

Neutral models of *de novo* gene emergence suggest that gene evolution has a preferred trajectory

Bharat Ravi Iyengar^{1,†}, Erich Bornberg-Bauer^{1,2}

¹Institute for Evolution and Biodiversity, University of Münster,
Hüfferstrasse 1, 48149 Münster, Germany

²Department of Protein Evolution, Max Planck Institute for Biology Tübingen,
Max-Planck-Ring 5, 72076 Tübingen, Germany

† Corresponding author: b.ravi@uni-muenster.de

1 Abstract

2 New protein coding genes can emerge from genomic regions that previously did not
3 contain any genes, via a process called *de novo* gene emergence. To synthesize a protein,
4 DNA must be transcribed as well as translated. Both processes need certain DNA se-
5 quence features. Stable transcription requires promoters and a polydenylation signal,
6 while translation requires at least an open reading frame (ORF). We develop mathe-
7 matical models based on mutation probabilities, and the assumption of neutral evolu-
8 tion, to find out how quickly genes emerge and are lost. We also investigate the effect
9 of the order by which DNA features evolve, and if sequence composition is biased by
10 mutation rate. We rationalize how genes are lost much more rapidly than they emerge,
11 and how genes with long ORFs preferentially arise in regions that are already tran-
12 scribed. Our study not only answers some fundamental questions on the topic of *de*
13 *nov* emergence but also provides a modeling framework for future studies.

14 Introduction

15 Organisms evolve new traits by expanding their functional genome. Evolution of new
16 genes is one of the ways by which new traits can emerge. The definition of a gene
17 is complicated, and has been changing constantly (Gerstein *et al.*, 2007). We use the
18 following working definition of a gene: a gene is a region of the genome that gives rise
19 to a functional product, that is an RNA or a protein. The next big challenge is to define
20 what a function is, and how to classify gene functions. One of the easiest definition of
21 a gene product's function is its ability to improve the survival of the organism under
22 one or more environmental conditions (Keeling *et al.*, 2019). Proteins perform the most
23 diverse kinds of molecular functions ranging from catalysis of biochemical reactions
24 to formation of cellular structures (Berg *et al.*, 2002), whereas RNAs that do not encode
25 proteins like mRNAs, or participate in protein synthesis like rRNAs and tRNAs, are
26 mostly involved in regulation of gene expression (Iyengar *et al.*, 2014; Statello *et al.*,
27 2021).

28 In this study we focus on the evolution of genes that encode new proteins. New protein
29 coding genes frequently arise via duplication of existing protein coding genes. These
30 duplicated genes then genetically and functionally diverge from their parent genes by
31 accumulating mutations (Long *et al.*, 2003; Näsvalld *et al.*, 2012). Recently it has been
32 shown that new protein coding genes can arise independent of gene duplication, from
33 genomic sequences that did not previously encode any protein. This phenomenon
34 is called *de novo* gene emergence (Tautz and Domazet-Lošo, 2011; Zhao *et al.*, 2014;
35 Schmitz and Bornberg-Bauer, 2017; Vakirlis *et al.*, 2017; Van Oss and Carvunis, 2019).
36 *De novo* protein coding genes can emerge from intergenic sequences, non-coding RNA
37 genes, introns, and even regions that partially overlap with existing protein coding
38 genes (Tautz and Domazet-Lošo, 2011; Zhao *et al.*, 2014; Vakirlis *et al.*, 2017; Van Oss
39 and Carvunis, 2019; Prabh and Rödelsperger, 2019). To express a protein, a genomic
40 sequence should be transcribed as well as translated. Therefore, the gene needs se-
41 quence features that enable both these processes.

42 The primary requirement for transcription is recruitment of RNA polymerase. This
43 process is facilitated by DNA sequences called (Lenhard *et al.*, 2012). A part of the pro-
44 moter, called the core promoter is the region that determines the start of transcription
45 (Haberle and Stark, 2018). Transcription can also be initiated without the requirement
46 of a defined promoter sequence, and this process can occur throughout the genome
47 Clark *et al.* (2011). Even though pervasive transcription is widespread in eukaryotic
48 genomes, most RNAs are quickly eliminated by RNA degrading enzymes (Schmid and

49 Jensen, 2018). Most eukaryotic mRNAs and long non-coding RNAs, are polyadeny-
50 lated at their 3' termini, which makes them resistant to degradation. Polyadenylation,
51 which also marks the end of transcription (Richard and Manley, 2009), is facilitated
52 by sequences known as polyadenylation signals (Proudfoot, 2011). In prokaryotes, the
53 end of transcription is determined by sequences known as terminators (Santangelo
54 and Artsimovitch, 2011). Termination at an appropriate site is another fundamental
55 requirement for a productive transcription.

56 The second major requirement for protein expression is translation of the transcribed
57 mRNA. The most fundamental requirement for translation is the presence of an open
58 reading frame (ORF). Usually, an mRNA needs additional features to initiate protein
59 synthesis. These features include ribosome binding sites in prokaryotes (Omotajo *et al.*,
60 2015), and Kozak consensus sequences in eukaryotes (Kozak, 1986; Acevedo *et al.*, 2018;
61 Noderer *et al.*, 2014).

62 Although gain of transcription and translation features guarantees *de novo* birth of a
63 protein coding gene, it does not ensure that the gene would persist in the genome for
64 many generations. The newly born gene can lose the features as easily as it gained
65 them unless it has been fixed in the genome, for example via evolutionary selection.
66 Specifically, if the protein synthesized by the *de novo* gene provides a fitness advantage
67 to the host organism, it will undergo positive selection (Keeling *et al.*, 2019). Simultane-
68 ously, the protein should have low toxicity and cost of synthesis, to survive purifying
69 selection. A common mechanism by which protein mediated toxicity occurs is by mis-
70 folding and aggregation of proteins (Bucciantini *et al.*, 2002; Hartl, 2017).

71 In this study we develop probabilistic models of *de novo* gene emergence in eukary-
72 otes. We focus on *de novo* genes emerging from non-genic sequences, also known as
73 "proto-genes" (Carvunis *et al.*, 2012; Van Oss and Carvunis, 2019). We model *de novo*
74 gene emergence as a two step process. In the first step, a non-genic DNA sequence
75 gains transcriptional features that allows it to express an untranslated (non-coding)
76 RNA. In the second step this non-coding RNA gene acquires translational features
77 that allows it to express a protein. Alternatively, translational features can emerge in
78 the non-genic DNA before transcriptional features. Specifically, we calculate the prob-
79 abilities by which polyadenylation signals and ORFs emerge, and are lost, based on
80 the rates at which different kinds of DNA mutations naturally occur (mutation bias).
81 With these probabilities we estimate whether transcriptional and translational features
82 preferentially evolve in a specific order. Using a similar approach, we calculate how
83 random mutations affect protein composition. In our models, we assume that the pro-
84 teins expressed from these proto-genes provide no fitness advantage and are not toxic

85 to the host organism. Thus our models are based on neutral evolution. Our models
86 predict that there is indeed a preferred sequence of DNA feature evolution during *de*
87 *novoo* emergence, even under the assumption of neutrality.

88 Results

89 In this work, we developed mathematical models to estimate the rates and probabilities
90 of *de novo* gene emergence, as well as gene loss. A proto-gene emerges from non-genic
91 DNA, when the latter mutates to gain sequence features necessary for transcription
92 and translation (Carvunis *et al.*, 2012; Van Oss and Carvunis, 2019). Both transcription
93 and translation are complex processes involving many biomolecular complexes that
94 work in concert. Here, we focus on the minimal requirements for these processes to oc-
95 cur. Experimental data show that transcription in eukaryotes can be initiated genome
96 wide, especially in regions proximal to enhancers, and more specifically, in regions de-
97 pleted of nucleosomes (Andersson and Sandelin, 2020). Thus transcription of many
98 proto-genes can be initiated at a basal transcription rate, and may not require specific
99 core promoters. Although transcription may be initiated without a core promoter, a
100 stable RNA product often requires a poly-A tail that also facilitates RNA export to
101 cytoplasm (Richard and Manley, 2009; Stewart, 2019). This assumption is supported
102 by experimental data, where the majority of detectably expressed proto-genes have
103 polyadenylated RNAs (Neme and Tautz, 2016; Witt *et al.*, 2019; Majic and Payne, 2020;
104 Blevins *et al.*, 2021). Thus, we define a poly-A signal as the primary requirement for
105 transcription such that it should exist but only at the end of a transcribed region. For
106 translation, we only require the gene to have an open reading frame (ORF). We only
107 focus on intronless genes in this study.

108 Using simple probability models, we calculate the likelihood of finding the sequence
109 features that facilitate *de novo* gene emergence by random chance, and the rate at which
110 these features are gained and lost due to random mutations. These probability models
111 are essentially described by two kinds of probability. First is the probability of finding
112 a DNA sequence feature (such as the poly-A signal or an ORF; Methods section 2.1).
113 This probability depends on the nucleotide distribution that can be roughly approxi-
114 mated by the GC-content. The second kind of probability is the transition probability,
115 that is the probability of a DNA sequence mutating to another (Methods section 2.2).
116 Transition probabilities that describe the gain and loss of DNA sequence features, de-
117 pend on mutation rate, mutation bias, and nucleotide composition. Our model of DNA
118 sequence evolution is in principle similar to a theoretical model described in a previous

119 study (Behrens and Vingron, 2010). In this study, the authors have calculated the time
120 required for different transcription factor binding sites to emerge in gene promoters.
121 We extend this principle to calculate not only the probability of sequence emergence
122 but also the likelihood that an emerged sequence will remain intact or be lost (Methods
123 section 3).

124 In our model we primarily focus on *de novo* gene emergence at one genomic locus. We
125 assume that transcription is initiated at this locus with a probability of 0.12. We based
126 this estimate on the fraction of intergenic open chromatin in *Drosophila melanogaster*
127 that is occupied by *cis*-regulatory elements (Methods section 5). We used a default
128 GC-content of 42% in our calculations which is reasonably close to the total genomic
129 GC-content of human (41%, Merchant *et al.*, 2007), *Drosophila melanogaster* (41.6%, Gra-
130 mates *et al.*, 2022), and *Saccharomyces cerevisiae* (38%, Wood *et al.*, 2002). Because GC
131 content can vary between different genomic loci, we also performed our calculations
132 with other values of GC-content. We note that nucleotide composition in any DNA re-
133 gion may not be uniformly distributed. Therefore, we also included calculations using
134 a more realistic approximation of nucleotide distribution, using frequencies of short
135 DNA sequences in the intergenic genome of *D.melanogaster*. Specifically, we calculated
136 the frequency of all DNA trimers for codon probabilities, and of all DNA hexamers for
137 poly-A signal probabilities (Methods section 5).

138 We used a mutation rate of 7.8×10^{-9} mutations per nucleotide position per generation,
139 which corresponds to spontaneous mutation rate estimated in a *D. melanogaster* pop-
140 ulation that was subjected to periodic population bottlenecks (mutation accumulation
141 line, Schrider *et al.*, 2013). We used mutation bias data (Table 1) from the same study
142 (Schrider *et al.*, 2013) which was in agreement with another previous study on human
143 pseudogenes (Zhang and Gerstein, 2003).

144 Using a similar approach, we predict how the nucleotide distribution and random
145 mutations shape protein composition and evolution.

146 Overall, our models define a null hypothesis under which protein coding DNA se-
147 quences evolve neutrally through mutation pressure alone.

