocks and one site per-gene per-site-block, of total length  $L = 2Pn$ . For  
 320 any genotype  $g \in G^{(P,P,n,1)}$  the match score  $m_{i,j}^{\text{XNOR}}$  of  $t_i$  and  $a_j$  in the XNOR  
 321 matching function is defined as:

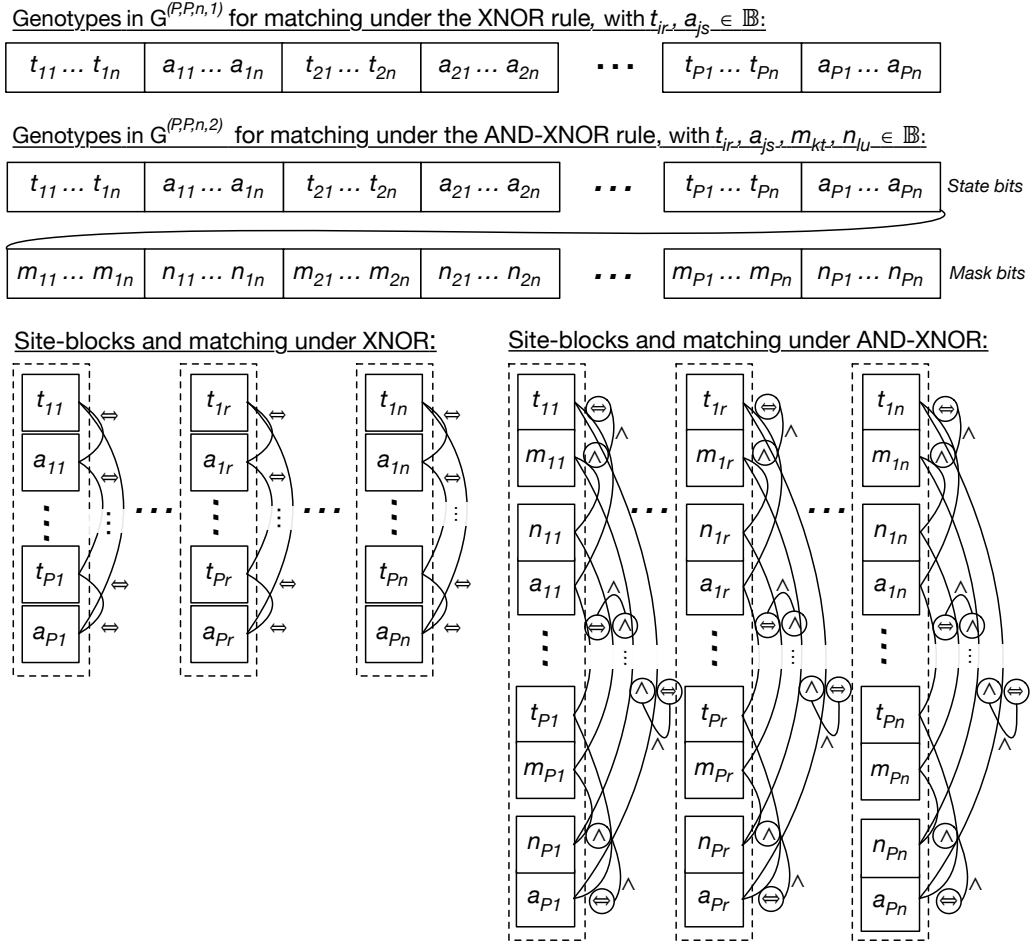


Figure 3: Genotype spaces, site-blocks and matching with the XNOR and AND-XNOR matching rules. (one or two-column figure)

$$m_{i,j}^{\text{XNOR}} = \sum_{r=1}^n t_{ir} \Leftrightarrow a_{jr}, \quad (7)$$

322 where  $(a \Leftrightarrow b) \equiv (a \odot b) \equiv \neg(a \oplus b)$  is the XNOR of  $a$  and  $b$ , true when  $(a \oplus b)$ ,  
 323 the XOR of  $a$  and  $b$ , is false. The XNOR match score  $m_{i,j}^{\text{XNOR}}$  of  $t_i$  and  $a_j$  is  
 324 inversely related to their Hamming distance  $d_H(t_i, a_j)$ :

$$m_{i,j}^{\text{XNOR}} = n - d_H(t_i, a_j). \quad (8)$$

325 We now introduce a third matching rule, which we call the *AND-XNOR*,  
 326 *MASKED-XNOR*, or *MASKED-MATCH* rule. Suppose that every macro-  
 327 molecular species adds one evolveable *mask bit* that switches on or off the ac-  
 328 cessibility for matching of exactly one of its *state bits* (Fig. 3). Mask-bits are  
 329 site-specific interaction modifiers. Now, with  $t_{ir}, a_{jr}, m_{ir}, n_{jr} \in \mathbb{B}$ ,  $1 \leq i, j \leq P$   
 330 and  $1 \leq r \leq n$ , we assign  $n$  state-bits to each of the  $P$  tRNA genes as before,  
 331 writing the state-bits of tRNA gene  $t_i$  as  $t_{i1} \dots t_{ir} \dots t_{in}$ , and in addition, we  
 332 assign  $n$  mask-bits to each of the  $P$  tRNA genes, writing the mask-bits of  
 333 tRNA gene  $t_i$  as  $m_{i1} \dots m_{ir} \dots m_{in}$ , so that  $m_{ir}$  is the mask-bit correspond-  
 334 ing to state-bit  $t_{ir}$ . Similarly, we assign  $n$  state-bits to aaRS gene and write  
 335 them as  $a_{j1} \dots a_{jr} \dots a_{jn}$ . In addition, we assign  $n$  mask-bits to each of the  
 336  $P$  aaRS genes, and write the mask-bits of aaRS gene  $a_j$  as  $n_{j1} \dots n_{jr} \dots n_{jn}$ ,  
 337 so that  $n_{jr}$  is the mask-bit corresponding to state-bit  $a_{jr}$ . Finally, we order  
 338 and concatenate genes into genotypes as follows (without loss of generality):  
 339  $g \equiv t_1 a_1 t_2 a_2 \dots t_P a_P m_1 n_1 m_2 n_2 \dots m_P n_P$ . Denote by  $G^{(P,P,n,2)}$  the set of all  
 340 possible binary genotypes with  $P$  tRNA genes and  $P$  aaRS genes interact-  
 341 ing over width  $n$  site-blocks, with 2 sites per-gene per-site-block, and a total  
 342 length  $L = 4Pn$ . For any genotype  $g \in G^{(P,P,n,2)}$  the match score  $m_{i,j}^{\text{AND-XNOR}}$   
 343 of  $t_i$  and  $a_j$  with AND-XNOR matching rule is defined:

$$m_{i,j}^{\text{AND-XNOR}} = \sum_{r=1}^n ((m_{ir} \wedge n_{jr}) \wedge (t_{ir} \Leftrightarrow a_{jr})), \quad (9)$$

344 where  $(a \wedge b)$  is the logical AND of  $a$  and  $b$ .

### 345 2.5. $P$ -ary interaction channels with multiplicative fitness over site-blocks

346 Let fitness depend multiplicatively on the match scores of corresponding  
 347 tRNA, aaRS species pairs (*i.e.* those that share the same index), and in-  
 348 versely on the match scores of non-corresponding tRNA, aaRS species pairs  
 349 (*i.e.* those with different indices). For example, if the fitness contributions of  
 350 a match between any cognate pair or of mismatch between any non-cognate  
 351 pair, one might define the viability fitness  $w(g)$  of genotype  $g \in G^{(P,P,n,1)}$  as:

$$w(g) = \frac{\prod_{i=1}^P \prod_{(j=1) \neq i}^P \phi^{m_{i,j}^{\text{XNOR}}}}{\phi^{(P-1)} \prod_{i=1}^P \phi^{m_{i,i}^{\text{XNOR}}}} = \frac{\prod_{i=1}^P \phi^{d_H(t_i, a_i)}}{\prod_{i=1}^P \prod_{(j=1) \neq i}^P \phi^{d_H(t_i, a_j)}}, \quad (10)$$

352 where  $0 < \phi \leq 1$  is a selection intensity parameter. The viabilities of eq. 10 are  
 353 positive and less than or equal to 1, and increase both as tRNAs and aaRSs of

354 the same index match while tRNAs and aaRSs of different indices mismatch.  
355 In the appendix, we show that the function in eq. 10 is multiplicative over site-  
356 blocks as previously defined, and that for all fitness functions multiplicative  
357 over site-blocks, the stationary distribution of fixed genotypes of [90] may  
358 readily be obtained as a product of the stationary frequencies of site-blocks.  
359 This result should be compared to Result 2 in [15], which states that in a  
360 large two-locus, two-allele haploid population in mutation-selection balance,  
361 a unique polymorphic equilibrium with full linkage equilibrium exists only in  
362 the absence of multiplicative epistasis.

### 363 2.6. From additive interaction energies to kinetic rate constants

364 As simple and tractable as the fitness function in eq. 10 may be, it is more  
365 realistic to suppose that the fitness of an aaRS-tRNA network is manifested  
366 through its translation of protein-coding genes. We therefore wish to create a  
367 decoding function that takes a match matrix as input and outputs a *decoding*  
368 *matrix* that specifies the conditional aminoacylation profile of every tRNA  
369 species.

370 We assume through the sequel that matches  $m_{i,j}$  between tRNA species  
371  $t_i$  and aaRS species  $a_j$  contribute additively to their binding energy in an  
372 aaRS-tRNA complex (whether activated or not), and that only one kinetic  
373 rate constant depends on this energy and varies from complex to complex  
374 with all other kinetic rate constants set equal (see next section). Table 5  
375 in Schimmel and Söll [94] displays kinetic data for aaRS-tRNA complexes  
376 with data from [96, 97] of about  $220 \text{ s}^{-1}$  for cognate aaRS-tRNA complexes  
377 and about  $1600 \text{ s}^{-1}$  for near-cognate interactions. We assumed a cognate  
378 dissociation rate constant of  $k_d^c = 220 \text{ s}^{-1}$  and a non-cognate dissociation rate  
379 constant of  $k_d^{nc} = 10\,000 \text{ s}^{-1}$  representing the background energy of interaction  
380 between tRNAs and aaRSs, also comparable to data in [98].

381 Define  $k$  as the number of matches required to diminish dissociation rate  
382 from  $k_d^{nc}$  to  $k_d^c$ , with  $1 \leq k \leq n$ . Following Johnson and Hummer [54], we  
383 calculate non-cognate and cognate equilibrium constants as reciprocals of the  
384 non-cognate and cognate dissociation rates. The dissociation rate constant  
385  $k_d^{i,j}$  between tRNA  $t_i$  and aaRS  $a_j$  with  $m_{i,j}$  matches,  $0 \leq m_{i,j} \leq n$  then may  
386 be defined

$$k_d^{i,j} = k_d^{nc} \exp[\nu m_{i,j}], \quad (11)$$

387 where  $\nu = (\log k_d^{nc} - \log k_d^c)/k$ .

