

1 Random peptides rich in small and disorder-  
2 promoting amino acids are less likely to be  
3 harmful

4 Luke Kosinski<sup>\*,††,1</sup>, Nathan Aviles<sup>†,††</sup>, Kevin Gomez<sup>‡</sup>, Joanna Masel<sup>§</sup>

5 \* Department of Molecular and Cellular Biology, University of Arizona, Tucson, AZ 85721

6 † Graduate Interdisciplinary Program in Statistics, University of Arizona, Tucson, AZ 85721

7 ‡Graduate Interdisciplinary Program in Applied Math, University of Arizona, Tucson, AZ 85721

8 § Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, 85721

9 †† These authors contributed equally to this work.

10 <sup>1</sup>Present address: Critical Path Institute, Tucson, AZ 85718

11 ORCID IDs: 0000-0002-8146-5955 (L. K.); 0000-0001-9998-4406 (N. A.); 0000-0002-6356-0318

12 (K. G.); 0000-0002-7398-2127 (J. M.)

- 13 Running title: Amino acid effects on fitness
- 14 Keywords: Experimental evolution, evolvability, fitness estimation, preadaptation, de novo
- 15 gene birth
- 16 Corresponding author: Joanna Masel, Department of Ecology and Evolutionary Biology,
- 17 University of Arizona, 1041 E Lowell St., Tucson, AZ 85721. E-mail: [masel@email.arizona.edu](mailto:masel@email.arizona.edu)
- 18 +1(520)626-9888

## 19 Abstract

20 Proteins are the workhorses of the cell, yet they carry great potential for harm via misfolding  
21 and aggregation. Despite the dangers, proteins are sometimes born *de novo* from non-coding  
22 DNA. Proteins are more likely to be born from non-coding regions that produce peptides that  
23 do little to no harm when translated than from regions that produce harmful peptides. To  
24 investigate which newborn proteins are most likely to “first, do no harm”, we estimate fitnesses  
25 from an experiment that competed *Escherichia coli* lineages that each expressed a unique  
26 random peptide. A variety of peptide metrics significantly predict lineage fitness, but this  
27 predictive power stems from simple amino acid frequencies rather than the ordering of amino  
28 acids. Amino acids that are smaller and that promote intrinsic structural disorder have more  
29 benign fitness effects. We validate that the amino acids that indicate benign effects in random  
30 peptides expressed in *E. coli* also do so in an independent dataset of random N-terminal tags in  
31 which it is possible to control for expression level. The same amino acids are also enriched in  
32 young animal proteins.

33

## 34 Introduction

35 Proteins are the workhorses of the cell, but they are dangerous. For example, the polypeptide  
36 backbone is the key structural feature of amyloids, putting all proteins at risk of forming  
37 insoluble aggregates (Chiti and Dobson 2017), and most proteins are expressed at or just  
38 beyond their solubility limits (Vecchi, et al. 2020). Despite these dangers, new protein-coding  
39 genes are nevertheless born *de novo* from essentially random sequences (McLysaght and  
40 Guerzoni 2015; Van Oss and Carvunis 2019; Vakirlis, Carvunis, et al. 2020). To be beneficial  
41 enough for *de novo* birth, a random peptide must first do no serious harm, i.e. it must not be  
42 detrimental to the basic functioning of a cell. Here we quantify the degree to which, and the  
43 summary statistics via which, a random peptide's propensity for harm can be predicted.

44 Neme et al. (2017) competed over 2 million *Escherichia coli* lineages, each containing a  
45 plasmid designed to express a unique random peptide, and tracked lineage frequencies over  
46 four days using deep DNA sequencing. This study has been criticized for providing too little  
47 support for the beneficial nature of the top candidates (Weisman and Eddy 2017; Knopp and  
48 Andersson 2018). But these criticisms do not detract from using the dataset to identify  
49 statistical predictors of serious harm versus relatively benign effect. Neme et al. (2017) used a  
50 strong promoter, so evaluation is of tolerance to high expression. Some fitness differences  
51 might be due to variation in expression e.g. due to auto-downregulation at the RNA level  
52 (Knopp and Andersson 2018) - we will return to this point in the last portion of the Results.  
53 Here we pursue analyses based on the hypothesis that the properties of the peptides  
54 contribute to variation in fitness among lineages.