148 **How likely is gene loss?**

149 If a gene is found in one population but not populations, then the most parsimonious
150 assumption is that this gene emerged only in one specific population. We asked what
151 is the chance that the gene was present in an ancestral population but was lost in all

152 lineages except one. To this end we calculated the probability of gene emergence and
 153 gene loss (Methods section 3). Specifically, the probability of gene gain can be de-
 154 fined as the sum of three probabilities. First is the probability that an ORF already
 155 exists and is not lost due to mutations ($P_{ORF-stay}$, Methods section 4), while a mutation
 156 causes transcription to emerge ($P_{RNA-gain}$, Methods section 5). Second is the probability
 157 that an ORF emerges ($P_{ORF-gain}$, Methods section 4) in a region of DNA that is already
 158 transcribed and continues to be transcribed ($P_{RNA-stay}$, Methods section 5). Third is the
 159 probability that neither of the two features already exist and both emerge at the same
 160 time due to mutations (this probability is very small and is negligible). We found that
 161 the probability of transcription gain is dependent on the existence of an ORF, such that
 162 it is higher when the ORF is already present. This dependence exists because an ORF
 163 does not have stop codons (TAA, TAG and TGA) in its sequence. This in turn, restricts
 164 the number of positions where a poly-A signal, that has a TAA in its sequence, can exist
 165 (Methods section 6). Because of the same reason, the probability of ORF gain is higher
 166 when the DNA is already transcribed. Therefore, we additionally calculated the con-
 167 ditional probabilities of transcription gain ($P_{RNA-gain|ORF}$) and ORF gain ($P_{ORF-gain|RNA}$)
 168 given the condition that an ORF and a transcript already exists, respectively. We thus
 169 define the total probability of gene gain ($P_{gene-gain}^*$) as:

$$P_{gene-gain}^* = P_{RNA-gain|ORF} \times P_{ORF-stay} + P_{ORF-gain|RNA} \times P_{RNA-stay} + P_{RNA-gain} \times P_{ORF-gain} \quad (1)$$

170 Equation 1 describes the total probability of gene gain that includes the probability
 171 that the gene was not previously present. To more precisely address the question of
 172 whether a gene that was absent in an ancestor but emerges in a descendant, we calcu-
 173 lated the corresponding conditional probability of gene gain:

$$P_{gene-gain} = P_{gene-gain}^* / (1 - P_{ORF} - P_{RNA}) \quad (2)$$

174 We note that the difference between the total probability ($P_{gene-gain}^*$) and the conditional
 175 probability ($P_{gene-gain}$) of gene gain, is very small.

176 Next, we calculated the probability of gene loss, given the gene (transcription and ORF)
 177 is already present ($P_{gene-loss}$). It is the sum of the probabilities that transcription ($P_{RNA-loss}$,
 178 Methods section 5) or the ORF is lost due to mutations ($P_{ORF-loss}$, Methods section 4).
 179 Specifically:

$$P_{gene-loss} = P_{RNA-loss} + P_{ORF-loss} \quad (3)$$

180 The probability that a gene is lost n times independently, is $(P_{gene-loss})^n$. To find out
 181 how many independent gene loss events are as likely as a single gene gain event, we
 182 calculate the ratio of logarithms (log-log ratio) of gene gain and gene loss for proto-
 183 genes with an ORF ranging from 30 – 300 codons in length.

$$\text{Gene losses per gene gain} = \log(P_{gene-gain}) / \log(P_{gene-loss}) \quad (4)$$

184 For example, a log-log ratio value of 2 would indicate that two independent gene loss
 185 events are as likely as a single gene gain event. We calculated this ratio for genes with
 186 different ORF lengths, because ORF gain and loss probabilities depend on the length
 187 of the ORF. We found that proto-genes with a GC-content of 42% and ORFs longer
 188 than 40 codons, can be lost two times independently in the time frame of one gene
 189 gain event (log-log ratio = 2, [Figure 1A](#)). For proto-genes with GC-contents of 34% and
 190 50%, the minimum number of codons required for a log-log ratio of 2, are 42 and 33
 191 respectively. To better understand how GC-content determines the likelihood of gene
 192 gain relative to gene loss, we calculated the log-log ratio of the two probability val-
 193 ues for proto-genes with different GC-content ([Figure S4](#)), and that contain ORFs of
 194 different lengths (codons). We found that the effect of GC-content on the log-log ra-
 195 tio depends on the length of the ORF. Specifically, the log-log ratio steadily increases
 196 with GC-content ranging from 30% to 60%, when the proto-gene contains an ORF with
 197 30 codons. When the ORF has 60 – 90 codons, the log-log ratio initially decreases
 198 with GC content and then increases. For ORFs containing more than 120 codons, the
 199 log-log ratio decreases with increasing GC-content (within the analysed range). **Be-**
 200 **cause GC-content may not accurately describe the nucleotide composition of a locus,**
 201 **we calculated gene gain and loss probabilities using frequencies of DNA trimers and**
 202 **hexamers from the intergenic genome of *D. melanogaster*. Specifically, we used the fre-**
 203 **quencies of hexamers to estimate the stationary, gain and loss probabilities of poly-A**
 204 **signals (which is also a hexamer) and the RNA. Likewise, we used the frequencies of**
 205 **trimers to estimate the stationary, gain and loss probabilities of different codons and**
 206 **the ORF. Using these values, we found that protogenes with more than 49 codons can**
 207 **be independently lost twice, in the time frame of one gene gain event (orange line,**
 208 **[Figure 1A](#)).**

209 **Natural selection and effective population size (N_e) determine the rate at which new**
 210 **genes spread in a population ([Kimura, 1962](#)). To understand the combined effect of**

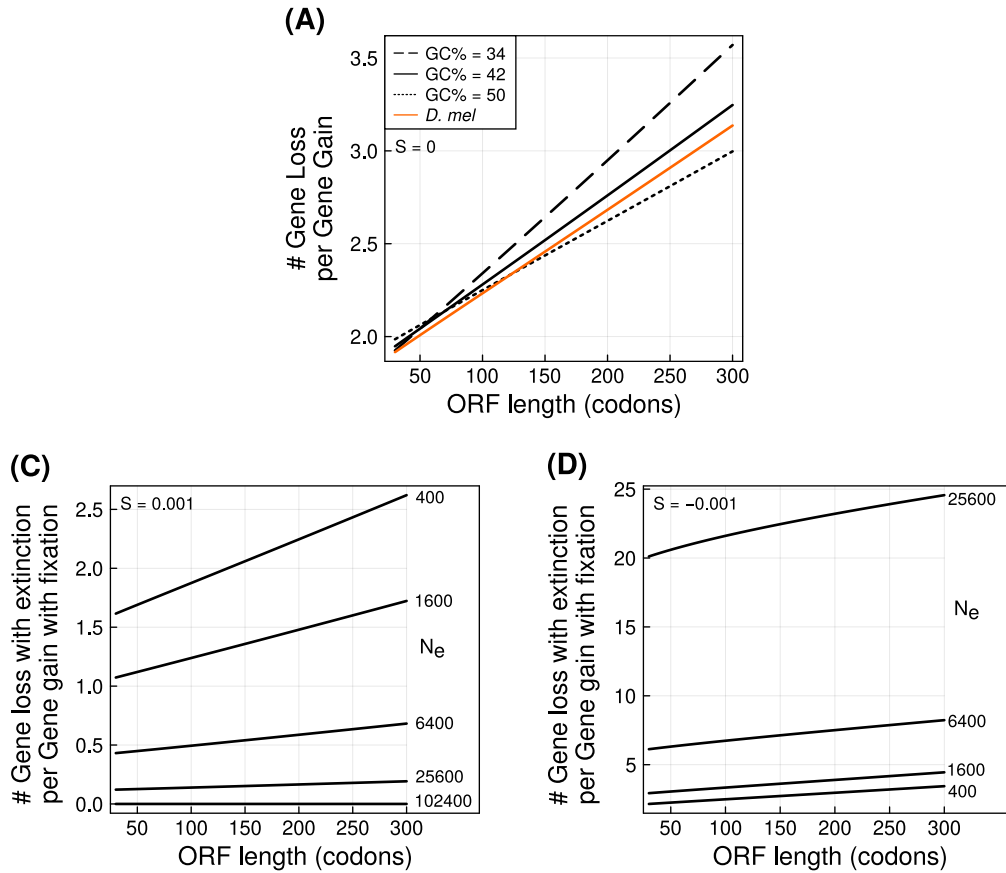


Figure 1: Genes are more likely to be independently lost twice, than being born once. The vertical axis in all panels shows the number of independent gene loss events with extinction (Equation 3) that can occur relative to one gene gain event with fixation (Equation 4, Table 2), under (A) no selection, (B) weak positive selection and (C) weak negative selection on the expression of the protein from a proto-gene. Horizontal axes show the number of codons in the ORF. Black lines in panel (A) denote the probability values estimated from overall GC content (34%, 42% and 50%) and the orange line denotes the probability values estimated from trimer and hexamer frequencies in *D. melanogaster* intergenic genome. In panels (B) and (C) we use the probability values from *D. melanogaster* intergenic genome. N_e denotes effective population size and s denotes the selection co-efficient.

211 mutations, and selection we calculated two fixation probabilities for populations con-
 212 sisting different number of diploid individuals. First is the probability that any one in-
 213 dividual in a population gains a gene ($2N_e \times P_{gene-gain}$), and this gene eventually spreads
 214 in entire the population via selection and population dynamics (P_{fix} , Table 2). The sec-
 215 ond probability describes the event that all individuals in a population have two copies
 216 of the proto-gene, one copy of the gene is lost in an individual ($2N_e \times P_{gene-loss}$), and the
 217 gene eventually goes extinct from the entire population due to population dynamics
 218 (P_{ext} , Table 2). We calculated the log-log ratio of the two probabilities under three sce-
 219 narios (Table 2), wherein the gene provides no fitness advantage (neutral, $s = 0$), the
 220 gene is marginally beneficial (positive, $s = 0.001$), and the gene is marginally dele-
 221 terious (negative, $s = -0.001$). We note that gene gain and loss probabilities due to

222 mutation are so small that it is unlikely that two or more individuals gain or lose a
223 gene simultaneously. Under the neutral scenario, the probability of fixation of an al-
224 lele is same as the probability of the mutation that gives rise to that allele (Figure 1A).
225 Extinction of a marginally beneficial gene is more likely than its fixation only when
226 populations are small or when the gene has a long ORF (Figure 1B). For example, a
227 gene of any length is less likely to go extinct than it is fixed, in populations with at
228 least 6400 individuals. Furthermore, a beneficial gene can be independently lost in
229 more than two small populations with 400 individuals, only when it has more than
230 133 codons. When a gene is mildly deleterious, its is very likely to go extinct such that
231 independent loss in more than two populations is more likely than its fixation in one
232 population (Figure 1C). We further note that time to fixation of a rare allele (for exam-
233 ple, a newly emerged protogene) would be much smaller ($4N_e$ generations; Kimura,
234 1962) than the time to gain ($10^{14} - 10^{20}$ generations) or lose it ($10^7 - 10^8$ generations).

235 Overall, our analysis suggests that a proto-gene expressed in only one population may
236 not necessarily mean that it emerged for the first time in this population. This is also
237 true for a model of gene emergence that requires a core promoter (TATA-box or Inr) for
238 initiating transcription (Figure S1).

239 Despite low probability of gene gain at one particular locus ($\sim 10^{-15}$ per generation),
240 the overall rate of gene gain throughout the genome could be much higher. To this end
241 we calculated the total gain rate of genes of any length, at any locus in *D. melanogaster*
242 intergenic genome. We found this rate to be 2.8×10^{-5} new genes per generation, that
243 corresponds to approximately 7.5×10^{-5} new genes per year (assuming *D. melanogaster*
244 generation time of two weeks; Fernández-Moreno *et al.*, 2007). We performed an anal-
245 ogous analysis on the rate of gain of new transcripts and found this rate to be 1.53
246 transcripts per year. This estimate is slightly higher than the rate of transcript gain pre-
247 dicted from *D. melanogaster* population transcriptomics data (0.13 – 0.34; Grandchamp
248 *et al.*, 2023).