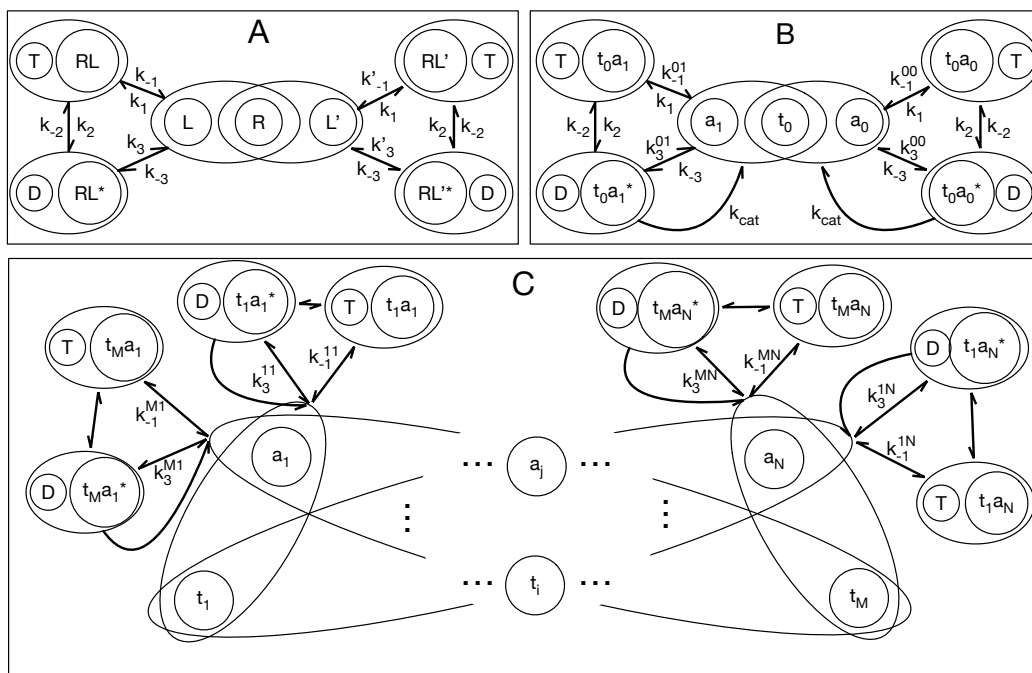


Figure 4: Application of kinetic proofreading schemes to compute decoding rates in a macroscopic interaction network of  $M \geq 2$  species of tRNAs and  $N \geq 2$  species of aaRSs, including kinetic proofreading of tRNAs by aaRSs but presently ignoring errors in amino acid selection by aaRSs or tRNA selection on ribosomes. A. The single-molecule two-cycle, two-state kinetic proofreading scheme of [61, fig. 1] for a receptor  $R$  that can preferentially select ligands of type  $L'$  over ligands of type  $L$ , assuming  $k_1 = k'_1$ ,  $k_{-3} = k'_{-3}$  and  $k_2 = k'_2$  are all pseudo-first-order rate constants and the concentrations of ligand species are equal, *i.e.*  $[L] = [L'] \gg [R]$ . B. The scheme from panel A redrawn from the perspective of a single tRNA species  $t_0$  alternatively aminoacylated (and instantaneously deacylated) by two aaRS species  $a_0$  and  $a_1$  of equal concentrations through catalytic steps with rate  $k_{cat} \ll k_3$ , thus  $[t_0] \gg [a_0] = [a_1] \gg 1$ . C. Generalization of the scheme in B to  $M$  species of tRNAs and  $N$  species of aaRSs. All corresponding rate constants are assumed equal across all interactions except those indicated. (one- or two-column figure)

388 *2.7. Decoding functions for macroscopic, well-mixed proofreading aaRS-tRNA*  
 389 *networks*

390 We now assume that matching feature-set-pairs contribute additively to  
 391 interaction energies between species pairs and transform interaction energies  
 392 into kinetic rates of dissociation, or off-rates, of aaRS-tRNA species-pair  
 393 complexes (in this section,  $k_{-1}^{i,j}$  is the same as  $k_d^{i,j}$  in eq. 11). We elab-  
 394 orate on the reaction scheme shown in fig. 4A to compute the decoding



395 rate/aminoacylation probability of one species of tRNA  $t_0$  interacting in well-  
396 mixed volume with two species of aaRSs  $a_0$  or  $a_1$  at equal concentration as in  
397 fig. 4B. We then generalize this to calculate the *maximal decoding probability*  
398  $c_{\max}(t_i \rightarrow a_j)$  that a tRNA of species  $t_i$ , with  $1 \leq i \leq M$  was last aminoacy-  
399 lated by an aaRS of species  $a_j$ , with  $1 \leq j \leq N$  in an aaRS-tRNA network  
400 of  $M$  tRNA species and  $N$  aaRS species with variable dissociation rate con-  
401 stants  $k_{-1}^{i,j}$  that vary between complexes of different aaRS-tRNA species pairs  
402 (fig. 4C). A comparable development was presented in [99], who were partic-  
403 ularly interested in the energy costs of proofreading.

404 Qian [61] re-cast the classic Hopfield kinetic proofreading model as the  
405 five-state Markov Chain shown in fig. 4A, describing a cell signalling receptor  
406  $R$  with a two-step activation scheme that discriminates against ligand  $L$   
407 in favor of ligand  $L'$  via off-rates (dissociation rates). The error rate per-  
408 receptor  $f$  is the ratio of activated receptor affinities with ligands  $L$  and  
409  $L'$ . Qian [61] computed the minimum error rate per-receptor  $f_{\min}$  for any  
410 set of kinetic constants in terms of the dissociation-rate-constant ratio  $\theta =$   
411  $k'_{-1}/k_{-1} < 1$  and an exponential function of the steady-state free energy of the  
412 cell  $\gamma = e^{(\Delta G_{DT}/RT)} \geq 1$ , associated with the (deliberately unbalanced) coupled  
413 reactions  $T \rightleftharpoons D$  in fig. 4, namely  $f_{\min} = \theta((1 + \sqrt{\gamma\theta})/(\sqrt{\gamma} + \sqrt{\theta}))^2$ . Qian [61] also  
414 re-derived the absolute lower thermodynamic limit over all possible kinetic-  
415 proofreading schemes [67], and the classical minimum per-receptor error-rate  
416  $f_{\min}$  in the two-state kinetic proofreading scheme shown in fig. 4A, with  
417  $\theta^2 \leq f_{\min} \leq \theta$  [56, 57]. These two bounds correspond to perfect proofreading  
418 (with infinite ATP) on the left and thermodynamic equilibrium/recent death  
419 on the right.

420 These results apply equally well to enzymes as the rate of catalysis ( $k_{cat}$   
421 in figs. 4B and 4C) vanishes. This is one of three conditions on the kinetic  
422 rate constants that achieve the minimum error rate  $f_{\min}$  [56, 61]. To achieve  
423 accuracy, enzymes and receptors add states from which they discard cognate  
424 substrates at appreciable rates so they can give non-cognate substrates more  
425 time to dissociate.

426 If the concentrations of aaRSs are large and equal to each other, the  
427 treatment of Qian [61] applies to Fig. 4B even though the roles of ligand and  
428 receptor are reversed. Let us define  $\theta_{001} \equiv (k_{-1}^{00}/k_{-1}^{01})$  as the ratio of dissociation  
429 rate constants of tRNA  $t_0$  with aaRS  $a_0$  and aaRS  $a_1$ , and similarly  $\theta_{011} \equiv$   
430  $(k_{-1}^{01}/k_{-1}^{01} = 1)$ . Then, at steady state, the relative rate of aminoacylation of  
431 tRNA  $t_0$  by aaRS  $a_1$  versus aaRS  $a_0$  may be written  $f_{\min} = [t_0 a_1^*]/[t_0 a_0^*]$ ,

432 bounded by  $\theta_{001}^2 \leq f_{\min} \leq \theta_{001}$ , and the time-averaged maximal decoding  
 433 probability  $c_{\max}(t_0 \rightarrow a_1)$  that tRNA  $t_0$  was last aminoacylated by aaRS  $a_1$   
 434 is:

$$c_{\max}(t_0 \rightarrow a_1) = [t_0 a_1^*] / ([t_0 a_0^*] + [t_0 a_1^*]), \quad (12)$$

435 with

$$H(\theta_{001}^2, \theta_{011}^2) / 2 \leq c_{\max}(t_0 \rightarrow a_1) \leq H(\theta_{001}, \theta_{011}) / 2, \quad (13)$$

436 where  $H(\alpha, \beta)$  is the harmonic average of  $\alpha$  and  $\beta$ . The maximal decoding  
 437 probability is maximal over all kinetic schemes of aminoacylation; however,  
 438 by the data processing inequality, it is also the maximal accuracy of transla-  
 439 tion over all error-rates in tRNA-selection by ribosomes.

440 More generally, let us define  $\theta_{ikj}$  as the ratio of dissociation rate constants  
 441 of tRNA  $t_i$  with aaRS  $a_k$  and aaRS  $a_j$  respectively, *i.e.*  $\theta_{ikj} \equiv k_{-1}^{ik} / k_{-1}^{ij}$ ,  
 442 with  $1 \leq i \leq M$  and  $1 \leq j, k \leq N$ . The maximal decoding probability  
 443  $c_{\max}(t_i \rightarrow a_j)$ , that a tRNA of species  $t_i$  was last aminoacylated by an aaRS  
 444 of species  $a_j$  in an aaRS-tRNA network of  $M$  species of tRNA and  $N$  species  
 445 of aaRS, is

$$c_{\max}(t_i \rightarrow a_j) = [t_i a_j^*] / \left( \sum_{k=1}^N [t_i a_k^*] \right), \quad (14)$$

446 with

$$(H_{k=1}^N \theta_{ikj}^2) / N \leq c_{\max}(t_i \rightarrow a_j) \leq (H_{k=1}^N \theta_{ikj}) / N, \quad (15)$$

447 where  $H_{k=1}^N \theta_{ikj}$  is the harmonic average over all  $\theta_{ikj}$ ,  $1 \leq k \leq N$ .

## 448 2.8. The Dependence of Load on Number of Encoded Amino Acids

449 Drawing on the terminology and concepts of earlier work [37, 92, 100, 101],  
 450 we present a highly simplified translational system to compare fitnesses and  
 451 stationary genotype frequencies of different matching rules. With reference  
 452 to Fig. 2B, we continue to assume  $P$  pairs of aaRS and tRNA species, as well  
 453 as  $P$  species of codons, amino acids, and site-types, so that tRNA species  
 454  $t_i$ ,  $1 \leq i \leq P$  always reads codon  $c_i$ , while aaRS  $a_i$  always charges amino acid  
 455  $aa_i$ , which has maximal fitness in sites of type  $s_i$ . With these assumptions,  
 456 the decoding probability  $c(aa_j | c_i)$  of decoding codon  $c_i$  as amino acid  $aa_j$  is  
 457 equal to  $c_{\max}(t_i \rightarrow a_j)$  of the last section,  $c(aa_j | c_i) \equiv c_{\max}(t_i \rightarrow a_j)$ . The

458 fitness  $w(aa_j|s_l)$  of amino acid  $aa_j$  in site-type  $s_l$ , with  $1 \leq j, l \leq P$  is  $\phi^{|x_j-x_l|}$ ,  
459 where  $x_i = (i-1)/(P-1)$ . The fitness  $w_l$  in site-type  $s_l$  is the expected fitness  
460 of translations of codons occupying that site-type, here exclusively codon  $c_l$ ,  
461 *i.e.*  $w_l = \sum_j^P w(aa_j|s_l)c(aa_j|c_l)$ . The fitness  $w_A(g)$  of genotype  $g$  selected for  
462 translational accuracy alone is the product of its fitnesses over all site-types:

$$w_A(g) = \prod_l^P w_l. \quad (16)$$

463 We implemented this model in a Python 3 script called “atINFLAT” for  
464 “aaRS-tRNA Interaction Network Fitness Landscape Topographer,” avail-  
465 able as supplementary data. It can compute the stationary genotype distri-  
466 butions of small networks and compute statistics such as fitnesses for indi-  
467 vidual genotypes from much larger networks.

468 It is easy to prove that binary codes with zero matches between any code-  
469 words have a maximum size of only two codewords [102, 103]. Thus, with the  
470 XNOR rule, in which tRNAs and aaRSs may potentially match or mismatch  
471 over their entire shared interface, the interactions of only two aaRS-tRNA  
472 pairs may be encoded perfectly without cross-matching. As predicted, when  
473 we used atINFLAT to compute maximum and average fitnesses on landscapes  
474 with and without proofreading, we found that both the maximum fitness de-  
475 creased and fixed-drift load increased when more than two cognate pairs were  
476 overloaded on the same interaction interface, reflecting an increasing cost of  
477 translational missense as more amino acids get encoded (Fig. 5).