55            Conveniently, computational predictors from peptide sequences alone are available for  
56            some properties, such as intrinsic structural disorder (ISD) and aggregation propensity. Because  
57            insoluble proteins have been implicated in toxicity and disease (Chiti and Dobson 2017) and  
58            peptides with high ISD are less prone to forming insoluble aggregates (Linding, et al. 2004;  
59            Angyan, et al. 2012), we hypothesize that highly disordered peptides are least likely to be  
60            strongly deleterious. Random sequences with high predicted disorder are well-tolerated *in vivo*  
61            (Tretyachenko, et al. 2017). Existing mouse (Wilson, et al. 2017) and *Drosophila* (Heames, et al.  
62            2020) proteins, which are the product of evolution, are predicted from their amino acid  
63            sequences to be more disordered than what would be translated from intergenic controls.

64            Younger protein-coding sequences should be particularly constrained to first do no  
65            harm, as they have had little time to evolve more sophisticated harm-avoidance strategies (Foy,  
66            et al. 2019). In support of the idea that high ISD is an accessible way to avoid harm, young  
67            animal and fungal domains (James, et al. 2021) and genes (Wilson, et al. 2017; Foy, et al. 2019;  
68            James, et al. 2021), and novel overprinted viral genes (Willis and Masel 2018) have higher  
69            predicted disorder than their older counterparts. Some studies have found that putative de  
70            novo protein candidates in *Saccharomyces* yeasts have lower rather than higher ISD (Carvunis,  
71            et al. 2012; Basile, et al. 2017; Vakirlis, et al. 2018), but this could be an artifact of  
72            proportionately greater inclusion of non-genes within the younger age classes. When Wilson et  
73            al. (2017) reanalyzed Carvunis et al.'s (2012) "proto-genes" of different ages, using more  
74            rigorous criteria to exclude non-genes from the data, the direction of the ISD trend was  
75            reversed. The same reversal of trend following a quality filter was also found by Vakirlis et al.  
76            (2018).

77           How much do amino acid frequencies matter compared to the order in which those  
78 amino acids are arranged? Prior research on young genes has suggested that high predicted ISD  
79 in that context is driven primarily by amino acid frequencies, with amino acid order playing a  
80 more minor role (Wilson, et al. 2017). Fortunately, the dataset of Neme et al. (2017) is large  
81 enough to look at the frequencies of each amino acid as predictors, rather than assume that  
82 existing prediction programs such as IUPred (Dosztányi, et al. 2005; Meszaros, et al. 2018) or  
83 Tango (Fernandez-Escamilla, et al. 2004; Linding, et al. 2004; Rousseau, et al. 2006) integrate all  
84 information about both amino acid frequencies and ordering in the best possible way. We can  
85 then test whether such programs have additional ability to predict peptide fitness, above and  
86 beyond the influence of amino acid frequencies. In doing so, we can estimate the relative roles  
87 of amino acid frequencies versus amino acid ordering in predicting fitness, as well as determine  
88 which amino acids have which effects.

89           Here we investigate the degree to which amino acid frequencies and amino acid  
90 ordering can predict the fitness effects of random peptides, and if so, which properties are  
91 most predictive. We also investigate whether the properties that help random peptides avoid  
92 harm in *E. coli* are also enriched in young eukaryotic proteins. With our work, we hope to  
93 further our understanding of how peptides avoid harm.

## 94   Methods

## 95 Data retrieval

96 Neme et al. (2017) performed seven experiments where *E. coli* lineages, each with a plasmid  
97 containing a unique random peptide, were grown and tracked using deep DNA sequencing. We  
98 downloaded sequencing counts from Dryad at <http://dx.doi.org/10.5061/dryad.6f356>, and  
99 obtained amino acid and nucleotide sequences directly from Rafik Neme. Experiment 7 was by  
100 far the largest with over 4 million reads, more than five times larger than the 2<sup>nd</sup> largest  
101 experiment and over 1.2 million reads more than all other experiments combined. Experiment  
102 7 contained all the peptides that the other six experiments classified as “increasing” or  
103 “decreasing,” and more. Small datasets from these other six experiments yield limited  
104 information because of the need to model changing mean fitness in a population, including not  
105 just the tracked lineages but also cells with an empty vector (see Estimating lineage fitness from  
106 random peptide sequencing counts section). We therefore chose to restrict our analysis to  
107 experiment 7. Experiment 7 consists of the numbers of reads of each random peptide sequence  
108 in 5 replicate populations of *E. coli* at 4 time points. We assume that fitness is identical across  
109 replicates, so we summed across all 5 replicates to obtain a total number of reads for each  
110 polypeptide at each time point.