249 Does *de novo* gene emergence follow a preferred trajectory of events?

250 For a proto-gene to emerge from a non-genic DNA sequence, both transcription and
251 ORF need to emerge. That is, probability of gene emergence is equal to the product
252 of probabilities of transcription gain and ORF gain. Thus it may appear that the order
253 of the occurrence of these two events does not matter. However, gene emergence is
254 much more likely when one of the two features already exists (Equation 1) and there-
255 fore it has two possible trajectories – ORF emerges first (ORF-first) or transcription

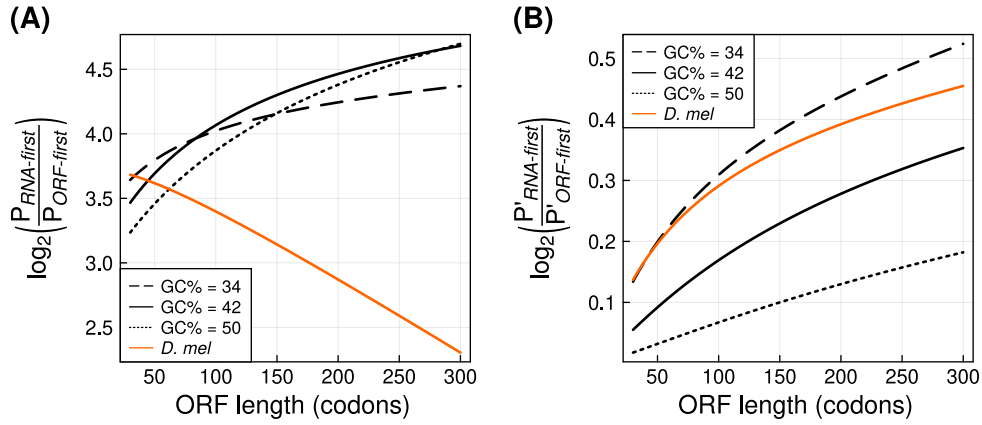


Figure 2: Proto-genes preferentially emerge RNA-first. The vertical axis shows the \log_2 transformed ratio of the probabilities of the RNA-first and the ORF-first trajectories, described as **(A)** a single-step process (Equations 5 & 6) and **(B)** a two-step process (Equations 7 & 8). A positive value suggests that RNA-first trajectory is more likely than ORF-first trajectory, and *vice versa*. In both panels, the horizontal axis denotes the number of codons in the ORF. Black lines denote the probability values estimated from overall GC content (34%, 42% and 50%) and the orange line denotes the probability values estimated from trimer and hexamer frequencies in *D. melanogaster* intergenic genome.

256 emerges first (RNA-first). Furthermore, emergence of an ORF is more likely when the
 257 proto-gene is already transcribed, and *vice versa* (Methods section 6). To understand
 258 whether *de novo* gene emergence has a preferred trajectory, we calculated two proba-
 259 bilities. First, the probability that an RNA exists and mutations cause a gain of an ORF
 260 but no disruption of transcription. This probability denotes the trajectory where RNA
 261 emerges first ($P_{RNA-first}$, Equation 5).

$$P_{RNA-first} = P_{RNA-stay} \times P_{ORF-gain|RNA} \quad (5)$$

262 The second probability ($P_{ORF-first}$, Equation 6) that denotes the trajectory where ORF
 263 emerges first, is the probability that an ORF already exists and mutations cause a gain
 264 of transcription but do not disrupt the ORF.

$$P_{ORF-first} = P_{ORF-stay} \times P_{RNA-gain|ORF} \quad (6)$$

265 We calculated the log-transformed ratio of $P_{RNA-first}$ and $P_{ORF-first}$, such that a positive
 266 value of the ratio would mean that ORF gain in an already transcribed DNA (RNA-
 267 first) is more likely than transcription gain for an untranscribed ORF (ORF-first). In
 268 other words, the RNA-first trajectory is more feasible. Likewise, a negative value of
 269 the ratio would suggest that the ORF-first trajectory is more feasible. We found that

270 proto-genes with ORFs of all the investigated lengths, preferentially emerge RNA-first
 271 (Figure 2A). More specifically, the likelihood of the RNA-first trajectory relative to that
 272 of the ORF-first trajectory increases with increasing ORF length. This is true for loci
 273 with three different values of GC content (34%, 42%, 50%). However, our complemen-
 274 tary analysis using DNA oligomer frequencies from *D. melanogaster* intergenic genome
 275 shows that the likelihood of RNA-first trajectory relative to that of ORF-first trajectory,
 276 reduces with increasing ORF length. Nonetheless, RNA-first trajectory is still the more
 277 feasible trajectory of *de novo* gene emergence. In a more stringent scenario of *de novo*
 278 emergence, where a core promoter (TATA box or Inr) is necessary for transcription,
 279 proto-genes containing short ORFs can emerge ORF-first (Figure S2)

280 A more stringent definition of the RNA-first trajectory would describe a two-step pro-
 281 cess. In the first step, an RNA emerges in an untranscribed region of DNA but an ORF
 282 does not emerge. In the second step an ORF emerges, while transcription stays intact
 283 (Equation 7). Likewise, a ORF-first trajectory can be defined by a two-step probability.
 284 In the first step ORF emerges in an untranscribed DNA region, and in the second step
 285 transcription emerges, while the ORF remains intact (Equation 8).

$$P'_{RNA-first} = \underbrace{P_{RNA-gain} \times (1 - P_{ORF} - P_{ORF-gain})}_{first\ step} \times \underbrace{P_{RNA-stay} \times P_{ORF-gain|RNA}/(1 - P_{ORF})}_{second\ step} \quad (7)$$

$$P'_{ORF-first} = \underbrace{P_{ORF-gain} \times (1 - P_{RNA} - P_{RNA-gain})}_{first\ step} \times \underbrace{P_{ORF-stay} \times P_{RNA-gain|ORF}/(1 - P_{RNA})}_{second\ step} \quad (8)$$

286 Even with this stringent definition, *de novo* emergence via the RNA-first trajectory is
 287 more probable than via the ORF-first trajectory (Figure 2B). For all the three GC-content
 288 values, as well the specific analysis with *D. melanogaster* genome, the two-step RNA-
 289 first trajectory becomes increasingly more probable than the two-step ORF-first trajec-
 290 tory, with increasing ORF length.

291 **Would extensive transcription loss suggest negative selection of toxic** 292 **proteins?**

293 Proto-genes that do not provide any fitness benefit to an organism may not be fixed
 294 in populations via natural selection. These genes may be lost due to mutation pres-
 295 sure. Some newborn proto-genes can also encode toxic proteins, that may aggregate or
 296 interfere with physiology in some other way. These genes would thus be eliminated

297 from the population genomes via negative selection. We note that an ORF is lost if the
 298 start codon is mutated, the stop codon is mutated to an amino acid encoding codon
 299 (non-stop mutation), or if an amino acid encoding codon is mutated to a stop codon
 300 (premature stop/non-sense mutation). However, it is likely that a non-stop mutation
 301 or a non-sense mutation, can still result in translation of a protein (extended or trun-
 302 cated, respectively). Furthermore, non-stop mutations can also lead to cellular toxicity
 303 (Choe *et al.*, 2016). Thus ORF loss does not ensure elimination of toxic proteins, which
 304 in turn suggests that transcription loss may more effectively inactivate the associated
 305 genes.

306 To understand which is the most probable mechanism of gene loss, we compared the
 307 probabilities of ORF loss ($P_{ORF-loss}$) and transcription loss ($P_{RNA-loss}$). We found that ORF
 308 loss is more probable than transcription loss, especially so when the ORFs are long
 309 (Figure 3A). If transcription is strictly dependent on a core promoter (TATA-box or
 310 Inr), small genes (30 codons) with low GC content (>34%) can be preferentially lost by
 311 loss of transcription (Figure S3). When the organismal fitness is affected only by the
 312 protein product of a proto-gene, then both RNA loss and ORF loss would affect the
 313 fitness equally. Thus the likelihood of gene extinction through ORF loss or RNA loss
 314 depends solely on the probability of these events.

315 Our analysis of gene loss mechanism considers the loss of one feature (for example
 316 transcription) without requiring the other feature to remain intact (for example the
 317 ORF). With more stringent analysis where we consider the loss of only one feature but
 318 not the other, we find a the that exclusive ORF loss is more likely than exclusive RNA

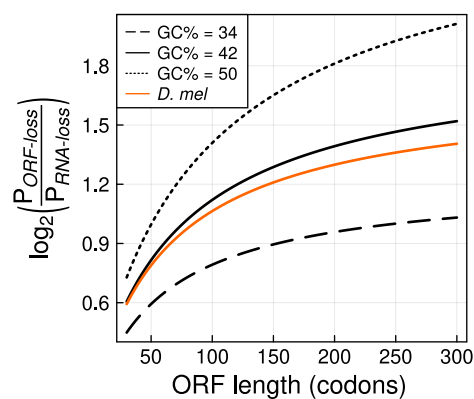


Figure 3: ORF loss in proto-genes is more probable than transcription loss. The horizontal axis denotes the number of codons in the ORF. The vertical axis shows the \log_2 transformed ratio of ORF loss and RNA loss probabilities, such that a positive value means ORF loss is more probable than RNA loss, and *vice versa*. Black lines denote the probability values estimated from overall GC content (34%, 42% and 50%) and the orange line denotes the probability values estimated from trimer and hexamer frequencies in *D. melanogaster* intergenic genome.

319 loss (Figure S5). The difference between the likelihood of the two gene loss mecha-
320 nisms is more remarkable when we consider exclusive loss of one feature.

321 Overall, our analysis suggests widespread transcription loss may be indicative of a
322 negative selection on protein expression.

323 Does mutation bias shape protein composition?

324 In previous sections, we showed that mutation bias affects the rate of *de novo* gene
325 emergence and loss. We next turned our attention to whether this bias affects the com-
326 position (and thereby the chemistry) of proteins encoded by proto-genes. To this end,
327 we first asked if the expected frequency of different amino acids is uniform. For a
328 GC-content of 42% and with DNA trimer frequencies from *D. melanogaster* intergenic
329 genome, we found that it is not uniform. Amino acids like leucine (L) and serine (S)
330 have a higher expected frequency than other amino acids. On the other hand, amino
331 acids like methionine (M) and tryptophan (W) are less probable (Figure 4A). This non-
332 uniformity is to a great extent determined by number of degenerate codons for an
333 amino acid. However, nucleotide composition also determines the frequency of an
334 amino acid. For example, leucine, arginine and serine, all have six codons each. How-
335 ever, leucine is more likely to be encoded than serine in a random stretch of DNA,
336 given a uniform GC-content of 42% (Figure 4A). The same is true for the frequencies of
337 these amino acids estimated from the DNA trimer distribution. Overall, our analysis is
338 roughly in agreement with many previous studies (Ohta and Kimura, 1971; Shen *et al.*,
339 2006; Gardini *et al.*, 2016). However, notable differences exist between our analysis
340 and that of the previous studies because the latter primarily focused on characterized
341 proteins and not random proto-gene derived protein sequences.