### 478 *2.9. Selection on both translational accuracy and rate is necessary to select* 479 *for masking to reduce cross-matching*

480 The symmetric P-ary interaction channel as we have defined it, selects  
481 only for translational accuracy and not on rate or energy expenditure. One  
482 can see this clearly with the help of a well-defined example using the AND-  
483 XNOR rule, and comparing the fitnesses of two genotypes  $g_H, g_M \in G^{(4,4,8,2)}$ .  
484 The first genotype,  $g_H$ , consists of four codewords from the Hamming [n=8,d=4]  
485 code [95] repeated twice, followed by all maskbits set:

$$g_H = (10000111)^2(01001011)^2(00101101)^2(00011110)^2 1^{64}. \quad (17)$$

486 Since all maskbits are set in  $g_H$ , all four tRNA species and all four aaRS  
487 species potentially match over their entire interfaces. The cognate match  
488 score for all pairs  $t_i, a_i$  is  $m_{i,i} = 8$  and the single non-cognate match score

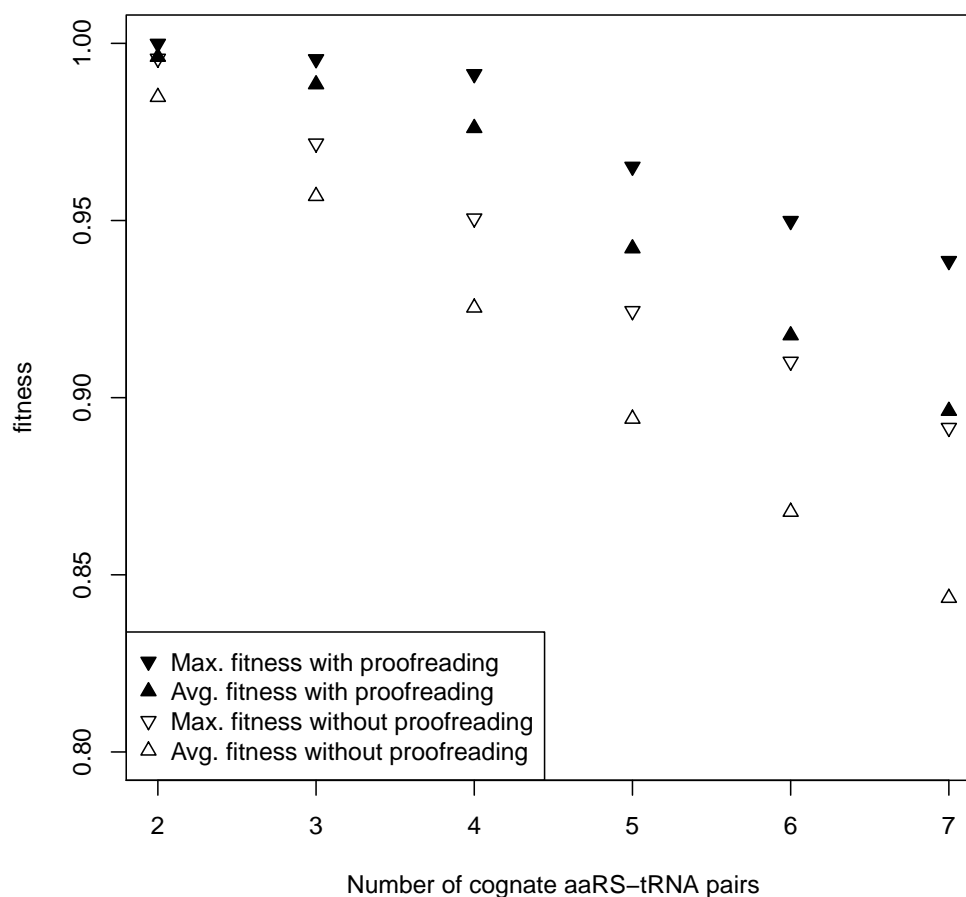


Figure 5: Decreasing average and maximal fitness of aaRS-tRNA networks as a function of encoded interactions under the XNOR rule, with and without proofreading. Parameters used here are  $n = 2$ ,  $k = 2$ ,  $\phi = 0.9$ , and  $\beta = 100$ . The fixed-drift loads are the differences between maximal and average fitnesses, which increase with the number of encoded interactions. Notice the discontinuities between  $P = 4$  and  $P = 5$ ; this is the transition where  $P > 2^n$ , the number of pairs exceeds available codewords. (one-column figure)

489 is  $m_{i,j} = 4$  for all pairs  $t_i, a_j$  with  $i \neq j$ . No binary codes of size  $n = 8$  can  
 490 achieve a larger minimum Hamming distance than four [104, 105].

491 A second genotype  $g_M$  may be constructed from any two tetramers and

492 their complements in the left and right halves of the interface, and using  
 493 masking to eliminate cross-matching. For example,

$$g_M = (1^{16}0^{16})^2(1^40^4)^2(0^41^4)^2, \quad (18)$$

494 which achieves cognate matches ( $m_{i,i} = 4$  for all pairs  $t_i, a_i$ ) and zero cross-  
 495 matching ( $m_{i,j} = 0$  for all pairs  $t_i, a_j$  with  $i \neq j$ ). With the standard fitness  
 496 function that we have been using in which fitness depends only on accu-  
 497 racy and not rate of translation and using  $k = 4$ , the fitnesses of these two  
 498 genotypes are exactly equal:

```

499 # atinflat version 0.8
500 # execution command:
501 # atinflat.py --pairs 4 --width 8 --match 4 --mask --phi 0.9
502 #           -g hamming-8-4.txt
503 genotype:
504 1000011110000111010010110100101100101101001011010001111000011110
505 111111111111111111111111111111111111111111111111111111111111111
506 | fitness: 0.9996721776752496
507 | match: [[8 4 4 4], [4 8 4 4], [4 4 8 4], [4 4 4 8]]
508 | proofread code: [[1. 0. 0. 0.], [0. 1. 0. 0.],
509 |                 [0. 0. 1. 0.], [0. 0. 0. 1.]]
510 genotype:
511 1111111111111111000000000000000000111111111111111000000000000000
512 1111000011110000111100001111000000001111000011110000111100001111
513 | fitness: 0.9996721776752496
514 | match: [[4 0 0 0], [0 4 0 0], [0 0 4 0], [0 0 0 4]]
515 | proofread code: [[1. 0. 0. 0.], [0. 1. 0. 0.],
516 |                 [0. 0. 1. 0.], [0. 0. 0. 1.]]
    
```

517 The example illustrates a key property of our macroscopic kinetic match  
 518 landscape model, which is that accuracy depends on relative dissociation  
 519 rate constants and concentrations, a prediction borne out by experimental  
 520 evidence [34, 94, 106]. We conjecture that these two genotypes have maximal  
 521 fitness because they both achieve the maximal possible distance of four be-  
 522 tween all code words — and they are not alone; many others in their neutral  
 523 network have the same fitness. Other genotypes with equal fitness to  $g_H$  and  
 524  $g_M$  include all those with the structure of  $g_H$  but substituting any four of the  
 525 16 Hamming [8,4] codewords in any order, in any one of  $2 \times 8!$  permutations

526 of codeword columns and codeword symbols (implying a degeneracy of more  
527 than  $3.522 \times 10^9$  Hamming code genotypes) as well as other non-linear per-  
528 fect binary codes [107] — all with every mask-bit set — and a much smaller  
529 number of those with the same structure as  $g_M$  and half of the mask-bits  
530 off, using one of only 255 combinations of two tetramer codewords and their  
531 complements besides those used in  $g_M$ .

532 Even though  $g_H$  and  $g_M$  achieve identical accuracy and fitness in the  
533 match landscape with  $w_A(g), g \in G^{(4,4,8,2)}$ , the rates of translation in cells  
534 with genotype  $g_H$  would be vastly slower than in cells with genotype  $g_M$ ,  
535 because the dissociation (discard) rate of cognate complexes is only between  
536  $1 \text{ s}^{-1}$  and  $2 \text{ s}^{-1}$  in the former, while in the latter it is the typical cognate rate  
537 that we assumed,  $220 \text{ s}^{-1}$ . In the classic kinetic proofreading schemes, this  
538 discard rate must be much greater than the actual rate of product formation  
539  $k_{\text{cat}}$  [56, 61] (but see [65, 66, 70]). For example, in tRNA-Ile of *Salmonella*  
540 *typhimurium* this rate is estimated to be  $5 \text{ s}^{-1}$  [108]. Furthermore, the overall  
541 rate of protein synthesis, which factors directly into growth rate [109], can  
542 be limited by the slowest rate of aminoacylation [110, 111]. As a result, both  
543 the accuracy and rate of translation are expected to factor into fitness [112].  
544 Because the fitnesses of  $g_H$  and  $g_M$  are exactly equal without taking transla-  
545 tional rate into account, incorporating any rate-dependent fitness factor that  
546 decreases with the cognate aminoacylation rate in our model will disadvan-  
547 tage those genotypes that maximize matching between cognate complexes.  
548 Selection for accuracy should then select for mask bits to turn off to reduce  
549 cross-matching and maintain the high non-cognate/cognate dissociation rate  
550 ratios required for accuracy at intermediate levels of cognate matching.

551 To test this prediction, we introduce an empirically parametrized fitness  
552 factor that crudely penalizes cognate aminoacylation rates when they are  
553 slower than the assumed cognate rate of  $220 \text{ s}^{-1}$ . In accordance with an  
554 observation of  $k_{\text{cat}} = 5 \text{ s}^{-1}$  [108] and a cognate dissociation/discard rate of  
555  $220 \text{ s}^{-1}$ , we define the average aminoacylation rate  $k_{\text{cat}}(g)$  of genotype  $g$  as  
556 proportional to the harmonic mean of cognate dissociation rates between  
557 cognate tRNAs and aaRSs:

$$k_{\text{cat}}(g) = \frac{1}{44} H_{i=1}^n k_d^{i,i}. \quad (19)$$

558 Controlled measurements with wild-type and mutant enzymes showed that  
559 only  $k_{\text{cat}}$  correlated with growth rate and the following measurements of  
560  $(k_{\text{cat}}, w)$  were observed, where  $w$  is growth rate in Luria Broth, written rela-

561 tive to wild-type [108, Table 3]:  $\{(0.19, 0.24), (0.6, 0.6), (5, 1)\}$ .

562 Using GNUPLOT 5.2 to fit two exponential viability functions  $w_1(k_{\text{cat}}) =$   
563  $A+B \exp(C_1 k_{\text{cat}})$  and  $w_2(k_{\text{cat}}) = 1 - \exp(C_2 k_{\text{cat}})$  to these data and also through  
564 the origin, we obtained the following fits:

$$w_1(k_{\text{cat}}) = 1.00127 - 1.005 \exp(-1.51245 k_{\text{cat}}) \quad (20)$$

$$w_2(k_{\text{cat}}) = 1 - \exp(-1.50576 k_{\text{cat}}), \quad (21)$$

565 both with a root mean square residual of less than 1%.

566 We defined a new fitness function  $w_{AR}(g)$  to select for both translational  
567 accuracy and rate as the product of two fitness factors:

$$w_{AR}(g) = w_A(g)w_2(k_{\text{cat}}(g)). \quad (22)$$

568 Using this new fitness function  $w_{AR}(g)$  and  $k = 4$ , we obtained the follow-  
569 ing results:

```
570 # atinflat version 0.8
571 # execution command:
572 # atinflat.py --pairs 4 --width 8 --match 4 --mask --phi 0.9
573 #           --rate -g hamming-8-4.txt
574 #
575 genotype:
576 100001111100001111010010111010010111001011101001011101000111110000111110
577 111111111111111111111111111111111111111111111111111111111111111111111111
578 | fitness: 0.0036361253612561006
579 | match: [[8 4 4 4], [4 8 4 4], [4 4 8 4], [4 4 4 8]]
580 | proofread code: [[1. 0. 0. 0.], [0. 1. 0. 0.],
581 |                 [0. 0. 1. 0.], [0. 0. 0. 1.]]
582 genotype:
583 11111111111111111111000000000000000001111111111111111111110000000000000000
584 111100001111000011110000111100000000111100001111000011110000111100001111
585 | fitness: 0.9991349818561294
586 | match: [[4 0 0 0], [0 4 0 0], [0 0 4 0], [0 0 0 4]]
587 | proofread code: [[1. 0. 0. 0.], [0. 1. 0. 0.],
588 |                 [0. 0. 1. 0.], [0. 0. 0. 1.]]
```

589 Even with  $k = 8$ , so the assumed cognate dissociation rate is only reached  
590 with a full eight matches, the masked genotype still has higher fitness:

```

591 # atinflat version 0.8
592 # execution command:
593 # atinflat.py --pairs 4 --width 8 --match 8 --mask --phi 0.9
594 #             --rate -g hamming-8-4.txt
595 #
596 genotype:
597 1000011110000111010010110100101100101101001011010001111000011110
598 111111111111111111111111111111111111111111111111111111111111111
599 | fitness: 0.9855421043520338
600 | match: [[8 4 4 4], [4 8 4 4], [4 4 8 4], [4 4 4 8]]
601 | proofread code: [[0.94 0.02 0.02 0.02], [0.02 0.94 0.02 0.02],
602                  [0.02 0.02 0.94 0.02], [0.02 0.02 0.02 0.94]]
603 genotype:
604 1111111111111111000000000000000001111111111111111000000000000000
605 1111000011110000111100001111000000001111000011110000111100001111
606 | fitness: 0.9860719918123261
607 | match: [[4 0 0 0], [0 4 0 0], [0 0 4 0], [0 0 0 4]]
608 | proofread code: [[0.94 0.02 0.02 0.02], [0.02 0.94 0.02 0.02],
609                  [0.02 0.02 0.94 0.02], [0.02 0.02 0.02 0.94]]

```

Hamming codes are efficient with respect to codeword length [95]. In this work, codewords are transmitted in parallel, so selection on code-word length  $g_M$  occurs through selection to avoid overly tight binding. Our results show that genetic match codes can be selected to sacrifice code-words to achieve shorter codeword length without cross-matching.