111 Following Neme et al. (2017), we took the 1061 peptides out of over one million that  
112 had  $\geq 5$  reads across all 5 replicates of experiment 7. Neme et al. (2017) used this cutoff because  
113 it is not possible to infer fitness with any reasonable resolution for individual peptides with  
114 fewer than five reads. The dramatic nature of this data reduction is unsurprising, firstly because  
115 each initial unique peptide was present in only one copy, and secondly because most peptides

116 are likely deleterious. We note therefore that our analyzed subset of peptides with at least five  
117 reads are certainly non-lethal, and likely less deleterious than the average random peptide.  
118 Nonetheless, we achieved enough resolution to distinguish between more and less harmful  
119 peptides, with remarkably large effect sizes considering the restricted fitness range.

120 We further excluded the six peptides that, while meeting the criterion of  $\geq 5$  reads, had  
121 all of those reads at the same timepoint, leaving 1055 peptides for analysis.

122

### 123 Estimating lineage fitness from random peptide sequencing counts

124 The expected number of reads  $\lambda_{it}$  of peptide  $i$  at times  $t=1,2,3,4$  was modeled as:

$$125 \quad \lambda_{it} = N_t p_{i0} \prod_{k=1}^t \frac{\omega_i}{W_{k-1}},$$

126 where  $N_t$  is the observed total number of reads,  $p_{i0}$  is the initial frequency of peptide  $i$  at the  
127 beginning of the experiment (prior to the round of selection used to produce the first measured  
128 timepoint  $t = 1$ ),  $\frac{\omega_i}{W_t}$  is the fitness of bacteria with peptide  $i$  at time  $t$  (i.e. their propensity to  
129 contribute to the next time point), and  $W_k$  is population mean fitness at time  $k$ , including  
130 bacteria containing empty vectors for which we have no direct count data.

131 The likelihoods of observed peptide counts were estimated from this expectation and  
132 two different error models. A Poisson distribution, which captures sampling error alone, was  
133 used to generate our initial estimates of  $p_{i0}$ ,  $\omega_i$ , and  $W_k$  (collectively yielding  $\lambda_{it}$ ) because it is



134 analytically tractable. Under a Poisson error function, the likelihood of observing  $n_{it}$  reads of  
135 peptide  $i$  at time  $t$  is

136 
$$f_{Poisson}(n_{it}|\lambda_{it}) = \frac{\lambda_{it}^{n_{it}} e^{-\lambda_{it}}}{n_{it}!}.$$

137 To also capture variance inflation  $\kappa$  due to PCR amplification, we used a negative binomial  
138 distribution in the Polya form:

139 
$$f_{NBP}(n_{it}|\lambda_{i,t},\kappa) = \left( \frac{\Gamma\left(n_{it} + \frac{\lambda_{i,t}}{\kappa - 1}\right)}{n_{it}! \Gamma\left(\frac{\lambda_{i,t}}{\kappa - 1}\right)} \right) \left(\frac{1}{\kappa}\right)^{\frac{\lambda_{i,t}}{\kappa - 1}} \left(1 - \frac{1}{\kappa}\right)^{n_{it}}$$

140 where  $\Gamma(\cdot)$  is the gamma function. We used the initial estimates of  $p_{i0}$ ,  $\omega_i$ , and  $W_k$  to  
141 numerically fit the negative binomial model. For the specifics of fitting the Poisson and negative  
142 binomial models, see Supporting Information. Weights were calculated, for use in downstream  
143 linear models, from this likelihood inference procedure, as the inverse of Fisher information  
144 (see Supporting Information).

145 An existing software package for estimating lineage fitness from sequencing counts is  
146 Fit-Seq (Li, et al. 2018), which captures the amplification of PCR error through a more  
147 sophisticated distribution for the number of reads that is derived in the supplementary  
148 information of Levy et al. (2015). However, Fit-Seq assumes that mean fitness is a simple  
149 average of all measured lineages' fitness, requiring all individuals to be tagged and measured.  
150 But Neme et al.'s (2017) experiment included lineages carrying an empty plasmid, i.e. with the  
151 selectable marker but no random peptide. Worse, the proportion of cells with an empty vector  
152 can be presumed to increase over time. In the absence of a reliable way to directly quantify

153 cells with empty vectors, we instead consider mean population fitness over time to be a set of  
154 independent parameters to be fitted.