342 Next, we aimed to find out if some amino acid substitutions are more probable than
343 the others. To this end, we calculated the gain probability (Equation 9) for every amino
344 acid. Specifically for every amino acid x , we calculated the average probability that it
345 substitutes any of the other nineteen amino acids, due to mutations. More precisely,
346 the total probability that an amino acid y will mutate to an amino acid x is the product
347 of the expected frequency of amino acid y (P_y), and the substitution rate from x to y
348 ($\mu_{y \rightarrow x}$). The gain probability for an amino acid x is then defined as:

$$G_x = \sum_{b \neq a} P_y \times \mu_{y \rightarrow x} \quad (9)$$

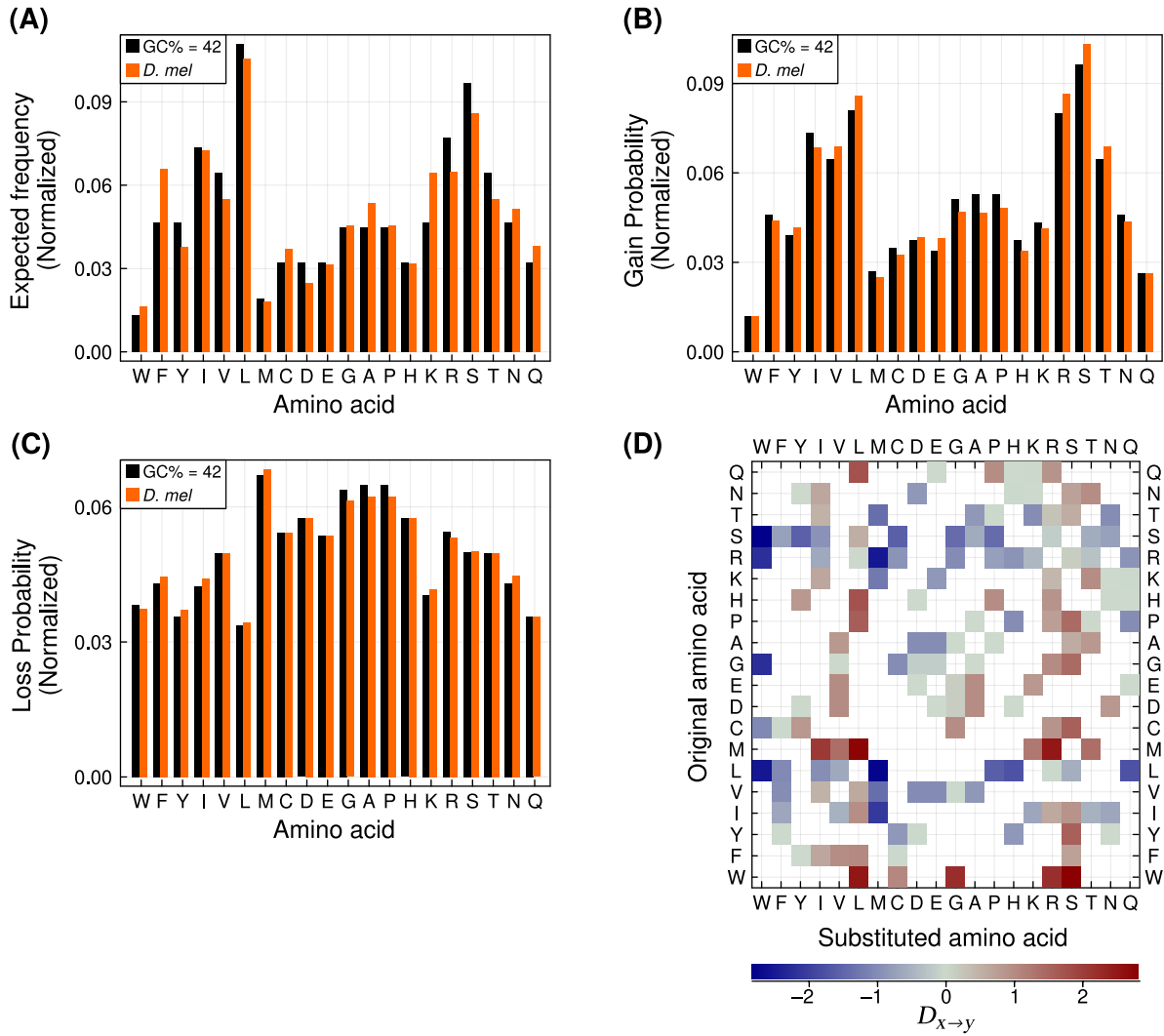


Figure 4: Mutation bias and GC-content shape protein composition. Panels (A), (B) and (C) show expected frequency, gain probability (Equation 9), and loss probability (Equation 10) for different amino acids (horizontal axes). The vertical axis in all these three subplots show the normalized values of these measures such that their sum over all amino acids is equal to one. Panel (D) shows the directionality (colored tiles, $D_{x \rightarrow y}$, Equation 11) of the substitution from an amino acid (vertical axis) to another (horizontal axis). Positive values of $D_{x \rightarrow y}$ are shown in shades of red, whereas negative values are shown in shades of blue. We only show substitutions that can be realized by a single nucleotide mutation but we include all possible nucleotide mutations in the directionality calculations for these substitutions. In all panels, amino acids are ordered based on their chemical similarity (Kim *et al.*, 2009).

349 We found that amino acid gain probabilities were also non-uniform across different
 350 amino acids (Figure 4B) but they had a strong positive correlation with the expected
 351 frequency (GC = 42%, Pearson $\rho = 0.955$, $P < 10^{-10}$; *D. melanogaster* trimer frequencies,
 352 Pearson $\rho = 0.841$, $P < 10^{-10}$).

353 Next, we calculated how easily is an amino acid mutated to any other amino acid.
 354 Specifically, the conditional loss probability for an amino acid (L_x), given the amino

355 acid already exists, is defined as:

$$L_x = \sum_{x \neq y} \mu_{x \rightarrow y} \quad (10)$$

356 The loss probability was also not uniformly distributed (Figure 4C) but it did not sig-
357 nificantly correlate with the expected frequency of the amino acids (GC = 42%, Pearson
358 $\rho = -0.255$, $P = 0.2818$; *D. melanogaster* trimer frequencies, Pearson $\rho = -0.356$, $P =$
359 0.1232).

360 Next, we investigated if some amino acid substitutions are more common than the
361 others. As expected, amino acid substitutions that require more than one nucleotide
362 change were much less probable than those needing just a single nucleotide change.
363 Therefore, we focused our analysis on amino acid substitutions whose total proba-
364 bilities were more than that of a double nucleotide mutation (note that this procedure
365 identifies amino acid substitutions that are possible via single nucleotide mutations but
366 it does not exclude multinucleotide mutations from the calculation of amino acid sub-
367 stitution probabilities). Although single nucleotide mutations are more probable than
368 double mutations, the different single mutations do not occur at the same rate due to
369 mutation rate bias. Thus an amino acid substitution (for example K→E) may not be as
370 likely to occur as the reverse substitution (E→K). That is substitutions between a pair of
371 amino acids may have a directionality. Although many previous studies have explored
372 the likelihood of different amino acid substitutions (Ohta and Kimura, 1971; Dayhoff
373 *et al.*, 1978; Henikoff and Henikoff, 1992; Gonnet *et al.*, 1992; Jones *et al.*, 1992; Whelan
374 and Goldman, 2001; Schneider *et al.*, 2005; Kosiol *et al.*, 2007; Le and Gascuel, 2008),
375 none to our knowledge have focused primarily on the directionality of these substitu-
376 tions. To verify if such a directionality exists, we calculated the log transformed ratio
377 of forward and reverse substitution probability, $D_{x \rightarrow y}$, for every pair of amino acids, x
378 and y . A positive value of $D_{x \rightarrow y}$ means that x to y substitution is more likely than y to
379 x substitution, and *vice versa*.

$$D_{x \rightarrow y} = \log_2 \left(\frac{\mu_{x \rightarrow y}}{\mu_{y \rightarrow x}} \right) \quad (11)$$

380 As we mentioned in the previous paragraph, we restricted this analysis to substitutions
381 that can be achieved via single nucleotide mutations. We found that most substitutions
382 indeed had a preferred direction ($D_{xy} \neq 0$, Figure 4D). The median absolute value of
383 D_{xy} was 0.905 and the maximum value was 2.836. That is, in 38 out of 75 amino acid

384 pairs we analysed, one direction of substitution was at least 1.87 times more probable
385 than the other. For example, E→K substitution is 1.8 times more probable than K→E
386 substitution. This may appear as a small number but many such unidirectional substi-
387 tutions occurring in an evolving protein could indicate a directional selection. **We note**
388 **that that actual frequencies of specific substitutions depend on the initial distribution**
389 **of different codons in the ORF sequence.**

390 A previous study found that mutations that cause hydrophobic amino acids to appear
391 on protein surfaces, can lead to evolution of protein dimers (Hochberg *et al.*, 2020). The
392 authors also suggested that such an evolutionary process may be widespread because
393 mutation bias tends to facilitate emergence of hydrophobic amino acids. They based
394 this argument on another study that estimated mutation rate bias in bacteria (Her-
395 shberg and Petrov, 2010). The authors (Hochberg *et al.*, 2020) further suggested that
396 hydrophobic amino acids are more frequently found in random protein sequences.

397 We asked if our model also makes similar predictions. To this end, we first calculated
398 the percentage of amino acids in a protein sequence that are expected to be hydropho-
399 bic. A popular hydrophobicity scale of amino acids is based on solubility of an amino
400 acid in water or ethanol (Kyte and Doolittle, 1982). However, this scale does not classify
401 tryptophan as a hydrophobic amino acid. Another hydrophobicity scale estimated by a
402 different study (Wimley and White, 1996) is more biologically realistic and is based on
403 rates of transfer of different amino acids between water and a hydrophobic medium.
404 This study analysed two hydrophobic media – a lipid bilayer and octanol, and also
405 considered the effect of peptide bonds in the calculations. Therefore we used the hy-
406 drophobicity scale by Wimley and White (1996) for classifying amino acids such that
407 hydrophobic amino acids have a negative free energy change of transfer from water
408 to hydrophobic media. Conversely, a hydrophilic amino acid would have a positive
409 value of the free energy change. Based the octanol hydrophobicity scale, we classified
410 the following amino acids as hydrophobic – cysteine (C), phenylalanine (F), isoleucine
411 (I), leucine (L), methionine (M), valine (V), tryptophan (W) and tyrosine (Y). With the
412 lipid bilayer hydrophobicity scale, valine is not classified as hydrophobic, possibly be-
413 cause it does not integrate well in the bilayer. Therefore we used the octanol scale for
414 our calculations.

415 We found that ~40% of amino acids in a protein sequence are expected to be hydropho-
416 bic (with a GC-content of 42% as well as with *D. melanogaster* trimer distribution). This
417 finding is broadly in agreement with the previous study (Hochberg *et al.*, 2020). We
418 note that cytosolic proteins may require more than 42% of their constituent amino acids
419 to be hydrophobic, in order to fold efficiently (Dill, 1985). Next, we asked if a protein

420 sequence in general tends to be hydrophobic. Specifically, we calculated expected hy-
 421 drophobicity ($\bar{\alpha}$, Equation 12) of a protein sequence, based on expected frequency of
 422 different amino acids (P_x , Figure 4) and their hydrophobicity values (α_x).

$$\bar{\alpha} = \sum_x P_x \alpha_x \quad (12)$$

423 We found that the expected hydrophobicity is higher than zero (0.38 with GC% = 42;
 424 0.36 with *D. melanogaster* trimer frequencies), which suggests that random protein se-
 425 quences are not on an average, hydrophobic.

426 To find out if random mutations indeed cause hydrophobic amino acids to occur in
 427 protein sequences, we calculated the probability that any mutation substitutes a non-
 428 hydrophobic amino acid with a hydrophobic amino acid ($P_{\alpha\text{-gain}}$, Equation 13). Like-
 429 wise, we calculated the probability that any mutation causes a hydrophobic to non-
 430 hydrophobic amino acid substitution ($P_{\alpha\text{-loss}}$, Equation 14).

$$P_{\alpha\text{-gain}} = \sum_{\substack{x \\ |\alpha_x > 0}} \sum_{\substack{y \\ |\alpha_y \geq 0}} P_y \times \mu_{y \rightarrow x} \quad (13)$$

$$P_{\alpha\text{-loss}} = \sum_{\substack{x \\ |\alpha_x > 0}} \sum_{\substack{y \\ |\alpha_y \geq 0}} P_x \times \mu_{x \rightarrow y} \quad (14)$$

431 Next, we calculated the ratio of $P_{\alpha\text{-gain}}$ and $P_{\alpha\text{-loss}}$, and found that it is slightly more
 432 than 1 (1.114 with GC% = 42; 1.015 with *D. melanogaster* trimer frequencies). This sug-
 433 gests that gain of a hydrophobic amino acid is slightly more probable than its loss.
 434 We performed a complementary analysis where we calculated the average change in
 435 hydrophobicity due to any random mutation ($\bar{\Delta\alpha}$), defined by:

$$\bar{\Delta\alpha} = \frac{\sum_x \sum_y \mu_{x \rightarrow y} \times (\alpha_y - \alpha_x)}{\sum_x \sum_y \mu_{x \rightarrow y}} \quad (15)$$

436 We found the average change in hydrophobicity to be slightly less than zero (-0.0798),
 437 which suggests that on an average mutations cause a very small increase in hydropho-
 438 bicity. We reiterate that hydrophobic amino acids have a negative value hydrophobic-

439 ity and thus a negative change in hydrophobicity also denotes a shift towards a more
440 hydrophobic protein. To focus exclusively on non-hydrophobic to hydrophobic amino
441 acid substitutions, we excluded amino acid pairs where both amino acids are either
442 hydrophobic or are non-hydrophobic. That is, we only analysed pairs of amino acids
443 that have opposite signs of hydrophobicity. With this focused analysis, we find that
444 the average change in hydrophobicity is comparable (-0.0793). For both the analyses,
445 we note that the actual change in hydrophobicity will depend on the composition of
446 the ancestral protein sequence.