Our results are general. In Fig. 6, we show the full stationary genotype distributions under two fitness functions  $w_A(g)$  and  $w_{AR}(g)$  on the smaller genotype space  $G^{(4,4,2,2)}$  and  $k = 1$ , showing that masking is systematically favored over the entire match landscape and increasingly so with genotype fitness, under combined selection on translational accuracy and rate. Thus, selection on both the specificity of association and rate of dissociation can partition macromolecular interaction interfaces to reduce cross-matching.

Natural selection increases and maintains information in genomes [113–116]. A useful measure of this information is the reduction in entropy of the stationary distribution of genotypes with that selection, relative to without it. For example, the maximum entropy of genotypes in  $G^{(4,4,2,2)}$  occurs on a perfectly flat fitness landscape in which all genotypes have equal fitness, and its value is the genome length in bits, 32. For the data in Fig. 6 with



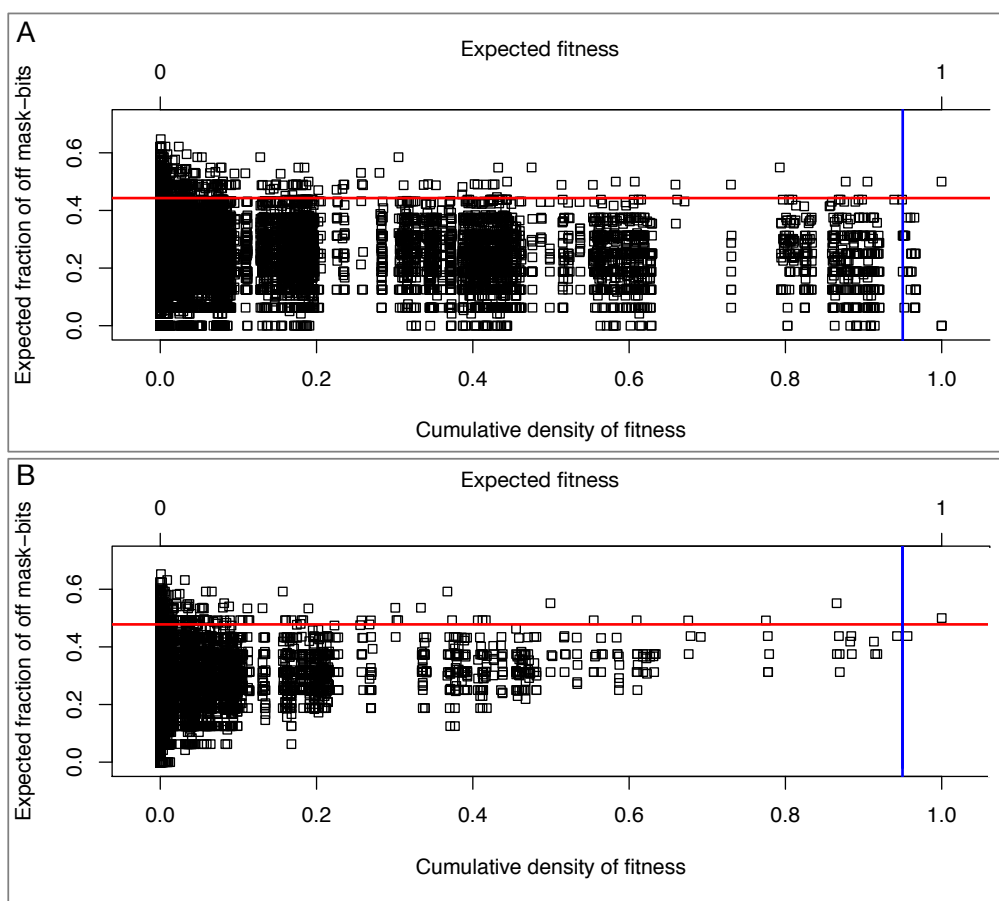


Figure 6: Expected fractions of masked sites in the steady-state fitness equivalence classes of  $2^{32}$  genotypes in  $G^{(4,4,2,2)}$  (points), expected fractions of masked sites (red lines) and expected fitnesses (blue lines) as functions of the stationary cumulative densities of fitness in match landscapes with perfect one-step kinetic proofreading,  $P = 4$ ,  $n = 2$ ,  $k = 1$ ,  $\phi = 0.9$ , and  $\beta = 100$ . A. Match landscape with selection for translational accuracy alone (fitness function  $w_A(g)$ ) with expected fitness 0.9501304 and expected fraction of masked sites 0.4428638. B. Match landscape with combined selection for translational accuracy and rate (fitness function  $w_{AR}(g)$ ), with expected fitness is 0.9498575 and expected fraction of masked sites 0.4780706. Machine error in these data, as judged by the integration of cumulative density functions, is less than  $10^{-11}$ . (two-column figure)

628 perfect kinetic proofreading and  $\beta = 100$ , we found that the entropy of the  
 629 stationary genotype distribution under selection for accuracy alone, through  
 630 the fitness function  $w_A(g)$ , is about 7.82 bits for a maximum information gain

631 of about 24.18 bits. The entropy of the stationary genotype distribution  
632 under combined selection for both accuracy and rate, through the fitness  
633 function  $w_{AR}(g)$ , is about 6.22 bits for a larger maximum information gain of  
634 about 25.78 bits. Thus, in a population of fixed size 101, about 1.6 more bits  
635 of information are gained under combined selection for rate and accuracy  
636 than under selection for accuracy alone. Without proofreading, the results  
637 are not very different: the maximum information gained under selection for  
638 accuracy alone is close to 23.5 bits, while under combined selection for both  
639 accuracy and rate, the maximum information gain is close to 25.26 bits.

### 640 3. Discussion

641 We have shown that combined selection on translational accuracy and rate  
642 is sufficient to select for divergence in tRNA-interaction interfaces by aaRSs.  
643 Our results do not contradict other hypotheses about this phenomenon [30].  
644 We used mask bits as interaction modifiers to demonstrate our main re-  
645 sult. When they mask or diminish interactions, these modifier bits may be  
646 interpreted as the presence of structural features such as *identity antideter-*  
647 *minants* that prevent or weaken interactions at specific locations, possibly by  
648 guiding and orienting interactions away from other interaction-determining  
649 features [34].

650 The notion of “matching” used in this work should not be taken literally.  
651 The essential feature of the XNOR rule is its provision of two ways to match  
652 (0/0 and 1/1), corresponding to the availability of alternative paired sets  
653 of features in biomacromolecules that promote assortative interactions. As-  
654 sortative interactions occur by means of both *complementarity* in the shapes  
655 and motions of cognate pairs of tRNA and aaRS species, and *identifiability* or  
656 distinctiveness in the shapes and motions of cognate and non-cognate pairs.  
657 Because of the symmetry of mutation that we assumed in this work, we could  
658 have equivalently named our landscape a “complementarity landscape” and  
659 obtained identical results using an XOR matching rule instead of the XNOR  
660 rule. It would then be simple, although vague and misleading, to interpret  
661 matching features as complementarily charged amino acid side-chains or com-  
662 plementary RNA nucleobases that interact directly. However, this would be  
663 oversimplified on multiple levels: first, because identifying features in tRNAs  
664 can depend only indirectly on underlying bases and residues through the  
665 overall shape and motion in what is called *indirect read-out* [117]; second,  
666 tRNAs are extensively post-transcriptionally modified, which also biochem-

667 ically integrates information from multiple sites in ways crucial for tRNA  
668 identity [118]; third, feature matching and mis-matching occurs in general  
669 through different sequence alphabets in RNA and proteins; and fourth, aaRS  
670 proteins are autocatalytically synthesized through the aaRS-tRNA network  
671 itself [119].

672 Thus, in the present work we analyzed only the simplest one of four  
673 increasingly complex variations on the general problem of evolution of a self-  
674 encoded aaRS-tRNA network. We define four connected notions to make  
675 our arguments: *description*, *self-description*, *self-encoded description*, and  
676 *self-encoded self-description*. By *description* we mean that when a mature,  
677 folded gene product evolves to complement the shape and motion of a fixed  
678 and unevolving ligand like a metabolite in order to specifically bind it, it “de-  
679 scribes” that metabolite. This notion of “description” depends on the com-  
680 plex genotype-phenotype maps of RNA and protein folding, and therefore can  
681 attain complex and emergent evolutionary dynamics [91, 120]. Nonetheless,  
682 by definition, descriptions are of evolutionarily fixed targets and therefore in-  
683 trinsically less rugged, with smaller neutral network size or degeneracy, than  
684 the match landscapes studied in the present work. We contend that evolving  
685 a description of an unevolvable metabolite ligand corresponds to discovering  
686 what might be called an *Easter egg* in sequence space. Under the assump-  
687 tion of symmetric mutation, the “all-ones” genotype studied in the Appendix  
688 of Sella [90] corresponds to selection to match any equivalent evolutionarily  
689 static Easter Egg in sequence space, of any arbitrary sequence neighborhood.

690 In the present work on the other hand, we analyzed the problem of *self-*  
691 *description*: specifically, we evolved co-inherited cognate tRNA-aaRS gene  
692 pairs to describe one another, so that their expressed products obtain com-  
693plementary and identifying shapes and motions with one another. More  
694generally, the notion of self-description represents the information acquired  
695in genes by natural selection about the shapes and motions of the prod-  
696ucts (or regulatory regions) of other genes (which correspond to “self” with  
697respect to the cell they are co-inherited in). During the evolutionary collec-  
698tivation of genes and gene products into genomes and cells hypothesized  
699by Woese and co-authors, genes acquired information via natural selection  
700about the shapes and motions of other gene products, in order to interact  
701specifically and/or conditionally with them. This *self-description* (or equiv-  
702alently *self-information*) is the epistatic “biological glue” that binds folded  
703macromolecules, cells and organisms together, enabling them to convert en-  
704ergy into work and execute complex emergent functions. Self-description

705 applies equally well to epistatic interactions within genes and gene-products,  
706 where it programs their folds, major modes of motion, and allosteric changes  
707 in shape and motion in response to changes in cellular state. The reason  
708 that these notions of self-description are all consistent is precisely because  
709 the self-descriptions of biological entities at multiple scales become integrated  
710 through major evolutionary transitions.

711 Bedian [119] also called what he modeled “self-description,” but he meant  
712 something entirely different: the mutual self-compatible encoding of a set of  
713 aaRS catalytic active sites capable of aminoacylating different amino acids  
714 onto distinct tRNAs, so that the collection of self-encoded aaRSs active sites  
715 can autocatalytically resynthesize themselves and each other. In our ter-  
716 minological framework, this is *self-encoded description*, because tRNAs are  
717 treated as fixed and unevolving targets, like amino acids. Bedian’s model,  
718 and subsequent extensions by Wills and co-workers, consider that these dif-  
719 ferent selectivities of different aaRSs depend on distinct sets of *critical sites*  
720 in each aaRS (where each critical site corresponds to one of our site-types). The  
721 distinct sets of critical sites of aaRSs may be thought of as multiple distinct  
722 Easter eggs in sequence space that all must be simultaneously discovered and  
723 compatibly mutually encoded for the network of aaRS active sites to nucleate.  
724 But aaRSs have both catalytic and tRNA-binding domains. Bedian, Wills  
725 and co-workers have so far not considered the problem of tRNA recognition  
726 by autocatalytically encoded aaRSs in their work, which generalizes what we  
727 studied here in what might be called *self-encoded self-description*. Full treat-  
728 ment of the problem, involving autocatalytically-encoded Easter eggs and  
729 Match Landscapes, is reserved for future investigations. Progress will allow  
730 a fuller investigation of even larger models to investigate the coevolution of  
731 genetic code and metabolism [121, 122].

732 We conjecture that our present results will hold for these more com-  
733 plex models. We offer an interpretation of “matching” for our present re-  
734 sults which applies to all of these more complex biological settings; namely,  
735 matching represents the self-information contained in self-descriptions, or the  
736 information contained in genes about the identifying shapes and motions of  
737 other co-inherited genes and gene products. Commensurately informative  
738 self-descriptions are expected to be *nearly neutral* with one another in the  
739 sense of [90] and references therein, and as shown for interaction interfaces  
740 previously [123]. The nearly-neutral evolution of interaction-determining fea-  
741 tures within a high-dimensional sequence space of equally fit solutions makes  
742 compensatory mutations much more likely than reversals. This explains both

743 why interaction-informative features can evolve and diverge even while under  
744 strong selection, and why interaction-determining features are phylogeneti-  
745 cally informative.