### 155 Clustering non-independent sequences

156 Upon visual inspection, we found that some peptide sequences were extremely similar, with  
157 only one or two amino acid differences; these data points will not contain independent  
158 information about the relationship between sequence and fitness. To account for non-  
159 independence, we clustered peptides by their Hamming distance, and either took only the  
160 peptide whose fitness had the highest weight within its cluster, or took weighted means within  
161 clusters, or included cluster in our regression models as a random effect term. Single-link  
162 clustering with Hamming distance cutoffs of 6 to 29 amino acids all produced an identical set of  
163 646 clusters for our 1055 peptides. The largest cluster had 228 random peptides, and the  
164 second largest had only 13. The vast majorities of clusters contained only 1 sequence (Dataset  
165 S1). A few peptides had mutations in their non-random regions; these mutations were counted  
166 in our Hamming distance measurements.

167       Such similar sequences are highly unlikely to arise by chance if the peptides were truly  
168 random;  $20^{50} \approx 10^{65}$  peptides are possible, far more than the  $\sim 2 \times 10^6$  observed. Because we  
169 analyze only peptides with at least 5 reads, replicated sequencing error is an unlikely cause. We  
170 see the same nearly-identical sequences appearing in every experimental replicate, suggesting  
171 either that mutations occurred during Neme et al.'s (2017) initial growth phase, or that the  
172 "random" peptides synthesized for the experiment are not entirely random. We note that

173 construction of the “random” peptide library involved ligations of a smaller set of “seed”  
174 sequences, introducing non-randomness at this stage.

### 175 Predictors of fitness

176 All peptides are exactly 65 amino acids long with 50 amino acids of random sequence, so there  
177 was no need to control for length.

#### 178 *GC content*

179 Many amino acid sequences mapped to several possible nucleotide sequences, as part of the  
180 same problem of mutation or non-random construction discussed above. To calculate one GC  
181 content for each random peptide, we calculated a simple average of GC content across all the  
182 nucleotide sequences in the dataset that map to the peptide with the largest weight in the  
183 cluster.

184 To calculate GC content for the over two million peptides with at least one sequencing  
185 read, we took a simple average of the GC content from the random portion of the peptides.

#### 186 *Disorder*

187 Protein disorder was measured using IUPred2 (Dosztányi, et al. 2005; Meszaros, et al. 2018) for  
188 amino acid sequences, and using disorder propensity (Theillet, et al. 2013) for individual amino  
189 acids. IUPred2 returns an ISD score between zero and one for each amino acid in a sequence,  
190 with higher scores indicating greater intrinsic disorder. To calculate an ISD score for each  
191 random peptide, we took the average of the scores for the whole sequence (i.e. including non-  
192 random parts). We used a square root transform because it produced a more linear relationship

193 with fitness than no transform. All measurements referring to ISD or IUPred used IUPred2  
194 except  $\Delta$ ISD, which used the original IUPred program – differences between the two are  
195 minimal (Meszaros, et al. 2018).

196 Disorder propensity gives each amino acid a score based on the frequency it is found in  
197 disordered proteins relative to ordered proteins (Theillet, et al. 2013). The disorder propensity  
198 score for a peptide was determined by averaging the disorder propensity scores for the amino  
199 acids in the random region. When we use the disorder propensity metric, we explicitly refer to  
200 it as “disorder propensity” and not as “ISD.”

#### 201 *Aggregation propensity*

202 Tango (Fernandez-Escamilla, et al. 2004; Linding, et al. 2004; Rousseau, et al. 2006) returns an  
203 aggregation score for each amino acid in a sequence. At least five sequential amino acids with a  
204 score greater than or equal to five indicates an aggregation-prone region. We scored peptide  
205 aggregation propensity as the number of amino acids within regions scored as aggregation-  
206 prone, including contributions from non-random regions.

#### 207 *Solubility*

208 CamSol (Sormanni, et al. 2015) returns a solubility score for each amino acid in a sequence, as  
209 well as a simple average of all scores for a sequence, which CamSol calls a “solubility profile.”  
210 We used the solubility profile of the full sequences, including non-random regions.