447 Overall our analyses show that random mutations may preferentially cause hydropho-
448 bic amino acids to accumulate in protein sequence but this effect may be very small.

449 Discussion

450 In this work, we addressed some fundamental questions about *de novo* emergence
451 of protogenes. Broadly, we asked how quickly these genes can emerge and be lost,
452 whether their birth and death has a preferred trajectory of mutational events, and how
453 the composition of their protein sequence is determined. To this end, we developed
454 mathematical models and used them to address specific questions. These models are
455 based on mutation probability and mutation rate bias, and represent neutral evolu-
456 tion where DNA sequences are not evolving under natural selection. Thus our models
457 provide an opportunity to test the alternative hypothesis that some proto-genes are
458 evolving under selection.

459 We show that when a protein product of a proto-gene doesn't affect organismal fitness,
460 the gene can be lost much more rapidly than they are gained, such that it can be inde-
461 pendently lost two to three times in the time required for it to emerge (Figure 1A). In a
462 hypothetical case, a proto-gene may be present in a population (A) but is absent in two
463 other geographically isolated populations (B and C). Two evolutionary scenarios could
464 explain this observation. In the first scenario, the proto-gene simply emerged for the
465 first time in population A. The second scenario posits that the proto-gene was present
466 in the ancestral population but was lost in both population B and population C. Our
467 models suggest that the second scenario is more likely than is assumed. **This inference**
468 **can be further extended to the level of closely related species where the divergence**
469 **times are relatively shorter than the timescale of multiple gain and loss events within**
470 **a lineage.** Future phylogenetic studies on proto-genes should consider this finding for
471 inferring the dynamics of gene gain and loss. **Proto-genes that encode proteins that are**
472 **beneficial to an organism's fitness are likely to be protected in the genome against gene**

473 loss (Lee *et al.*, 2019). Our models indeed show that loss of even moderately beneficial
474 genes is unlikely in large populations (Figure 1B).

475 Within the limitations of our model assumptions, we answer the long standing ques-
476 tion of which trajectory of mutational events leads to *de novo* emergence. That is,
477 whether ORF emerges first or transcription emerges first (Schmitz and Bornberg-Bauer,
478 2017). We show that, in the absence of selection, genes are likely to emerge in existing
479 transcripts (transcription emerges first, Figure 2). Long non-coding RNAs may be a
480 potential source of proto-genes but this evolutionary process may be constrained by
481 a different set of factors. A previous study has shown that new transcripts frequently
482 arise in regions overlapping existing genes. These new transcripts can harbor ORFs
483 which in turn leads to *de novo* gene emergence (Blevins *et al.*, 2021). This is an example
484 of RNA-first trajectory. However, the ORFs present in these genes may be under in-
485 direct selection pressure that acts on the overlapping gene sequences. Therefore, they
486 may be less frequently mutated or lost.

487 Although our model is based on the assumption of neutral evolution, one can speculate
488 the effects of positive selection on the proto-gene. For example, if an ORF encodes a
489 protein that is beneficial to the organism, then selection would cause the fixation of
490 the gene as soon as the it emerges via RNA-first trajectory. On the other hand, the
491 selection would have no effect on an ORF that is not transcribed, as long as it is does
492 not harbor a functional DNA sequence (such as regulatory elements), and is not under
493 indirect selection (for example, if it overlaps with an existing gene). The ORF may
494 mutate before transcription emerges and the corresponding mutated protein may not
495 be beneficial anymore. Thus positive selection may further increase the likelihood of
496 the RNA-first trajectory.

497 We show that transcription loss is less likely than the loss of long ORFs (Figure 3).
498 However, it may be difficult to find direct evidence for such a phenomenon. Specifi-
499 cally, it is difficult to infer from data if transcription was lost or never emerged. If the
500 protein encoded by an ORF is indeed toxic, it would be silenced or purged rapidly in
501 microevolutionary time scales. One way of inferring negative selection from transcrip-
502 tion loss may involve correlation of ORF sequence variation and transcription status.
503 Mutations within an ORF sequence are much more probable than mutations that dis-
504 rupt transcription or an ORF. This is so because the entire ORF sequence can have a
505 higher number of mutable sites than a promoter. If an ORF variant is untranscribed in
506 many populations (species or taxa), it is possibly toxic.

507 A large volume of work exists on the likelihood of amino acid substitutions and many
508 of these models are widely used in phylogenetics studies (Ohta and Kimura, 1971;

509 Dayhoff *et al.*, 1978; Henikoff and Henikoff, 1992; Gonnet *et al.*, 1992; Jones *et al.*, 1992;
510 Whelan and Goldman, 2001; Schneider *et al.*, 2005; Kosiol *et al.*, 2007; Le and Gascuel,
511 2008). However, commonly used amino acid substitution models such as PAM and
512 BLOSUM (Dayhoff *et al.*, 1978; Henikoff and Henikoff, 1992) do not focus on direction-
513 ality of mutations (the substitution matrices are symmetric). Our aim was to highlight
514 how proteins can evolve under mutation pressure and that mutation bias causes some
515 substitutions to occur more frequently than the others. Furthermore, we do not base
516 our substitution probability calculations on known protein sequences but rather on
517 random protein sequences that are more likely to be encoded in proto-genes. Using
518 our models, we show that for a pair of amino acids, one direction of substitution can
519 be significantly more likely than the other. If a substitution occurs more frequently
520 than it is expected to occur in our neutral model, then directional evolution can be a
521 possible explanation for the observation. We also show that random mutations can
522 lead to a small increase in a protein's hydrophobicity. Hydrophobic amino acids can
523 facilitate folding (Dill, 1985), but can also cause proteins to aggregate (Hochberg *et al.*,
524 2020). Thus the effect of hydrophobicity increasing mutations on a protein's stability
525 and toxicity, is dependent on where they occur in the protein structure.

526 Like all theoretical models, our models are not without limitations. First, we assume
527 that nucleotides are uniformly distributed according to GC-content. This is not true as
528 nucleotide composition can vary significantly across the genome. **Since our model pri-**
529 **marily focuses on gene emergence at one locus, it can be used to predict evolutionary**
530 **probabilities for different genomic loci. Furthermore, several predictions of our models**
531 **are not qualitatively affected by GC-content. More realistic nucleotide composition dis-**
532 **tribution estimated from *D. melanogaster* intergenic sequences, also qualitatively agree**
533 **with GC-content based predictions. Thus, most predictions are qualitatively robust to**
534 **small differences in parameters. Second, in our models, we assumed uniform mutation**
535 **rate throughout the genome. However, recent studies have shown that mutation rates**
536 **are also not uniform across the genome, and can vary significantly (Monroe *et al.*, 2022).**
537 Our models can be adapted to investigate specific genomic loci, where estimates of the
538 local nucleotide composition and mutation rate can more accurately determine the *de*
539 *nov* gene evolution in these loci. Third, our models are based on parameters such as
540 mutation rate and mutation bias, which we obtained from published data on human
541 and *D.melanogaster* genomes (Zhang and Gerstein, 2003; Schrider *et al.*, 2013). The ac-
542 curacy of the models' prediction will depend on the accuracy of parameters. However,
543 mutation rate itself only determines the values of different probabilities. Our results,
544 which are based on probability ratios, will not significantly change with a different
545 value of mutation rate. Mutation bias may have more significant effect on the predic-
546 tions than mutation rate. Existing data suggest that different species have a different

547 spectrum of mutations (Cano *et al.*, 2022), especially for species that are highly distant
548 from each other. Thus our inferences based on our chosen parameters will not uni-
549 versally apply to all species. Fourth, we assume that occurrence of a point mutation
550 is independent of another mutation. This is not the case for tandem mutations, where
551 multinucleotide mutations occur more frequently than their expected frequency un-
552 der the assumption of mutational independence (Harris and Nielsen, 2014). Fourth,
553 we did not specifically model the effect of many regulatory sequence features such
554 as enhancers, transcription factor binding sites, 5' and 3' untranslated regions, and
555 Kozak consensus sequence, that can influence the evolution of proto-genes. For exam-
556 ple *de novo* emergence is very likely to occur near enhancers (Majic and Payne, 2020).
557 To model this phenomenon in our analyses, we assumed that DNA loci feasible for
558 transcript emergence are close to known *cis*- regulatory elements (Gallo *et al.*, 2005;
559 Halfon *et al.*, 2007; Gallo *et al.*, 2010; Rivera *et al.*, 2018; Keränen *et al.*, 2022). Specif-
560 ically, we assumed that transcription can be initiated at a probability that is equal to
561 the proportion of such regulatory elements in the intergenic genome. This is a rough
562 estimate but it can be updated when more specific data becomes available. For ex-
563 ample, ATAC-seq data can provide reliable estimates of nucleosome depleted regions,
564 that can act as transcription initiation sites (Buenrostro *et al.*, 2013). This analysis can be
565 supplemented with data from techniques like native elongating transcript sequencing
566 (NET-seq; Churchman and Weissman, 2012) and global run-on sequencing (GRO-seq;
567 Lopes *et al.*, 2017), that can detect nascent transcripts. Our modeling framework also
568 allows incorporation of any number of sequence features in the calculations as long
569 as their sequences can be defined (Methods section 3). Finally, we ignore insertions,
570 deletions and transpositions as mechanisms of mutation. Insertions and deletions (in-
571 dels) are reported to be less frequent than substitutions (Schridder *et al.*, 2013), and both
572 indels and transpositions cause larger sequence alterations than point mutations. Our
573 modeling methodology also allows incorporation of these mechanisms as a source of
574 mutations.

575 A famous aphorism about theoretical models says “all models are wrong but some are
576 useful” (Box, 1979). Our models are not an exception. They may not be 100% accurate
577 but they make useful predictions that can lead to a more focused experimental vali-
578 dation. Many of our predictions are qualitatively robust to small parameter changes.
579 Some predictions are also in agreement with experimental data. For example, our pre-
580 diction of transcript gain rate is only slightly inflated than the estimated rates using
581 an empirical analysis of *D. melanogaster* populations (Grandchamp *et al.*, 2023). This
582 difference could have possibly resulted from our relaxed requirements for a promoter,
583 or other unknown factors that may have affected population evolution (for example,
584 the actual effective population size and environmental conditions; Lange *et al.*, 2021).

585 When supplemented with experimental data analysis our models can be made more
586 complex to accommodate diverse molecular mechanisms driving gene evolution, and
587 can provide more accurate predictions. Therefore our work opens up an opportunity
588 for theoretical and computational biologists to design analyses using our modeling
589 framework, that more accurately describe their system of interest.

590 Materials and Methods

591 1 Mutation probabilities

Substitution	Probability(μ)
A : T \rightarrow T : A	0.056
A : T \rightarrow G : C	0.243
A : T \rightarrow C : G	0.074
G : C \rightarrow A : T	0.483
G : C \rightarrow T : A	0.075
G : C \rightarrow C : G	0.069

Table 1: Mutation bias probabilities for different nucleotide mutations. A : T denotes an A-T base pair in a double stranded DNA. Thus A \rightarrow G mutation on one DNA strand would cause a T \rightarrow C mutation on the complementary strand. We describe the other mutations in the same way.