746 Our theory that macromolecular interactions are encoded through sets  
747 of complementary and identifying features extends the universal principle  
748 of heredity clarified by Watson and Crick, through which all possible ge-  
749 netic sequences may be replicated by virtue of complementarity [124]. The  
750 relativity of the notions of complementarity and identity in the definition  
751 of self-description implies that macromolecular interactions are governed by  
752 symbolic representations, as discussed by Maynard Smith [115]. That is,  
753 within the context of a specific cell, arbitrary molecular shapes and motions  
754 are symbolically associated with specific functions. The notion of symbolic  
755 association is defined not only by the absence of relationship between the  
756 form and meaning of signals [115], but also by its cryptographic nature, in  
757 that it requires coordinated information to decode signals correctly within a  
758 large space of equally unambiguously expressive alternatives.

759 The statistical Tree of Life became more tree-like after the Darwinian  
760 Transition precisely because through this transition, cells evolved languages  
761 of self-encoded descriptions and self-descriptions critical to their fitness as  
762 cells. These genetic and cellular languages are symbolic, cryptographic, open-  
763 endedly expressive, and increasingly constrained from changing by the in-  
764 creasingly complex corpus of descriptions and self-descriptions they encode.  
765 Since languages evolve in a statistically tree-like manner [125, 126], so did  
766 the advent of these cellular and genetic languages caused cells to evolve in a  
767 statistically tree-like manner. Furthermore, the large degeneracy of equiva-  
768 lent self-descriptions implies that such a language may be surprisingly easy  
769 to originate spontaneously, yet once originated, will be heavily constrained  
770 to change only in ancestrally compatible ways [12].

771 It is easy to imagine that macromolecular interaction codes, like lan-  
772 guages, evolve to be both expressive and unambiguous, that is, to encode  
773 more and more interactions in robust and error-tolerant (and ambiguity-  
774 reducing) ways. The coding theory analogy to the universality of replication  
775 by complementarity lies in the notion of *non-trivial perfect codes*. Perfect  
776 codes uniquely cover all of a finite sequence space with a maximum number  
777 of code-words spaced a minimum distance apart, so that every single pos-  
778 sible code-word can be received unambiguously and decoded correctly even  
779 after one or more symbols in the code-word were altered. While we expect  
780 biological codes to be generally far from perfect, the theory of perfect codes

781 may be a useful reference point from which to relax assumptions, and seems  
782 relevant to the stochastic setting of gene expression. In this context, it is  
783 of interest to note that surprisingly few varieties of small non-trivial perfect  
784 codes exist (where *non-trivial* means a code with more than one code word,  
785 not using every possible sequence as a codeword, nor the  $P$ -ary repetition  
786 code) [102]. For symbolic alphabets of prime power size, all non-trivial per-  
787 fect codes have codeword sizes, lengths, and minimum distance parameters  
788 equal to those of either Hamming Codes or Golay Codes [102, 127]. However,  
789 the Golay codes are too large to be relevant to the problem of perfect coding  
790 of 20 or fewer aaRS-tRNA cognate interactions. The RNA alphabet is of  
791 prime power size, namely four. The Hamming code  $\mathcal{H}_r(h)$  over an alphabet  
792 of size  $r$  with positive integer index parameter  $h$  has  $M = r^{n-h}$  codewords of  
793 length  $n = (r^h - 1)/(r - 1)$  and minimum Hamming distance between code-  
794 words of 3, allowing correction of single-symbol errors. It is of interest to  
795 note that  $\mathcal{H}_4(2)$  contains four codewords of size 5,  $\mathcal{H}_4(3)$  contains 16 code-  
796 words of size 21, and  $\mathcal{H}_4(4)$  contains 64 codewords of size 85. The  $\mathcal{H}_4(3)$   
797 perfect codeword length of 21 is surprisingly close to the size of a postulated  
798 primordial tRNA hairpin [24, 128, 129] with acceptor stem length of 7 and  
799 anticodon loop of length 7, while the  $\mathcal{H}_4(4)$  perfect codeword length of 85 is  
800 surprisingly close to the typical lengths of tRNAs today.

801 We can use our theory to roughly calculate the probability  $p(n, P, d, H)$   
802 that an aaRS-tRNA network will evolve  $P$  matching codewords of minimum  
803 distance  $d$  over an interface of length  $n$  in a system with  $M$  mutually dissim-  
804 ilar tRNA replicators and  $N$  mutually dissimilar aaRS ribozyme replicators  
805 (with  $P \leq M, N$ ), and aaRS-tRNA per-site background and target symbol-  
806 pair frequencies defined by the expected relative entropy  $H$ . Counting all  
807 possible pairs between tRNA and aaRS genes, and assuming that tRNAs  
808 have evolveable anticodons, this probability is

$$p(n, P, d, H) = \binom{M}{P} \binom{N}{P} \mathcal{N}_2(n, P, d) (1 - \exp(-E(n, P))), \quad (23)$$

809 where  $\mathcal{N}_2(n, P, d)$  is the number of binary codes of length  $n$ , size  $P$  and mini-  
810 mum Hamming distance  $d$ ,  $E(n, P) = kMNn^2 2^{-nPH}$  is the expected number  
811 of random sequences achieving normalized score  $nPH$  in a search space of  
812 size  $nM \times nN$ , from Karlin-Altschul theory (and in which  $k$  is a correction  
813 factor for edge effects) [130], and where the expected relative entropy per-site  
814  $H$  may be computed by enumerating over all pairs of RNA bases, assuming

815 a specific base composition common to all genes, and an expected target  
816 similarity corresponding to one or fewer errors per  $n$  symbols. Although  
817 a unique finite number of codes  $\mathcal{N}_2(n, P, d)$  exists over any finite sequence  
818 space, no expression for its value is known [127]. However this number must  
819 be much larger than the number of ways to choose  $P$  codewords from any  
820 Hamming-code of length  $n$  and size  $Q \geq P$ , provided Hamming codes of  
821 that length and size exist, because of the existence of a potentially large, yet  
822 unknown number of non-linear codes with Hamming parameters [127]. The  
823 number of distinct Hamming codes of length  $n$  over an alphabet of size  $q$  is  
824  $q!n!$  [127]. Further investigation is needed, but we believe that  $p(n, P, d, H)$   
825 may be surprisingly large.

826 The theory by which we computed stationary genotype distributions can  
827 incorporate up to three kinds of mutational asymmetry [88] such as GC-bias,  
828 transition bias, or transcription- or strand-dependent mutation, all relevant  
829 to problems in the evolution of the genetic code. It should be expected  
830 that incorporating asymmetric mutation will break symmetries in the fit-  
831 ness of genotypes and will change the expected composition of interaction-  
832 determining features.

833 An importantly unrealistic assumption in the present work is that of large  
834 aaRS concentrations in our macroscopic model of aminoacylation. The sto-  
835 chastic dynamics of cellular-scale aminoacylation coupled to the sink of trans-  
836 lating ribosomes is complex, exhibiting phenomena such as ultra-sensitivity [111].  
837 We have implemented a mesoscopic version of aminoacylation kinetics using  
838 Gillespie's direct method [131], results with which will be published else-  
839 where. Although our results do not depend on how translational rate is  
840 implemented, our model can fruitfully be integrated into a fully stochastic  
841 model of translation such as in Shah et al. [132]. In future work we will in-  
842 corporate these and other extensions into new models for the coevolution of  
843 genetically encoded descriptions and self-descriptions with codon meanings  
844 and metabolism in structured populations, to better understand evolution  
845 through the Darwinian Transition.

#### 846 4. Acknowledgments

847 DHA and AC-H were supported by the National Science Foundation  
848 (INSPIRE-1344279). DHA was supported by a Julius Kuhn Guest Profes-  
849 sorship at Martin Luther Universität in Halle-Wittenberg. Computational  
850 research in this report was performed on the MERCED HPC cluster sup-

851 ported by the National Science Foundation (ACI-1429783). Our funding  
 852 sources had no role in study design, analysis, interpretation, writing or sub-  
 853 mission. We thank Harish Bhat, Ivo Grosse, Suzanne Sindi, Kyle Kauffman,  
 854 Michael Frisch, and Emily Jane McTavish for valuable discussions.

855 **Appendix A. Decomposition of the steady state solution of fixed**  
 856 **genotypes with multiplicative fitness components**

857 **Remark 1.** Define  $\mathbb{V}$  to be the set of possible values in a site.  $\mathbb{V}$  could be  
 858 the set of nucleotides, the set of amino acids, etc. In this particular study,  
 859  $\mathbb{V} = \{0, 1\}$

860 **Remark 2.** Define  $G^{(M,N,n,p)}$  ( $M$  not necessarily unequal to  $N$ ) to be the set  
 861 of all possible genotypes of width  $n \in \mathbb{N}$  and  $pn$  sites per-gene, with  $p \in \mathbb{N}$ .  
 862 If  $\forall g \in G^{(M,N,n,p)}$  have length  $L$ , then  $|G^{(M,N,n,p)}| = |\mathbb{V}|^L$ .

863 **Remark 3.** Consider a genotype,  $g \in G^{(M,N,n,p)}$ . Let  $\mathbb{T}$  be the set of tRNA  
 864 genes in  $g$  with  $|\mathbb{T}| = M$ ,  $2 \leq M < \infty$  and let  $\mathbb{A}$  be the set of aaRS genes in  
 865  $g$  with  $|\mathbb{A}| = N$ ,  $2 \leq N < \infty$ . The lengths of genes  $t \in \mathbb{T}$  and  $a \in \mathbb{A}$  are all  
 866 equal to  $np \forall t, a$ . Let  $p(M + N) = L_b$ . Define **block**  $b_i^g \in \mathbb{V}^{L_b}$ ,  $i \in \{1, 2, \dots, n\}$   
 867 to be the sequence of  $p$  ordered values starting at the  $j^{th}$  site across all  $t$  and  
 868  $a$  genes in genotype  $g$ , with  $j = (i - 1)p$ . For a genotype  $g \in G^{(M,N,n,p)}$ , there  
 869 will be  $n$  blocks  $b_i^g$ , and each will be  $L_b$  long, it is possible for  $b_i^g = b_j^g$  for  
 870  $1 \leq i \neq j \leq n$ , and  $G^{(M,N,1,p)} = \mathbb{V}^{L_b}$  is the set of all possible types of blocks.

871 **Theorem 1.** Let  $w_g$  be the viability of genotype  $g \in G^{(M,N,1,p)}$ ,  $\mathcal{N}$  be the  
 872 population size, and  $\beta = \mathcal{N} - 1$  for the Moran process,  $\beta = 2(\mathcal{N} - 1)$  for  
 873 the haploid Wright-Fisher process, and  $\beta = 2\mathcal{N} - 1$  for the diploid Wright-  
 874 Fisher process. Given that the stationary frequency  $P_g^*$  of genotype  $g$  is  $P_g^* =$

875  $\frac{w_g^\beta}{\sum_{h \in G^{(M,N,1,p)}} w_h^\beta}$ , and that the viability  $W_\kappa$  is multiplicative across blocks in a

876 genotype  $\kappa \in G^{(M,N,n>1,p)}$  (i.e.  $W_\kappa = \prod_{i=1}^n w_{b_i^\kappa}$ ), then the stationary frequency

877  $P_\kappa^*$  of genotype  $\kappa$  is

$$P_\kappa^* = \prod_{i=1}^n P_{b_i^\kappa}^* \quad (\text{A.1})$$



878 PROOF (PROOF OF A.1).  $P_\kappa^* = \prod_{i=1}^n P_{b_i^\kappa}^*$

879 By definition,

$$P_\kappa^* = \frac{W_\kappa^\beta}{\sum_{\xi \in G^{(M,N,n,p)}} W_\xi^\beta} \quad (\text{A.2})$$

880 By the multiplicativity property this becomes

$$P_\kappa^* = \frac{\prod_{i=1}^n w_{b_i^\kappa}^\beta}{\sum_{\xi \in G^{(M,N,n,p)}} \prod_{j=1}^n w_{b_j^\xi}^\beta}. \quad (\text{A.3})$$

881 It needs to be shown that  $\frac{\prod_{i=1}^n w_{b_i^\kappa}^\beta}{\sum_{\xi \in G^{(M,N,n,p)}} \prod_{j=1}^n w_{b_j^\xi}^\beta} = \prod_{i=1}^n \frac{w_{b_i^\kappa}^\beta}{\sum_{g \in G^{(M,N,1,p)}} w_g^\beta}$ . Essen-

882 tially, the proof breaks down to whether  $\sum_{\xi \in G^{(M,N,n,p)}} \prod_{j=1}^n w_{b_j^\xi}^\beta = \left( \sum_{g \in G^{(M,N,1,p)}} w_g^\beta \right)^n$

883 Start with,

$$\sum_{\xi \in G^{(M,N,n,p)}} \prod_{j=1}^n w_{b_j^\xi}^\beta = \sum_{\xi \in G^{(M,N,n,p)}} w_{b_1^\xi}^\beta \cdot w_{b_2^\xi}^\beta \cdot \dots \cdot w_{b_n^\xi}^\beta. \quad (\text{A.4})$$

884 Since  $G^{(M,N,1,p)}$  is the set of all possible blocks,  $b_j^\xi$ , and no combination of  $L_b$   
 885 length genotypes across blocks is impossible, there are  $B^n$  possible sequences  
 886 for genotypes  $\xi \in G^{(M,N,n,p)}$ , where  $B = |G^{(M,N,1,p)}|$ . This is consistent with  
 887 the cardinality of  $G^{(M,N,n,p)}$  since  $L = L_b n$  and thus  $B^n = |\mathbb{V}|^L = |\mathbb{V}|^{L_b n} =$   
 888  $|G^{(M,N,1,p)}|^n$ . Since we are summing over all possible genotypes  $\xi \in G^{(M,N,n,p)}$ ,  
 889 and since different genotypes in  $G^{(M,N,n,p)}$  with the same blocks but in differ-  
 890 ent orders will have the same viability, then every viability term will be of  
 891 the form  $\binom{n}{n_{g_1}, n_{g_2}, \dots, n_{g_B}} w_{g_1}^{n_{g_1} \beta} w_{g_2}^{n_{g_2} \beta} \dots w_{g_B}^{n_{g_B} \beta}$  where each  $g_i \in G^{(M,N,1,p)}$  is (pos-  
 892 sibly arbitrarily) ordered from 1 to  $B$  and  $n_{g_i} \in \mathbb{W}$  is the number of blocks  
 893 of genotype  $\xi$  that are  $g_i$ . Since every genotype is represented, (A.4) is a  
 894 multinomial and can be rewritten  $\left( \sum_{g \in G^{(M,N,1,p)}} w_g^\beta \right)^n$ . If this were not the case  
 895 and one of the viability coefficients was less than the expected multinomial

896 coefficient, then that could only mean that at least one genotype was not  
897 being counted. If one had a coefficient larger than expected it would have to  
898 mean that at least one genotype was being counted more than once. There-  
899 fore to prove (A.1), plug this multinomial representation into (A.3),

$$900 \quad P_{\kappa}^* = \frac{\prod_{i=1}^n w_{b_i^{\kappa}}^{\beta}}{(\sum_{g \in G^{(M,N,1,p)}} w_g^{\beta})^n} = \prod_{i=1}^n \frac{w_{b_i^{\kappa}}^{\beta}}{\sum_{g \in G^{(M,N,1,p)}} w_g^{\beta}} = \prod_{i=1}^n P_{b_i^{\kappa}}^*$$
$$901 \quad \therefore P_{\kappa}^* = \prod_{i=1}^n P_{b_i^{\kappa}}^* \quad \square$$

- 902 [1] C. R. Woese, Interpreting the universal phylogenetic tree, Proceedings  
903 of the National Academy of Sciences 97 (2000) 8392–8396.
- 904 [2] C. R. Woese, G. J. Olsen, M. Ibba, D. Söll, Aminoacyl-tRNA syn-  
905 thetases, the genetic code, and the evolutionary process, Microbiology  
906 and molecular biology reviews: MMBR 64 (2000) 202–236.
- 907 [3] C. R. Woese, On the evolution of cells, Proceedings of the National  
908 Academy of Sciences 99 (2002) 8742–8747.
- 909 [4] K. Vetsigian, C. Woese, N. Goldenfeld, Collective evolution and the  
910 genetic code, Proceedings of the National Academy of Sciences 103  
911 (2006) 10696–10701.
- 912 [5] E. Roberts, A. Sethi, J. Montoya, C. R. Woese, Z. Luthey-Schulten,  
913 Molecular signatures of ribosomal evolution, Proceedings of the Na-  
914 tional Academy of Sciences of the United States of America 105 (2008)  
915 13953–13958.
- 916 [6] T. Ruusala, D. Andersson, M. Ehrenberg, C. G. Kurland, Hyper-  
917 accurate ribosomes inhibit growth., The EMBO Journal 3 (1984) 2575–  
918 2580.
- 919 [7] C. G. Kurland, Translational accuracy and the fitness of bacteria,  
920 Annual Review of Genetics 26 (1992) 29–50.
- 921 [8] W. Ran, P. G. Higgs, Contributions of Speed and Accuracy to Trans-  
922 lational Selection in Bacteria, PLOS ONE 7 (2012) e51652.

- 923 [9] L. Ribas de Pouplana, M. A. S. Santos, J.-H. Zhu, P. J. Farabaugh,  
924 B. Javid, Protein mistranslation: friend or foe?, *Trends in Biochemical*  
925 *Sciences* 39 (2014) 355–362.
- 926 [10] E. V. Koonin, The Biological Big Bang model for the major transitions  
927 in evolution, *Biology Direct* 2 (2007) 21.
- 928 [11] M. L. Katz, C. Shapiro, Systems Competition and Network Effects,  
929 *The Journal of Economic Perspectives* 8 (1994) 93–115.
- 930 [12] F. H. Crick, The origin of the genetic code, *Journal of Molecular*  
931 *Biology* 38 (1968) 367–379.
- 932 [13] E. Szathmry, J. M. Smith, The major evolutionary transitions, *Nature*  
933 374 (1995) 227–232.
- 934 [14] M. M. Desai, D. Weissman, M. W. Feldman, Evolution Can Favor  
935 Antagonistic Epistasis, *Genetics* 177 (2007) 1001–1010.
- 936 [15] U. Liberman, M. Feldman, On the evolution of epistasis III: the haploid  
937 case with mutation, *Theoretical Population Biology* 73 (2008) 307–316.
- 938 [16] E. V. Koonin, Horizontal gene transfer: essentiality and evolvability  
939 in prokaryotes, and roles in evolutionary transitions, *F1000Research* 5  
940 (2016).
- 941 [17] M. C. Rivera, R. Jain, J. E. Moore, J. A. Lake, Genomic evidence  
942 for two functionally distinct gene classes, *Proceedings of the National*  
943 *Academy of Sciences of the United States of America* 95 (1998) 6239–  
944 6244.
- 945 [18] R. Jain, M. C. Rivera, J. A. Lake, Horizontal gene transfer among  
946 genomes: the complexity hypothesis, *Proceedings of the National*  
947 *Academy of Sciences of the United States of America* 96 (1999) 3801–  
948 3806.
- 949 [19] H. B. Fraser, A. E. Hirsh, L. M. Steinmetz, C. Scharfe, M. W. Feldman,  
950 Evolutionary rate in the protein interaction network, *Science (New*  
951 *York, N.Y.)* 296 (2002) 750–752.

- 952 [20] C. S. Francklyn, E. A. First, J. J. Perona, Y.-M. Hou, Methods for  
953 kinetic and thermodynamic analysis of aminoacyl-tRNA synthetases,  
954 Methods (San Diego, Calif.) 44 (2008) 100–118.
- 955 [21] C. Francklyn, K. Musier-Forsyth, S. A. Martinis, Aminoacyl-tRNA  
956 synthetases in biology and disease: new evidence for structural and  
957 functional diversity in an ancient family of enzymes, RNA (New York,  
958 N.Y.) 3 (1997) 954–960.
- 959 [22] Y. I. Wolf, L. Aravind, N. V. Grishin, E. V. Koonin, Evolution of  
960 aminoacyl-tRNA synthetases—analysis of unique domain architectures  
961 and phylogenetic trees reveals a complex history of horizontal gene  
962 transfer events, Genome Research 9 (1999) 689–710.
- 963 [23] P. O’Donoghue, Z. Luthey-Schulten, On the evolution of structure  
964 in aminoacyl-tRNA synthetases, Microbiology and molecular biology  
965 reviews: MMBR 67 (2003) 550–573.
- 966 [24] P. Schimmel, R. Giegé, D. Moras, S. Yokoyama, An operational RNA  
967 code for amino acids and possible relationship to genetic code, Pro-  
968 ceedings of the National Academy of Sciences of the United States of  
969 America 90 (1993) 8763–8768.
- 970 [25] S. Shaul, D. Berel, Y. Benjamini, D. Graur, Revisiting the operational  
971 RNA code for amino acids: Ensemble attributes and their implications,  
972 RNA (New York, N.Y.) 16 (2010) 141–153.
- 973 [26] G. Eriani, M. Delarue, O. Poch, J. Gangloff, D. Moras, Partition of  
974 tRNA synthetases into two classes based on mutually exclusive sets of  
975 sequence motifs, Nature 347 (1990) 203–206.
- 976 [27] S. Cusack, Aminoacyl-tRNA synthetases, Current Opinion in Struc-  
977 tural Biology 7 (1997) 881–889.
- 978 [28] J. R. Brown, W. F. Doolittle, Root of the universal tree of life based  
979 on ancient aminoacyl-tRNA synthetase gene duplications, Proceedings  
980 of the National Academy of Sciences of the United States of America  
981 92 (1995) 2441–2445.

- 982 [29] L. Ribas de Pouplana, P. Schimmel, Two classes of tRNA synthetases  
983 suggested by sterically compatible dockings on tRNA acceptor stem,  
984 *Cell* 104 (2001) 191–193.
- 985 [30] L. Ribas de Pouplana, P. Schimmel, Aminoacyl-tRNA synthetases:  
986 potential markers of genetic code development, *Trends in Biochemical*  
987 *Sciences* 26 (2001) 591–596.
- 988 [31] D. S. Goodsell, S. Dutta, C. Zardecki, M. Voigt, H. M. Berman, S. K.  
989 Burley, The RCSB PDB Molecule of the Month: Inspiring a Molecular  
990 View of Biology, *PLOS Biology* 13 (2015) e1002140.
- 991 [32] M. Sprinzl, F. Grueter, A. Spelzhaus, D. H. Gauss, Compilation of  
992 tRNA sequences., *Nucleic Acids Research* 8 (1980) r1–r22.
- 993 [33] R. Giegé, E. Touzé, B. Lorber, A. Théobald-Dietrich, C. Sauter,  
994 Crystallogensis Trends of Free and Liganded Aminoacyl-tRNA Syn-  
995 thetases, *Crystal Growth & Design* 8 (2008) 4297–4306.
- 996 [34] R. Giegé, M. Sissler, C. Florentz, Universal rules and idiosyncratic  
997 features in tRNA identity, *Nucleic Acids Research* 26 (1998) 5017–  
998 5035.
- 999 [35] E. Freyhult, V. Moulton, D. H. Ardell, Visualizing bacterial tRNA  
1000 identity determinants and antideterminants using function logos and  
1001 inverse function logos, *Nucleic Acids Research* 34 (2006) 905–916.
- 1002 [36] D. H. Ardell, S. G. E. Andersson, TFAM detects co-evolution of tRNA  
1003 identity rules with lateral transfer of histidyl-tRNA synthetase, *Nucleic*  
1004 *Acids Research* 34 (2006) 893–904.
- 1005 [37] D. H. Ardell, G. Sella, On the evolution of redundancy in genetic codes,  
1006 *Journal of Molecular Evolution* 53 (2001) 269–281.
- 1007 [38] E. Freyhult, Y. Cui, O. Nilsson, D. H. Ardell, New computational  
1008 methods reveal tRNA identity element divergence between Proteobac-  
1009 teria and Cyanobacteria, *Biochimie* 89 (2007) 1276–1288.
- 1010 [39] K. C. H. Amrine, W. D. Swingley, D. H. Ardell, tRNA signatures re-  
1011 veal a polyphyletic origin of SAR11 strains among alphaproteobacteria,  
1012 *PLoS computational biology* 10 (2014) e1003454.