592 We calculated nucleotide substitution probabilities based on mutation rate and muta-
593 tion rate bias data. Specifically, we used a mutation rate (λ) of 7.8×10^{-9} mutations per
594 nucleotide position per generation (Schrider *et al.*, 2013). We derived our mutation bias
595 parameters from two published studies, the first on *Drosophila melanogaster* (Schrider
596 *et al.*, 2013), and the second on humans (Zhang and Gerstein, 2003). Table 1 shows the
597 exact values of mutation bias probabilities that we used in this study.

598 2 Probabilities of finding, gaining, and losing a DNA sequence

599 2.1 Probability of finding a DNA sequence

600 We calculated the probability of finding a DNA sequence based on global nucleotide
601 frequency distributions, given by the GC-content. Specifically, the probability of find-
602 ing either a G or a C is:

$$S = \frac{0.5 \times \text{GC}\%}{100}$$

603 The probability of finding an A or a T is: $W = 0.5 - S$. Using these values, we cal-
604 culated the probability of finding a DNA sequence motif by chance. For example, the
605 probability of finding the sequence ATG would be: $W \times W \times S$.

606 We also estimated the probability of finding specific DNA sequences in a reference
607 genome. Specifically, we calculated the frequencies of all 64 trimers and all 4096 hex-
608 amers in the genomic regions of *D. melanogaster* that exist in open chromatin, and do
609 not contain any known genes or regulatory elements (see [Methods section 5](#)).

610 2.2 Probability of gaining a DNA sequence

611 We calculated the probability of gaining a DNA sequence due to mutations using GC-
612 content, mutation rate and mutation bias. Specifically, we calculated the probability
613 that a DNA sequence does not exist, and it emerges due to specific nucleotide muta-
614 tions. More precisely, this probability is the product of two other probabilities. The
615 first is the probability of finding a DNA sequence (x) that is not the sequence of inter-
616 est. The second probability is that this sequence x mutates to the sequence of interest.
617 To explain this calculation better, we use the example of CTA mutating to ATG. The first
618 probability, that is the probability of finding CTA by chance is SW^2 . CTA mutates to
619 ATG via two nucleotide mutations ($C \rightarrow A$ and $A \rightarrow G$). Thus the probability of this DNA
620 change would be the probability of two nucleotide mutations (λ^2) multiplied by two
621 mutation bias probabilities ($G : C \rightarrow T : A$ and $A : T \rightarrow G : C$). Overall, the chance of CTA mu-
622 tating to ATG would be:

$$SW^2 \times \lambda^2 \times \mu_{G:C \rightarrow T:A} \times \mu_{A:T \rightarrow G:C}$$

623 Next, we calculated the probability that every nucleotide triplet that is not ATG, mu-
624 tates to ATG. This can happen via one, two, or three nucleotide mutations. The sum of
625 all these mutation probabilities is the probability of ATG gain.

626 Using the same principle we calculated the gain probability of any DNA sequence
627 motif (of any length or composition). We excluded insertions and deletions as a mech-
628 anism of gain of small DNA sequences that we analysed in this study.

629 **2.3 Probability of losing a DNA sequence**

630 We calculated the loss of a DNA sequence motif using the same principle we used for
631 calculating gain probabilities. However, we defined loss probability as a conditional
632 probability, that is we assume that the DNA sequence of interest already exists in the
633 genome. For example, the loss probability of a specific ATG sequence would be the
634 sum of probabilities of ATG mutating to any of the other 63 nucleotide triplets (via one,
635 two, or three nucleotide mutations). We use conditional loss probabilities by default,
636 because usually one is interested in finding out how quickly an existing DNA sequence
637 can erode.

638 We used this method to calculate the loss probability of any DNA sequence motif, and
639 we excluded insertions and deletions from this calculation.

640 **3 Probabilities of finding, gaining, and losing DNA features**

641 Usually, a specific function is encoded in DNA by several DNA sequences. For exam-
642 ple, translation stop is encoded by three codons (TGA, TAG, TAA). We use the term DNA
643 features to mean a set of DNA sequences that are associated with the same function.
644 For every such DNA feature set, there is a complementary set of non-features, that is
645 DNA sequences that are not associated with the feature's function. For example, the
646 non-feature set of stop codons would be all the other 61 codons.

647 The probability that a DNA feature exists, is the sum of probabilities of every DNA
648 sequence in that set ([Methods section 2.1](#)).

649 The probability that a DNA feature is gained via mutations, is the sum of probabilities
650 of every non-feature sequence mutating to any feature sequence. If F denotes the
651 feature set, and $\mu_{y \rightarrow x}$ denotes the probability of a DNA sequence y mutating to a DNA
652 sequence x (see [Methods section 2.2](#)), then:

$$P_{feature-gain} = \sum_{x \in F} \sum_{y \notin F} P_y \times \mu_{y \rightarrow x} \quad (16)$$

653 The probability that a DNA feature is lost via mutations is a conditional probability
 654 that given a feature exists, it mutates to any of the non-feature sequences.

$$P_{feature-loss} = \frac{\sum_{x \in F} \sum_{y \notin F} P_x \times \mu_{x \rightarrow y}}{\sum_{x \in F} P_x} \quad (17)$$

655 Because a feature set usually has many DNA sequences, a mutation can change a fea-
 656 ture sequence such that the resulting sequence is also a part of the feature set. Thus we
 657 defined the probability ($P_{feature-stay}$) that a feature does not erode because of mutations.
 658 Specifically, it is the sum of two probabilities. First is the probability that no mutation
 659 occurs in the DNA sequence (P_0), and the second probability describes the event where
 660 the mutated sequence remains a part of the feature set.

$$P_{feature-stay} = P_0 + \sum_{x \in F} \sum_{\substack{y \in F \\ y \neq x}} \mu_{x \rightarrow y} \quad (18)$$

661 The probability that no mutation occurs (P_0) in a DNA sequence of length k is described
 662 by Poisson distribution.

$$P_0 = 1 - e^{-k\lambda}$$

663 Because the mutation rate is biased (Table 1), the probability that no mutation occurs
 664 in a DNA sequence depends on its composition.

665 The probability that an A or a T mutates (λ_{AT}), is thus described as:

$$\lambda_{AT} = 2 \times \lambda \times (\mu_{A:T \rightarrow T:A} + \mu_{A:T \rightarrow G:C} + \mu_{A:T \rightarrow C:G})$$

666 Likewise, the probability that a G or a C mutates (λ_{GC}) is:

$$\lambda_{GC} = 2 \times \lambda \times (\mu_{G:C \rightarrow A:T} + \mu_{G:C \rightarrow T:A} + \mu_{G:C \rightarrow C:G})$$

667 (Note that the general mutation rate, λ , is an average of λ_{AT} and λ_{GC} .)

668 Thus the probability that a sequence of length k , containing m number of A and T, does
669 not mutate is:

$$P_0 = (1 - e^{-m\lambda_{AT}}) \times (1 - e^{-(k-m)\lambda_{GC}})$$

670 We also calculated all the above-defined probabilities ($P_{feature-gain}$, $P_{feature-loss}$ and $P_{feature-stay}$),
671 using DNA trimer and hexamer frequencies from *D. melanogaster*. In this case the prob-
672 ability of finding a nucleotide sequence depends on the trimer/hexamer distributions
673 instead of GC-content, but the probability of mutational changes are only dependent
674 on mutation bias. Trimers and hexamers would contain codons and polyA signals,
675 respectively.

676 4 Probabilities of finding, gaining, and losing an ORF

677 4.1 Probability of finding an ORF

678 A reading frame is a nucleotide sequence with a length that is a multiple of three. A
679 reading frame that begins with a start codon (ATG), and ends with one of the three stop
680 codons (TAG, TGA, TAA) is an open reading frame (ORF). This necessarily means that
681 there are no stop codons within the sequence. Thus the probability of finding an ORF
682 containing k codons including start and stop codons (P_{ORF}) is:

$$P_{ORF}(k) = P_{ATG} \times P_{stop} \times (1 - P_{stop})^{k-2} \quad (19)$$

683 Here, P_{ATG} and P_{stop} are the probabilities of finding a start codon, and a stop codon by
684 chance, respectively.

685 4.2 Probability of gaining an ORF

686 As we defined in the previous section, an ORF has three requirements (start codon,
687 stop codon, and no premature stop codon in the sequence). Thus an ORF can emerge
688 due to mutations via three mechanisms. In each of these mechanisms, one require-
689 ment is initially absent whereas the other two are present. Then mutations cause the
690 missing requirement to emerge while not disrupting the other two requirements. More
691 specifically, the ORF can be gained via the following three mechanisms:

- 692 1. Gain of a start codon ($P_{ATG-gain}$) while a stop codon continues to exist at the end
693 of a reading frame ($P_{stop-stay}$), and there is no emergence of stop codon within the
694 reading frame ($1 - P_{stop} - P_{stop-gain}$).
- 695 2. Gain of a stop codon ($P_{stop-gain}$), while a start codon continues to exist at the be-
696 ginning of a reading frame ($P_{ATG-stay}$), and there is no emergence of stop codon
697 within the reading frame.
- 698 3. Loss of a premature stop codon, at any of the $k - 2$ codon positions within the
699 reading frame ($P_{stop-gain}$). At the same time start and stop codons remain undis-
700 rupted by mutations, and no stop codon emerges at any of the other $k - 3$ posi-
701 tions.

702 Thus we define the probability of ORF gain ($P_{ORF-gain}$) as:

$$\begin{aligned}
 P_{ORF-gain}(k) = & P_{ATG-gain} \times P_{stop-stay} \times (1 - P_{stop} - P_{stop-gain})^{k-2} \\
 & + P_{ATG-stay} \times P_{stop-gain} \times (1 - P_{stop} - P_{stop-gain})^{k-2} \\
 & + P_{ATG-stay} \times P_{stop-stay} \times P_{stop-loss} \times (k - 2) \times (1 - P_{stop} - P_{stop-gain})^{k-3} \quad (20)
 \end{aligned}$$

703 4.3 Probability of ORF loss

704 ORF can be lost when any of its three requirements are lost. We thus define the condi-
705 tional probability of ORF loss as:

$$P_{ORF-loss}(k) = P_{ATG-loss} + P_{stop-loss} + (k - 2) \times \frac{P_{stop-gain}}{1 - P_{stop}} \quad (21)$$

706 The last term in this equation describes the conditional probability of stop-gain, given
707 the assumption that no stop codon exists within the ORF.

708 4.4 Probability that ORF remains intact

709 We assumed that an ORF of a certain length remains intact if none of the necessary
710 features are lost. However, the ORF sequence can mutate to cause non-synonymous
711 changes in the translated protein sequence. This condition applies to all the three ORF
712 probabilities described above. We define the probability that an ORF remains intact
713 ($P_{ORF-stay}$) as:

$$P_{ORF-stay}(k) = P_{ATG-stay} \times P_{stop-stay} \times (1 - P_{stop} - P_{stop-gain})^{k-2} \quad (22)$$

714 5 Probabilities of finding, gaining, and losing transcription

715 5.1 Probability of transcription

716 In this study we defined a transcription model that is not dependent on specific pro-
717 moter sequences. Instead we defined a probability that transcription is initiated at a
718 genomic locus. To this end, we first extracted the genomic coordinates of intergenic se-
719 quences in *D. melanogaster* genome (release 6.49, FlyBase; Gramates *et al.*, 2022). Using
720 the data on the genomic locations of open chromatin in *D. melanogaster* (Milon *et al.*,
721 2014), we obtained open chromatin regions within intergenic genome. Next, we calcu-
722 lated the fraction of this open-intergenic genome occupied by *cis*- regulatory elements
723 (Gallo *et al.*, 2005; Halfon *et al.*, 2007; Gallo *et al.*, 2010; Rivera *et al.*, 2018; Keränen
724 *et al.*, 2022). A significant proportion of these elements include enhancers. Enhancers,
725 that can promote transcription in their neighboring genomic regions (Andersson and
726 Sandelin, 2020), are also known to promote *de novo* gene emergence (Majic and Payne,
727 2020). Thus we reasoned that proportion of *cis*- regulatory elements in open inter-
728 genic chromatin would roughly approximate the likelihood of transcription initiation
729 ($P_{RNAinit}$). We do not model the exact sequences because data is not yet available on
730 exact sequence requirements and mutational robustness of enhancers and other *cis*-
731 regulatory elements. Thus poly-A signal is the primary sequence determinant of tran-
732 scription probabilities. Relatedly, open intergenic chromatin regions, that are not inter-
733 rupted by *cis*- regulatory elements are the potential sites for *de novo* gene emergence.