- 1013 [40] E. V. Koonin, A. S. Novozhilov, Origin and evolution of the genetic  
1014 code: the universal enigma, *Iubmb Life* 61 (2009) 99–111.
- 1015 [41] S. E. Massey, The neutral emergence of error minimized genetic codes  
1016 superior to the standard genetic code, *Journal of Theoretical Biology*  
1017 408 (2016) 237–242.
- 1018 [42] S. N. Rodin, S. Ohno, Two types of aminoacyl-trna synthetases could  
1019 be originally encoded by complementary strands of the same nucleic  
1020 ACID, *Origins of life and evolution of the biosphere* 25 (1995) 565–  
1021 589.
- 1022 [43] S. N. Rodin, A. S. Rodin, On the origin of the genetic code: signatures  
1023 of its primordial complementarity in tRNAs and aminoacyl-tRNA syn-  
1024 thetases, *Heredity* 100 (2008) 341–355.
- 1025 [44] C. W. Carter, L. Li, V. Weinreb, M. Collier, K. Gonzalez-Rivera,  
1026 M. Jimenez-Rodriguez, O. Erdogan, B. Kuhlman, X. Ambroggio,  
1027 T. Williams, S. N. Chandrasekharan, The Rodin-Ohno hypothesis that  
1028 two enzyme superfamilies descended from one ancestral gene: an un-  
1029 likely scenario for the origins of translation that will not be dismissed,  
1030 *Biology Direct* 9 (2014) 11.
- 1031 [45] J. A. G. M. de Visser, J. Hermisson, G. P. Wagner, L. Ancel Meyers,  
1032 H. Bagheri-Chaichian, J. L. Blanchard, L. Chao, J. M. Cheverud, S. F.  
1033 Elena, W. Fontana, G. Gibson, T. F. Hansen, D. Krakauer, R. C.  
1034 Lewontin, C. Ofria, S. H. Rice, G. von Dassow, A. Wagner, M. C.  
1035 Whitlock, Perspective: Evolution and detection of genetic robustness,  
1036 *Evolution; International Journal of Organic Evolution* 57 (2003) 1959–  
1037 1972.
- 1038 [46] C. O. Wilke, J. L. Wang, C. Ofria, R. E. Lenski, C. Adami, Evolution  
1039 of digital organisms at high mutation rates leads to survival of the  
1040 flattest, *Nature* 412 (2001) 331–333.
- 1041 [47] O. G. Berg, P. H. von Hippel, Selection of DNA binding sites by  
1042 regulatory proteins. Statistical-mechanical theory and application to  
1043 operators and promoters, *Journal of Molecular Biology* 193 (1987)  
1044 723–750.

- 1045 [48] U. Gerland, J. D. Moroz, T. Hwa, Physical constraints and functional  
1046 characteristics of transcription factor-DNA interaction, *Proceedings of*  
1047 *the National Academy of Sciences of the United States of America* 99  
1048 (2002) 12015–12020.
- 1049 [49] J. Berg, M. Lässig, A. Wagner, Structure and evolution of protein  
1050 interaction networks: a statistical model for link dynamics and gene  
1051 duplications, *BMC Evolutionary Biology* 4 (2004) 51.
- 1052 [50] V. Mustonen, J. Kinney, C. G. Callan, M. Lässig, Energy-dependent  
1053 fitness: A quantitative model for the evolution of yeast transcription  
1054 factor binding sites, *Proceedings of the National Academy of Sciences*  
1055 *of the United States of America* 105 (2008) 12376–12381.
- 1056 [51] A. Y. Tulchinsky, N. A. Johnson, W. B. Watt, A. H. Porter, Hybrid  
1057 incompatibility arises in a sequence-based bioenergetic model of tran-  
1058 scription factor binding, *Genetics* 198 (2014) 1155–1166.
- 1059 [52] T. Friedlander, R. Prizak, C. C. Guet, N. H. Barton, G. Tkačik, In-  
1060 trinsic limits to gene regulation by global crosstalk, *Nature Communi-*  
1061 *cations* 7 (2016) 12307.
- 1062 [53] E. J. Deeds, O. Ashenberg, J. Gerardin, E. I. Shakhnovich, Robust pro-  
1063 teinprotein interactions in crowded cellular environments, *Proceedings*  
1064 *of the National Academy of Sciences* 104 (2007) 14952–14957.
- 1065 [54] M. E. Johnson, G. Hummer, Nonspecific binding limits the number of  
1066 proteins in a cell and shapes their interaction networks, *Proceedings*  
1067 *of the National Academy of Sciences of the United States of America*  
1068 108 (2011) 603–608.
- 1069 [55] M. E. Johnson, G. Hummer, Interface-Resolved Network of Protein-  
1070 Protein Interactions, *PLoS Computational Biology* 9 (2013).
- 1071 [56] J. J. Hopfield, Kinetic Proofreading: A New Mechanism for Reducing  
1072 Errors in Biosynthetic Processes Requiring High Specificity, *Proceed-*  
1073 *ings of the National Academy of Sciences* 71 (1974) 4135–4139.
- 1074 [57] J. Ninio, Kinetic amplification of enzyme discrimination, *Biochimie* 57  
1075 (1975) 587–595.

- 1076 [58] K.-W. Leong, U. Uzun, M. Selmer, M. Ehrenberg, Two proofreading  
1077 steps amplify the accuracy of genetic code translation, *Proceedings of*  
1078 *the National Academy of Sciences of the United States of America* 113  
1079 (2016) 13744–13749.
- 1080 [59] T. Yamane, J. J. Hopfield, Experimental evidence for kinetic proof-  
1081 reading in the aminoacylation of tRNA by synthetase., *Proceedings of*  
1082 *the National Academy of Sciences of the United States of America* 74  
1083 (1977) 2246–2250.
- 1084 [60] H. Mellenius, M. Ehrenberg, Transcriptional accuracy modeling sug-  
1085 gests two-step proofreading by RNA polymerase, *Nucleic Acids Re-*  
1086 *search* 45 (2017) 11582–11593.
- 1087 [61] H. Qian, Reducing Intrinsic Biochemical Noise in Cells and Its Ther-  
1088 modynamic Limit, *Journal of Molecular Biology* 362 (2006) 387–392.
- 1089 [62] W. S. Hlavacek, A. Redondo, H. Metzger, C. Wofsy, B. Goldstein,  
1090 Kinetic proofreading models for cell signaling predict ways to escape  
1091 kinetic proofreading, *Proceedings of the National Academy of Sciences*  
1092 *of the United States of America* 98 (2001) 7295–7300.
- 1093 [63] T. W. McKeithan, Kinetic proofreading in T-cell receptor signal trans-  
1094 duction, *Proceedings of the National Academy of Sciences* 92 (1995)  
1095 5042–5046.
- 1096 [64] V. Barone, M. Lang, S. F. G. Krens, S. J. Pradhan, S. Shamipour,  
1097 K. Sako, M. Sikora, C. C. Guet, C.-P. Heisenberg, An Effective Feed-  
1098 back Loop between Cell-Cell Contact Duration and Morphogen Signal-  
1099 ing Determines Cell Fate, *Developmental Cell* 43 (2017) 198–211.e12.
- 1100 [65] A. Murugan, D. A. Huse, S. Leibler, Speed, dissipation, and error in  
1101 kinetic proofreading, *Proceedings of the National Academy of Sciences*  
1102 109 (2012) 12034–12039.
- 1103 [66] K. Banerjee, A. B. Kolomeisky, O. A. Igoshin, Elucidating interplay  
1104 of speed and accuracy in biological error correction, *Proceedings of*  
1105 *the National Academy of Sciences of the United States of America* 114  
1106 (2017) 5183–5188.



- 1107 [67] M. Ehrenberg, C. Blomberg, Thermodynamic constraints on kinetic  
1108 proofreading in biosynthetic pathways, *Biophysical Journal* 31 (1980)  
1109 333–358.
- 1110 [68] C. G. Kurland, The role of guanine nucleotides in protein biosynthesis.,  
1111 *Biophysical Journal* 22 (1978) 373–392.
- 1112 [69] C.-M. Zhang, J. J. Perona, K. Ryu, C. Francklyn, Y.-M. Hou, Distinct  
1113 kinetic mechanisms of the two classes of Aminoacyl-tRNA synthetases,  
1114 *Journal of Molecular Biology* 361 (2006) 300–311.
- 1115 [70] P. Sartori, S. Pigolotti, Kinetic versus Energetic Discrimination in  
1116 Biological Copying, *Physical Review Letters* 110 (2013) 188101.
- 1117 [71] Wright, S., The roles of mutation, inbreeding, crossbreeding and se-  
1118 lection in evolution., *Proc. 6th Int. Congress on Genetics, Ithaca, NY,*  
1119 *USA* 1 (1932) 356–366.
- 1120 [72] S. E. Ahnert, Structural properties of genotype-phenotype maps, *Journal*  
1121 *of the Royal Society Interface* 14 (2017).
- 1122 [73] Dykhuizen, D, Recommendation of [Lunzer M et al., *Science* 2005 ,  
1123 310(5747):499-501], 2005.
- 1124 [74] K. Crona, D. Greene, M. Barlow, The peaks and geometry of fitness  
1125 landscapes, *Journal of Theoretical Biology* 317 (2013) 1–10.
- 1126 [75] F. J. Poelwijk, S. Tnase-Nicola, D. J. Kiviet, S. J. Tans, Reciprocal sign  
1127 epistasis is a necessary condition for multi-peaked fitness landscapes,  
1128 *Journal of Theoretical Biology* 272 (2011) 141–144.
- 1129 [76] Kauffman, S.A., *The Origins of Order. Self-Organization and Selection*  
1130 *in Evolution.*, Oxford University Press, Oxford, U.K., 1993.
- 1131 [77] S. Kauffman, S. Levin, Towards a general theory of adaptive walks on  
1132 rugged landscapes, *Journal of Theoretical Biology* 128 (1987) 11–45.
- 1133 [78] M. L. Siegal, A. Bergman, Waddington’s canalization revisited: Devel-  
1134 opmental stability and evolution, *Proceedings of the National Academy*  
1135 *of Sciences* 99 (2002) 10528–10532.

- 1136 [79] T. MacCarthy, A. Bergman, Coevolution of robustness, epistasis, and  
1137 recombination favors asexual reproduction, *Proceedings of the National*  
1138 *Academy of Sciences* 104 (2007) 12801–12806.
- 1139 [80] A. Orlenko, P. B. Chi, D. A. Liberles, Characterizing the roles of  
1140 changing population size and selection on the evolution of flux control  
1141 in metabolic pathways, *BMC evolutionary biology* 17 (2017) 117.
- 1142 [81] M. Nei, Modification of Linkage Intensity by Natural Selection, *Ge-*  
1143 *netics* 57 (1967) 625–641.
- 1144 [82] Feldman, M. W., Selection for linkage modification: I. Random mating  
1145 populations., *Theoretical Population Biology* 3 (1972) 324–346.
- 1146 [83] L. Altenberg, U. Liberman, M. W. Feldman, Unified reduction princi-  
1147 ple for the evolution of mutation, migration, and recombination, *Pro-*  
1148 *ceedings of the National Academy of Sciences* 114 (2017) E2392–E2400.
- 1149 [84] C. O. Wilke, C. Adami, Interaction between directional epistasis and  
1150 average mutational effects, *Proceedings of the Royal Society of London*  
1151 *B: Biological Sciences* 268 (2001) 1469–1474.
- 1152 [85] P.-A. Gros, H. L. Nagard, O. Tenaillon, The Evolution of Epistasis  
1153 and Its Links With Genetic Robustness, Complexity and Drift in a  
1154 Phenotypic Model of Adaptation, *Genetics* 182 (2009) 277–293.
- 1155 [86] D. M. Weinreich, R. A. Watson, L. Chao, Perspective: Sign epistasis  
1156 and genetic constraint on evolutionary trajectories, *Evolution; Inter-*  
1157 *national Journal of Organic Evolution* 59 (2005) 1165–1174.
- 1158 [87] D. B. Weissman, M. M. Desai, D. S. Fisher, M. W. Feldman, The  
1159 rate at which asexual populations cross fitness valleys, *Theoretical*  
1160 *Population Biology* 75 (2009) 286–300.
- 1161 [88] G. Sella, A. E. Hirsh, The application of statistical physics to evolu-  
1162 tionary biology, *Proceedings of the National Academy of Sciences* 102  
1163 (2005) 9541–9546.
- 1164 [89] D. M. McCandlish, A. Stoltzfus, Modeling evolution using the proba-  
1165 bility of fixation: history and implications, *The Quarterly Review of*  
1166 *Biology* 89 (2014) 225–252.

- 1167 [90] G. Sella, An exact steady state solution of Fisher's geometric model  
1168 and other models, *Theoretical population biology* 75 (2008) 30–4.
- 1169 [91] D. D. Pollock, G. Thiltgen, R. A. Goldstein, Amino acid coevolu-  
1170 tion induces an evolutionary Stokes shift, *Proceedings of the National*  
1171 *Academy of Sciences of the United States of America* 109 (2012) E1352–  
1172 1359.
- 1173 [92] G. Sella, D. H. Ardell, The impact of message mutation on the fitness  
1174 of a genetic code, *Journal of Molecular Evolution* 54 (2002) 638–651.
- 1175 [93] P. A. P. Moran, The survival of a mutant gene under selection. II,  
1176 *Journal of the Australian Mathematical Society* 1 (1960) 485–491.
- 1177 [94] P. R. Schimmel, D. Söll, Aminoacyl-tRNA synthetases: general fea-  
1178 tures and recognition of transfer RNAs, *Annual Review of Biochemistry*  
1179 48 (1979) 601–648.
- 1180 [95] R. W. Hamming, Error detecting and error correcting codes, *The Bell*  
1181 *System Technical Journal* 29 (1950) 147–160.
- 1182 [96] R. Rigler, U. Pachmann, R. Hirsch, H. G. Zachau, On the interaction of  
1183 seryl-tRNA synthetase with tRNA Ser. A contribution to the problem  
1184 of synthetase-tRNA recognition, *European Journal of Biochemistry* 65  
1185 (1976) 307–315.
- 1186 [97] D. Riesner, A. Pingoud, D. Boehme, F. Peters, G. Maass, Distinct  
1187 steps in the specific binding of tRNA to aminoacyl-tRNA synthetase.  
1188 Temperature-jump studies on the serine-specific system from yeast and  
1189 the tyrosine-specific system from *Escherichia coli*, *European Journal*  
1190 *of Biochemistry* 68 (1976) 71–80.
- 1191 [98] I. A. Vasil'eva, N. A. Moor, Interaction of aminoacyl-tRNA synthetases  
1192 with tRNA: general principles and distinguishing characteristics of the  
1193 high-molecular-weight substrate recognition, *Biochemistry. Biokhimiia*  
1194 72 (2007) 247–263.
- 1195 [99] M. A. Savageau, R. R. Freter, On the evolution of accuracy and cost  
1196 of proofreading tRNA aminoacylation, *Proceedings of the National*  
1197 *Academy of Sciences of the United States of America* 76 (1979) 4507–  
1198 4510.

- 1199 [100] D. H. Ardell, G. Sella, No accident: genetic codes freeze in error-  
1200 correcting patterns of the standard genetic code, *Philosophical Trans-*  
1201 *actions of the Royal Society of London. Series B, Biological Sciences*  
1202 357 (2002) 1625–1642.
- 1203 [101] G. Sella, D. H. Ardell, The coevolution of genes and genetic codes:  
1204 Crick’s frozen accident revisited, *Journal of Molecular Evolution* 63  
1205 (2006) 297–313.
- 1206 [102] A. Tietäväinen, On the Nonexistence of Perfect Codes over Finite  
1207 Fields, *SIAM Journal on Applied Mathematics* 24 (1973) 88–96.
- 1208 [103] F. I. Solov’eva, Perfect binary codes: bounds and properties, *Discrete*  
1209 *Mathematics* 213 (2000) 283–290.
- 1210 [104] H. Helgert, R. Stinaff, Minimum-distance bounds for binary linear  
1211 codes, *IEEE Transactions on Information Theory* 19 (1973) 344–356.
- 1212 [105] M. Best, A. Brouwer, F. MacWilliams, A. Odlyzko, N. Sloane, Bounds  
1213 for binary codes of length less than 25, *IEEE Transactions on Infor-*  
1214 *mation Theory* 24 (1978) 81–93.
- 1215 [106] M. J. Rogers, D. Söll, Inaccuracy and the recognition of tRNA,  
1216 *Progress in Nucleic Acid Research and Molecular Biology* 39 (1990)  
1217 185–208.
- 1218 [107] F. Solov’eva, On perfect binary codes, *Discrete Applied Mathematics*  
1219 156 (2008) 1488–1498.
- 1220 [108] W. Paulander, S. Maisnier-Patin, D. I. Andersson, Multiple mech-  
1221 anisms to ameliorate the fitness burden of mupirocin resistance in  
1222 *Salmonella typhimurium*, *Molecular Microbiology* 64 (2007) 1038–1048.
- 1223 [109] M. Ehrenberg, C. G. Kurland, Costs of accuracy determined by a  
1224 maximal growth rate constraint, *Quarterly Reviews of Biophysics* 17  
1225 (1984) 45–82.
- 1226 [110] J. Elf, D. Nilsson, T. Tenson, M. Ehrenberg, Selective Charging of  
1227 tRNA Isoacceptors Explains Patterns of Codon Usage, *Science* 300  
1228 (2003) 1718–1722.

- 1229 [111] J. Elf, M. Ehrenberg, Near-Critical Behavior of Aminoacyl-tRNA Pools  
1230 in *E. coli* at Rate-Limiting Supply of Amino Acids, *Biophysical Journal*  
1231 88 (2005) 132–146.
- 1232 [112] M. Johansson, M. Lovmar, M. Ehrenberg, Rate and accuracy of bac-  
1233 terial protein synthesis revisited, *Current Opinion in Microbiology* 11  
1234 (2008) 141–147.
- 1235 [113] M. Kimura, Natural selection as the process of accumulating genetic  
1236 information in adaptive evolution, *Genetics Research* 2 (1961) 127–140.
- 1237 [114] J. Felsenstein, On the Biological Significance of the Cost of Gene  
1238 Substitution, *The American Naturalist* 105 (1971) 1–11.
- 1239 [115] J. Maynard Smith, The 1999 Crafoord Prize Lectures. The idea of  
1240 information in biology, *The Quarterly Review of Biology* 74 (1999)  
1241 395–400.
- 1242 [116] M. C. Donaldson-Matasci, C. T. Bergstrom, M. Lachmann, The fitness  
1243 value of information, *Oikos* (Copenhagen, Denmark) 119 (2010) 219–  
1244 230.
- 1245 [117] J. J. Perona, Y.-M. Hou, Indirect Readout of tRNA for Aminoacyla-  
1246 tion, *Biochemistry* 46 (2007) 10419–10432.
- 1247 [118] E. M. Novoa, M. Pavon-Eternod, T. Pan, L. Ribas de Pouplana, A  
1248 Role for tRNA Modifications in Genome Structure and Codon Usage,  
1249 *Cell* 149 (2012) 202–213.
- 1250 [119] V. Bedian, Self-description and the origin of the genetic code, *Bio*  
1251 *Systems* 60 (2001) 39–47.
- 1252 [120] M. C. Cowperthwaite, L. A. Meyers, How Mutational Networks Shape  
1253 Evolution: Lessons from RNA Models, *Annual Review of Ecology,*  
1254 *Evolution, and Systematics* 38 (2007) 203–230.
- 1255 [121] J. T.-F. Wong, Coevolution theory of the genetic code at age thirty,  
1256 *BioEssays: News and Reviews in Molecular, Cellular and Developmen-*  
1257 *tal Biology* 27 (2005) 416–425.

- 1258 [122] M. Di Giulio, An Autotrophic Origin for the Coded Amino Acids is  
1259 Concordant with the Coevolution Theory of the Genetic Code, *Journal*  
1260 *of Molecular Evolution* 83 (2016) 93–96.
- 1261 [123] M. Lynch, K. Hagner, Evolutionary meandering of intermolecular in-  
1262 teractions along the drift barrier, *Proceedings of the National Academy*  
1263 *of Sciences of the United States of America* 112 (2015) E30–38.
- 1264 [124] J. D. Watson, F. H. Crick, Genetical implications of the structure of  
1265 deoxyribonucleic acid, *Nature* 171 (1953) 964–967.
- 1266 [125] L. L. Cavalli-Sforza, Genes, peoples, and languages, *Proceedings of*  
1267 *the National Academy of Sciences of the United States of America* 94  
1268 (1997) 7719–7724.
- 1269 [126] A. Bouchard-Ct, D. Hall, T. L. Griffiths, D. Klein, Automated re-  
1270 construction of ancient languages using probabilistic models of sound  
1271 change, *Proceedings of the National Academy of Sciences of the United*  
1272 *States of America* 110 (2013) 4224–4229.
- 1273 [127] S. Roman, *Coding and Information Theory*, Graduate Texts in Math-  
1274 ematics, Springer-Verlag, 1992.
- 1275 [128] M. Di Giulio, Was it an ancient gene codifying for a hairpin RNA  
1276 that, by means of direct duplication, gave rise to the primitive tRNA  
1277 molecule?, *Journal of Theoretical Biology* 177 (1995) 95–101.
- 1278 [129] J. Widmann, M. Di Giulio, M. Yarus, R. Knight, tRNA creation by  
1279 hairpin duplication, *Journal of Molecular Evolution* 61 (2005) 524–530.
- 1280 [130] S. Karlin, S. F. Altschul, Methods for assessing the statistical signifi-  
1281 cance of molecular sequence features by using general scoring schemes,  
1282 *Proceedings of the National Academy of Sciences of the United States*  
1283 *of America* 87 (1990) 2264–2268.
- 1284 [131] D. T. Gillespie, Exact stochastic simulation of coupled chemical reac-  
1285 tions, *The Journal of Physical Chemistry* 81 (1977) 2340–2361.
- 1286 [132] P. Shah, Y. Ding, M. Niemczyk, G. Kudla, J. B. Plotkin, Rate-Limiting  
1287 Steps in Yeast Protein Translation, *Cell* 153 (2013) 1589–1601.

Table 1: Symbols, parameter values and references for the present work.

Symbol	Meaning	Value (Reference)
$P$	# of cognate tRNA-aaRS gene/species pairs	
$M$	# tRNA genes/species when $M \neq N$	
$N$	# aaRS genes/species when $M \neq N$	
$n$	width of interaction interface in site-blocks	
$k_d^{nc}$	dissociation rate constant with 0 matches	$10\,000\text{ s}^{-1}$
$k_d^c$	dissociation rate constant with $k$ matches	$220\text{ s}^{-1}$ [94]
$k$	# matches to diminish from $k_d^{nc}$ to $k_d^c$ , $1 \leq k \leq n$	
$G^{(M,N,n,p)}$	genotype space with $p$ sites per-gene per-site-block	
$g, g_H, g_M$	genomes $g \in G^{(P,P,n,p)}$	
$t_i$	tRNA gene/species, $1 \leq i \leq P$ or $i \in \{0, 1\}$	
$a_j$	aaRS gene/species, $1 \leq j \leq P$ or $j \in \{0, 1\}$	
$t_{ir}, a_{jr}$	state-bits of $t_i$ and $a_j$ , $1 \leq r \leq n$	
$m_{ir}, n_{jr}$	mask-bits of $t_{ir}$ and $a_{jr}$	
$L$	genome length, $L \in \{2Pn, 4Pn\}$	
$R$	matching rule, $R \in \{\text{XNOR}, \text{AND}, \text{AND-XNOR}\}$	
$m_{i,j}, m_{i,j}^R$	number of matches between $t_i$ and $a_j$ with rule $R$	
$k_d^{i,j}$	dissociation rate constant of $t_i$ and $a_j$	
$c_{\max}(t_i \longrightarrow a_j)$	maximal decoding probability	
$c_i$	codons, $1 \leq i \leq P$ or $i \in \{0, 1\}$	
$aa_j$	amino acids, $1 \leq j \leq P$ or $j \in \{0, 1\}$	
$s_l$	site-types, $1 \leq l \leq P$ or $l \in \{0, 1\}$	
$c(aa_j c_i)$	decoding probability	
$\gamma$	steady-state free energy of the cell	[61]
$\iota$	free energy of a match (viz. $\epsilon$ in [47, 54])	
$\theta$	dissociation rate constant ratio	
$\phi$	max. missense fitness cost per site-type	[37]
$\psi$	nonsense fitness cost per site-type	
$\delta$	ambiguity fitness cost per site-type	
$w$	viability fitness factor or term	
$\epsilon$	multiplicative epistasis per site-block	[15]
$\beta$	size of a haploid Moran population minus 1	[90, 93]
$d_h(\cdot, \cdot)$	Hamming distance	[95]
$H(\cdot, \cdot), H_{k=1}^N \cdot$	harmonic average	

Table 2: Viabilities of the symmetric multiplicative binary interaction channel with one site-block under two different matching rules.

Genotype ( $t_0a_0t_1a_1$ )	XNOR <sup>a</sup>	AND <sup>a</sup>	Viabilities <sup>b</sup>	
	( $\Leftrightarrow$ )	( $\wedge$ )	$w_{\Leftrightarrow}$	$w_{\wedge}$
0000	⊠	::	$\delta^2$	$\psi^2$
0001	⋈	::	$\phi$	$\psi^2$
0010	⋈	::	$\psi\delta$	$\psi^2$
0011	∥∥	∣	1	$\psi$
0100	∧	::	$\phi$	$\psi^2$
0101	::	::	$\psi^2$	$\psi^2$
0110	×	∖	$\phi^2$	$\psi\phi$
0111	∨	∨	$\psi\delta$	$\psi\delta$
1000	∨	::	$\psi\delta$	$\psi^2$
1001	×	/	$\phi^2$	$\psi\delta$
1010	::	::	$\psi^2$	$\psi^2$
1011	∧	∧	$\phi$	$\phi$
1100	∥∥	∣	1	$\psi$
1101	∨	∨	$\psi\delta$	$\psi\delta$
1110	⋈	⋈	$\phi$	$\phi$
1111	⊠	⊠	$\delta^2$	$\delta^2$

<sup>a</sup> Iconic representation of network phenotypes expressed for each genotype with each rule.

<sup>b</sup>  $0 < \psi < \phi < \delta = (\phi + 1)/2 < 1$ .