734 A poly-A signal marks the end of a transcript. The transcribed DNA region should
735 thus have only one poly-A signal at its end. A transcribed DNA of length l nucleotides,
736 contains $l - 5$ hexamers (6nt long sequences). Each of these hexamers can be a poly-A
737 signal with a probability, P_{polyA} . Probability that none of these hexamers is a poly-A
738 signal is $(1 - P_{polyA})^{l-5}$. Since we focus on protein coding genes in this work, we de-
739 scribe the length of transcript in terms of the length of protein it can encode. Assuming
740 that untranslated regions have no effect on protein synthesis, the length of a transcript
741 in number of nucleotides should be at least three times the length of the encoded pro-
742 tein in number of amino acids. Thus the probability that a DNA region produces a
743 transcript that can harbor an ORF with k codons ($P_{RNA}(k)$) is:

$$P_{RNA}(k) = P_{RNAinit} \times P_{polyA} \times (1 - P_{polyA})^{3k-5} \quad (23)$$

744 5.2 Probability of gaining transcription

745 As we mentioned in the previous section, two requirements need to be met to pro-
746 duce a transcript of any given length. First, a poly-A signal needs to be present at the
747 end of the DNA region to be transcribed, and second, no poly-A signal should exist
748 within this DNA region. Transcription can be gained if one of these requirements are
749 already met and the third required feature emerges due to mutations. At the same time
750 mutations should not destroy the existing features. It is possible that both the required
751 features are missing and they emerge simultaneously due to mutations. However such
752 an event is highly improbable. Thus transcription gain can occur via two mechanisms.
753 In the first mechanism a poly-A signal is gained at the end of a sequence ($P_{polyA-gain}$),
754 while no poly-A sites are present at any $l - 5$ (or $3k - 5$) sites within the sequence, and
755 none emerge ($P_{nopolyA-stay} = 1 - P_{polyA} - P_{polyA-gain}$). In the second mechanism, a poly-A
756 signal is present at the end of the sequence, and is not lost due to mutations. Mean-
757 while, one poly-A signal at any of the $3k - 5$ sites within the sequence is lost ($P_{polyA-loss}$),
758 whereas the other $3k - 6$ sites do not encode a poly-A signal and continue to remain
759 so.

760 The combined probability of these two transcription gain mechanism is defined as:

$$P_{RNA-gain}(k) = P_{RNAinit} \times \left(\begin{array}{l} P_{polyA-gain} \times (P_{nopolyA-stay})^{3k-5} \\ + P_{polyA-stay} \times (P_{nopolyA-stay})^{3k-6} \times (3k - 6) \times P_{polyA-loss} \end{array} \right) \quad (24)$$

761 5.3 Probability of losing transcription

762 Transcription is lost when either the poly-A signal is lost or if one is gained within
763 the transcribed sequence. The probability of transcription loss ($P_{RNA-loss}$), given that
764 already transcription exists, is thus defined as:

$$P_{RNA-loss}(k) = P_{polyA-loss} + (3k - 5) \times P_{polyA-gain} \quad (25)$$

765 5.4 Probability that transcription remains intact

766 Transcription remains intact if neither the poly-A signal is not lost and none is gained
767 within the sequence. Probability of this event ($P_{RNA-stay}$) is defined as:

$$P_{RNA-stay}(k) = P_{RNAinit} \times P_{polyA-stay} \times (P_{nopolyA-stay})^{(3k-5)} \quad (26)$$

768

769 6 Interdependence of transcription and ORF probabilities

770 All the four poly-A signal variants contain a TAA in their sequence (position 3 – 5),
771 which is also a stop codon. Since stop codons cannot exist inside an ORF, presence of an
772 ORF reduces the possible number of poly-A signal sites in a DNA region. Specifically,
773 a poly-A signal can exist in any of the three reading frames but it can overlap with
774 the ORF only in one frame. Thus, out of $3k - 5$ positions for a 6-mer (length of a
775 poly-A signal) in an ORF with k codons, k positions cannot harbor a poly-A signal
776 (Equation 23). Furthermore a poly-A signal cannot overlap with the start codon, thus
777 reducing two more positions from all sites where a poly-A signal can exist (one position
778 is already counted – if a poly-A signal overlaps with the first codon in the second
779 frame, then the second codon becomes a stop codon). Overall, the possible number
780 poly-A sites on an ORF with k codons is $2k - 5$. Thus for a gene to emerge via gain of
781 transcription, the DNA should not contain poly-A sites at these $2k - 5$ positions. As
782 a consequence, the probability of transcription gain for a specific transcript length, is
783 higher when the DNA region contains an ORF than when it does not (Equation 24). We
784 have ignored the probability of poly-A signal overlapping with a stop codon because
785 it makes a very small difference in the probability of RNA gain. This is so because at

786 least one stop codon (TAA) allows a poly-A signal to overlap with it in three frames.

787 Conversely to the effect of an existing ORF on transcription gain, stop codons are less
 788 probable in an existing transcript that does not contain any poly-A signal sequences
 789 within its sequence, than in an untranscribed DNA region. Specifically, the probability
 790 that a 6-mer contains a stop codon in the third position is same as the probability of
 791 finding a stop codon (P_{stop}). If poly-A signals are excluded from these 6-mers, then the
 792 probability of finding a stop codon is $P_{stop} - P_{poly-A}$. Conversely, lack of poly-A signals
 793 will increase the probability of amino acid coding codons and thereby their probability
 794 of mutating into a stop codon. Furthermore, the likelihood that an ORF is undisrupted
 795 my premature nonsense mutations, is lower when the DNA is already transcribed.
 796 Specifically:

$$P_{stop-gain|RNA} = \frac{P_{stop-gain} \times (1 - P_{stop} + P_{poly-A})}{1 - P_{stop}}$$

797 Overall the evolutionary dynamics of transcription and ORF gain, are not independent
 798 of each other. That is, presence of one feature makes the gain of the other more likely.
 799 The dependence becomes more prominent with increasing number of codons in the
 800 ORF of the proto-gene.

801 7 Calculation of fixation probabilities

802 We estimated gene fixation probabilities using Kimura's model (Kimura, 1962). Specific-
 803 ically we defined the probabilities of fixation and extinction of a proto-gene as shown
 804 in Table 2:

Case	Gain of gene (Fixation, P_{fix})	Loss of gene (Extinction, P_{ext})	Selection coefficient (s)
No selection on protogene	$\frac{1}{2N_e}$	$\frac{1}{2N_e}$	–
Selection on protogene	$\frac{1 - e^{-2s}}{1 - e^{-4N_e s}}$	$\frac{1 - e^{2s}}{1 - e^{4N_e s}}$	± 0.001

Table 2: Gene fixation/extinction probabilities with or without selection. Here N_e denotes effective population size and s denotes the selection co-efficient denoting the fitness advantage provided by the protein encoded in the proto-gene (positive valued for a beneficial protein, and negative valued for a deleterious protein).

805 **8 Probability of amino acid substitutions**

806 We calculated the probability of amino acid substitutions based on nucleotide sub-
807 stitution in the corresponding codons. In this case, multiple codons that code for the
808 same amino acid constitute a feature set ([Methods section 3](#)). Similarly when analysing
809 hydrophobic to non-hydrophobic substitutions, all codons that represent hydrophobic
810 amino acids form a feature set (and *vice versa*).

811 **Data availability**

812 We performed all calculations using Julia programming language, and all scripts are
813 freely available on GitHub: [BharatRaviIyengar/DeNovoEvolution](#). Specifically, our
814 model is implemented in the script `DeNovoProb.jl` which in turn uses the script
815 `nucleotidedefuncts.jl` for some basic functions.

816 **Acknowledgments**

817 We thank Pete Czuppon (IEB, University of Münster), for discussions on mathematical
818 models.

819 **References**

- 820 Acevedo, J. M., Hoermann, B., Schlimbach, T., and Teleman, A. A. 2018. Changes in global
821 translation elongation or initiation rates shape the proteome via the Kozak sequence. *Scien-
822 tific Reports*, 8(1): 4018.
- 823 Andersson, R. and Sandelin, A. 2020. Determinants of enhancer and promoter activities of
824 regulatory elements. *Nature Reviews Genetics*, 21(2): 71–87.
- 825 Behrens, S. and Vingron, M. 2010. Studying the evolution of promoter sequences: A waiting
826 time problem. *Journal of Computational Biology*, 17(12): 1591–1606.
- 827 Berg, J. M., Tymoczko, J. L., and Stryer, L. 2002. *Biochemistry*. W.H. Freeman, New York.
- 828 Blevins, W. R., Ruiz-Orera, J., Messeguer, X., and others 2021. Uncovering de novo gene birth
829 in yeast using deep transcriptomics. *Nature Communications*, 12(1): 604.

- 830 Box, G. 1979. Robustness in the strategy of scientific model building. In *Robustness in Statistics*,
831 pages 201–236. Academic Press.
- 832 Bucciantini, M., Giannoni, E., Chiti, F., and others 2002. Inherent toxicity of aggregates implies
833 a common mechanism for protein misfolding diseases. *Nature*, 416(6880): 507–511.
- 834 Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y., and Greenleaf, W. J. 2013. Transposition
835 of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-
836 binding proteins and nucleosome position. *Nature Methods*, 10(12): 1213–1218.
- 837 Cano, A. V., Rozhoňová, H., Stoltzfus, A., McCandlish, D. M., and Payne, J. L. 2022. Mutation
838 bias shapes the spectrum of adaptive substitutions. *Proceedings of the National Academy of
839 Sciences*, 119(7): e2119720119.
- 840 Carvunis, A.-R., Rolland, T., Wapinski, I., and others 2012. Proto-genes and de novo gene birth.
841 *Nature*, 487(7407): 370–374.
- 842 Choe, Y.-J., Park, S.-H., Hassemer, T., and others 2016. Failure of RQC machinery causes protein
843 aggregation and proteotoxic stress. *Nature*, 531(7593): 191–195.
- 844 Churchman, L. S. and Weissman, J. S. 2012. Native elongating transcript sequencing (NET-seq).
845 *Current Protocols in Molecular Biology*, 98(1).
- 846 Clark, M. B., Amaral, P. P., Schlesinger, F. J., and others 2011. The reality of pervasive transcrip-
847 tion. *PLoS Biology*, 9(7): e1000625.
- 848 Dayhoff, M., Schwartz, R., and Orcutt, B. 1978. Model of evolutionary change in proteins. In
849 *Atlas of protein sequence and structure*, volume 5, pages 345–352. National Biomedical Research
850 Foundation, Silver Spring MD.
- 851 Dill, K. A. 1985. Theory for the folding and stability of globular proteins. *Biochemistry*, 24(6):
852 1501–1509.
- 853 Fernández-Moreno, M. A., Farr, C. L., Kaguni, L. S., and Garesse, R. 2007. *Drosophila*
854 *melanogaster* as a model system to study mitochondrial biology. In *Methods in Molecular
855 Biology*, pages 33–49. Humana Press.
- 856 Gallo, S. M., Li, L., Hu, Z., and Halfon, M. S. 2005. *REDfly*: a regulatory element database for
857 *Drosophila*. *Bioinformatics*, 22(3): 381–383.
- 858 Gallo, S. M., Gerrard, D. T., Miner, D., and others 2010. *REDfly* v3.0: toward a comprehen-
859 sive database of transcriptional regulatory elements in drosophila. *Nucleic Acids Research*,
860 39(Database): D118–D123.
- 861 Gardini, S., Cheli, S., Baroni, S., and others 2016. On nature’s strategy for assigning genetic
862 code multiplicity. *PLOS ONE*, 11(2): e0148174.

- 863 Gerstein, M. B., Bruce, C., Rozowsky, J. S., and others 2007. What is a gene, post-ENCODE?
864 History and updated definition. *Genome Research*, 17(6): 669–681.
- 865 Gonnet, G. H., Cohen, M. A., and Benner, S. A. 1992. Exhaustive matching of the entire protein
866 sequence database. *Science*, 256(5062): 1443–1445.
- 867 Gramates, L. S., Agapite, J., Attrill, H., and others 2022. FlyBase: a guided tour of highlighted
868 features. *Genetics*, 220(4).
- 869 Grandchamp, A., Czuppon, P., and Bornberg-Bauer, E. 2023. High turnover of *de novo* tran-
870 scripts in *Drosophila melanogaster*. *bioRxiv*.
- 871 Haberle, V. and Stark, A. 2018. Eukaryotic core promoters and the functional basis of transcrip-
872 tion initiation. *Nature Reviews Molecular Cell Biology*, 19(10): 621–637.
- 873 Halfon, M. S., Gallo, S. M., and Bergman, C. M. 2007. REDfly 2.0: an integrated database of
874 cis-regulatory modules and transcription factor binding sites in drosophila. *Nucleic Acids
875 Research*, 36(Database): D594–D598.
- 876 Harris, K. and Nielsen, R. 2014. Error-prone polymerase activity causes multinucleotide muta-
877 tions in humans. *Genome Research*, 24(9): 1445–1454.
- 878 Hartl, F. U. 2017. Protein misfolding diseases. *Annual Review of Biochemistry*, 86(1): 21–26.
- 879 Henikoff, S. and Henikoff, J. G. 1992. Amino acid substitution matrices from protein blocks.
880 *Proceedings of the National Academy of Sciences*, 89(22): 10915–10919.
- 881 Hershberg, R. and Petrov, D. A. 2010. Evidence That Mutation Is Universally Biased towards
882 AT in Bacteria. *PLOS Genetics*, 6(9): 1–13.
- 883 Hochberg, G. K. A., Liu, Y., Marklund, E. G., and others 2020. A hydrophobic ratchet entrenches
884 molecular complexes. *Nature*, 588(7838): 503–508.
- 885 Iyengar, B. R., Choudhary, A., Sarangdhar, M., and others 2014. Non-coding RNA interact to
886 regulate neuronal development and function. *Frontiers in Cellular Neuroscience*, 8: 47.
- 887 Jones, D. T., Taylor, W. R., and Thornton, J. M. 1992. The rapid generation of mutation data
888 matrices from protein sequences. *Bioinformatics*, 8(3): 275–282.
- 889 Keeling, D. M., Garza, P., Nartey, C. M., and Carvunis, A.-R. 2019. Philosophy of Biology: The
890 meanings of ‘function’ in biology and the problematic case of *de novo* gene emergence. *eLife*,
891 8: e47014.
- 892 Keränen, S. V. E., Villahoz-Baleta, A., Bruno, A. E., and Halfon, M. S. 2022. REDfly: An inte-
893 grated knowledgebase for insect regulatory genomics. *Insects*, 13(7): 618.
- 894 Kim, Y., Sidney, J., Pinilla, C., Sette, A., and Peters, B. 2009. Derivation of an amino acid simi-
895 larity matrix for peptide:MHC binding and its application as a bayesian prior. *BMC Bioinforma-
896 tics*, 10(1).

- 897 Kimura, M. 1962. ON THE PROBABILITY OF FIXATION OF MUTANT GENES IN a POPU-
898 LATION. *Genetics*, 47(6): 713–719.
- 899 Kosiol, C., Holmes, I., and Goldman, N. 2007. An empirical codon model for protein sequence
900 evolution. *Molecular Biology and Evolution*, 24(7): 1464–1479.
- 901 Kozak, M. 1986. Point mutations define a sequence flanking the AUG initiator codon that
902 modulates translation by eukaryotic ribosomes. *Cell*, 44(2): 283–292.
- 903 Kyte, J. and Doolittle, R. F. 1982. A simple method for displaying the hydropathic character of
904 a protein. *Journal of Molecular Biology*, 157(1): 105–132.
- 905 Lange, J. D., Bastide, H., Lack, J. B., and Pool, J. E. 2021. A population genomic assessment of
906 three decades of evolution in a natural *Drosophila* population. *Molecular Biology and Evolution*,
907 39(2).
- 908 Le, S. Q. and Gascuel, O. 2008. An improved general amino acid replacement matrix. *Molecular*
909 *Biology and Evolution*, 25(7): 1307–1320.
- 910 Lee, Y. C. G., Ventura, I. M., Rice, G. R., and others 2019. Rapid evolution of gained essential de-
911 velopmental functions of a young gene via interactions with other essential genes. *Molecular*
912 *Biology and Evolution*, 36(10): 2212–2226.
- 913 Lenhard, B., Sandelin, A., and Carninci, P. 2012. Metazoan promoters: emerging characteristics
914 and insights into transcriptional regulation. *Nature Reviews Genetics*, 13(4): 233–245.
- 915 Long, M., Betrán, E., Thornton, K., and Wang, W. 2003. The origin of new genes: glimpses from
916 the young and old. *Nature Reviews Genetics*, 4(11): 865–875.
- 917 Lopes, R., Agami, R., and Korkmaz, G. 2017. GRO-seq, a tool for identification of transcripts
918 regulating gene expression. In *Methods in Molecular Biology*, pages 45–55. Springer New York.
- 919 Majic, P. and Payne, J. L. 2020. Enhancers facilitate the birth of *de novo* genes and gene integra-
920 tion into regulatory networks. *Molecular Biology and Evolution*, 37(4): 1165–1178.
- 921 Merchant, S. S., Prochnik, S. E., Vallon, O., and others 2007. The *Chlamydomonas* Genome
922 Reveals the Evolution of Key Animal and Plant Functions. *Science*, 318(5848): 245–250.
- 923 Milon, B., Sun, Y., Chang, W., and others 2014. Map of open and closed chromatin domains in
924 drosophila genome. *BMC Genomics*, 15(1).
- 925 Monroe, J. G., Srikant, T., Carbonell-Bejerano, P., and others 2022. Mutation bias reflects natural
926 selection in *Arabidopsis thaliana*. *Nature*, 602(7895): 101–105.
- 927 Näsval, J., Sun, L., Roth, J. R., and Andersson, D. I. 2012. Real-time evolution of new genes by
928 innovation, amplification, and divergence. *Science*, 338(6105): 384–387.

- 929 Neme, R. and Tautz, D. 2016. Fast turnover of genome transcription across evolutionary time
930 exposes entire non-coding DNA to de novo gene emergence. *eLife*, 5.
- 931 Noderer, W. L., Flockhart, R. J., Bhaduri, A., and others 2014. Quantitative analysis of mam-
932 malian translation initiation sites by FACS-seq. *Molecular Systems Biology*, 10(8): 748.
- 933 Ohta, T. and Kimura, M. 1971. Amino acid composition of proteins as a product of molecular
934 evolution. *Science*, 174(4005): 150–153.
- 935 Omotajo, D., Tate, T., Cho, H., and Choudhary, M. 2015. Distribution and diversity of ribosome
936 binding sites in prokaryotic genomes. *BMC Genomics*, 16(1): 604.
- 937 Prabh, N. and Rödelsperger, C. 2019. *De Novo*, divergence, and mixed origin contribute to
938 the emergence of orphan genes in *Pristionchus* nematodes. *G3 Genes|Genomes|Genetics*, 9(7):
939 2277–2286.
- 940 Proudfoot, N. J. 2011. Ending the message: poly(A) signals then and now. *Genes & Development*,
941 25(17): 1770–1782.
- 942 Richard, P. and Manley, J. L. 2009. Transcription termination by nuclear RNA polymerases.
943 *Genes & Development*, 23(11): 1247–1269.
- 944 Rivera, J., Keränen, S. V. E., Gallo, S. M., and Halfon, M. S. 2018. REDfly: the transcriptional
945 regulatory element database for *Drosophila*. *Nucleic Acids Research*, 47(D1): D828–D834.
- 946 Santangelo, T. J. and Artsimovitch, I. 2011. Termination and antitermination: RNA polymerase
947 runs a stop sign. *Nature Reviews Microbiology*, 9(5): 319–329.
- 948 Schmid, M. and Jensen, T. H. 2018. Controlling nuclear RNA levels. *Nature Reviews Genetics*,
949 19(8): 518–529.
- 950 Schmitz, J. and Bornberg-Bauer, E. 2017. Fact or fiction: updates on how protein-coding genes
951 might emerge de novo from previously non-coding DNA . *F1000Research*, 6(57).
- 952 Schneider, A., Cannarozzi, G. M., and Gonnet, G. H. 2005. Empirical codon substitution matrix.
953 *BMC Bioinformatics*, 6(1).
- 954 Schrider, D. R., Houle, D., Lynch, M., and Hahn, M. W. 2013. Rates and Genomic Consequences
955 of Spontaneous Mutational Events in *Drosophila melanogaster*. *Genetics*, 194(4): 937–954.
- 956 Shen, S., Kai, B., Ruan, J., and others 2006. Probabilistic analysis of the frequencies of amino
957 acid pairs within characterized protein sequences. *Physica A: Statistical Mechanics and its
958 Applications*, 370(2): 651–662.
- 959 Statello, L., Guo, C.-J., Chen, L.-L., and Huarte, M. 2021. Gene regulation by long non-coding
960 RNAs and its biological functions. *Nature Reviews Molecular Cell Biology*, 22(2): 96–118.

- 961 Stewart, M. 2019. Polyadenylation and nuclear export of mRNAs. *Journal of Biological Chemistry*,
962 294(9): 2977–2987.
- 963 Tautz, D. and Domazet-Lošo, T. 2011. The evolutionary origin of orphan genes. *Nature Reviews*
964 *Genetics*, 12(10): 692–702.
- 965 Vakirlis, N., Hebert, A. S., Opulente, D. A., and others 2017. A molecular portrait of de novo
966 genes in yeasts. *Molecular Biology and Evolution*, 35(3): 631–645.
- 967 Van Oss, S. B. and Carvunis, A.-R. 2019. De novo gene birth. *PLOS Genetics*, 15(5): 1–23.
- 968 Whelan, S. and Goldman, N. 2001. A general empirical model of protein evolution derived
969 from multiple protein families using a maximum-likelihood approach. *Molecular Biology and*
970 *Evolution*, 18(5): 691–699.
- 971 Wimley, W. C. and White, S. H. 1996. Experimentally determined hydrophobicity scale for
972 proteins at membrane interfaces. *Nature Structural & Molecular Biology*, 3(10): 842–848.
- 973 Witt, E., Benjamin, S., Svetec, N., and Zhao, L. 2019. Testis single-cell RNA-seq reveals the
974 dynamics of de novo gene transcription and germline mutational bias in drosophila. *eLife*, 8.
- 975 Wood, V., Gwilliam, R., Rajandream, M.-A., and others 2002. The genome sequence of
976 *Schizosaccharomyces pombe*. *Nature*, 415(6874): 871–880.
- 977 Zhang, Z. and Gerstein, M. 2003. Patterns of nucleotide substitution, insertion and deletion in
978 the human genome inferred from pseudogenes. *Nucleic Acids Research*, 31(18): 5338–5348.
- 979 Zhao, L., Saelao, P., Jones, C. D., and Begun, D. J. 2014. Origin and spread of *de Novo* genes in
980 *Drosophila melanogaster* populations. *Science*, 343(6172): 769–772.