#### 211 *Amino acid frequencies*

212 We counted frequencies among the 50 amino acids in the random portion of each peptide.

213 The values for all the above predictors for each peptide are listed in Dataset S1.

## 214 Statistics

215 All statistical tests were carried out in R version 3.6.3 (R Core Team 2019), with figures  
216 generated using “ggplot2” (Wickham 2016). Weighted linear mixed models were implemented  
217 using the “lmer” function from the “lme4” package (Bates, et al. 2015), with cluster as a  
218 random effect. See Supporting Information for details, including justification of a log-transform  
219 for fitness. When  $R^2$  values were needed, we instead averaged peptides within the same cluster  
220 into a combined datapoint, allowing us to avoid the use of random effect term. We calculated  
221  $R^2$  and adjusted  $R^2$  values using the base R “lm” function. Adjusted  $R^2$  is a modification of  $R^2$  to  
222 penalize additional predictors, and is calculated using the formula:

$$223 \quad R_{adj}^2 = 1 - (1 - R^2) \frac{n-1}{n-p-1},$$

224 where  $n$  are the number of data points and  $p$  are the number of predictors. Raw P-values are  
225 reported unless otherwise noted, i.e. without correction for multiple comparisons.

## 226 Data and code availability

227 All code and supplemental tables are available on GitHub at  
228 <https://github.com/MaselLab/RandomPeptides>. The original Neme et al. (2017) data can be  
229 found at Dryad <http://dx.doi.org/10.5061/dryad.6f356>, and the original sequences are available  
230 at the European Nucleotide Archive (ENA) under the project number PRJEB19640.

## 231 Results

## 232 Estimating the fitness effects of random peptides

233 Assessing predictors of the fitness effects of random peptides requires those fitness effects to  
234 be measured accurately and precisely. Neme et al. (2017) tracked lineage frequencies over four  
235 days, and categorized a peptide as increasing or decreasing in frequency by comparing the DNA  
236 sequencing counts of day 4 to day 1 using DESeq2 (Love, et al. 2014).

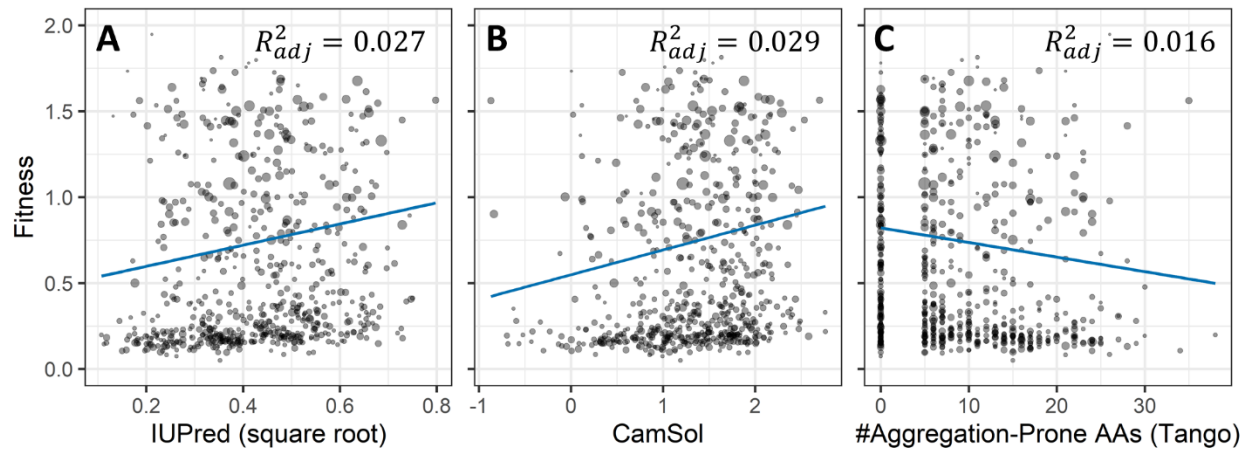
237 We reanalyze the same data, instead using a custom maximum likelihood framework  
238 (see Materials & Methods) to quantitatively estimate “fitness” and its associated confidence  
239 interval / weight. “Fitness” here refers to allele frequency changes over an entire cycle of  
240 population growth and dilution, rather than per generation. Our method classifies peptides  
241 quantitatively rather than qualitatively. It accounts for the fact that mean population fitness  
242 increases over the four days (see Materials and Methods). Our use of all available data within  
243 an appropriate maximum likelihood framework should make our method more sensitive and  
244 specific for identifying benign vs harmful peptides (see Supplementary Text).

245 Note that some peptides are pseudoreplicates (see Materials & Methods). There were  
246 646 total clusters, of which there was statistical support for increases in frequency for the  
247 highest-weighted peptide in 138 clusters, and for decreases in 488 clusters. Some of our  
248 statistics use cluster as a random effect within a linear mixed model. When fixed-effect models  
249 are used, such as to generate interpretable  $R^2$  values, we collapse each cluster into a single  
250 pseudo-datapoint with value given by the weighted mean and weight given by the sum of  
251 weights.

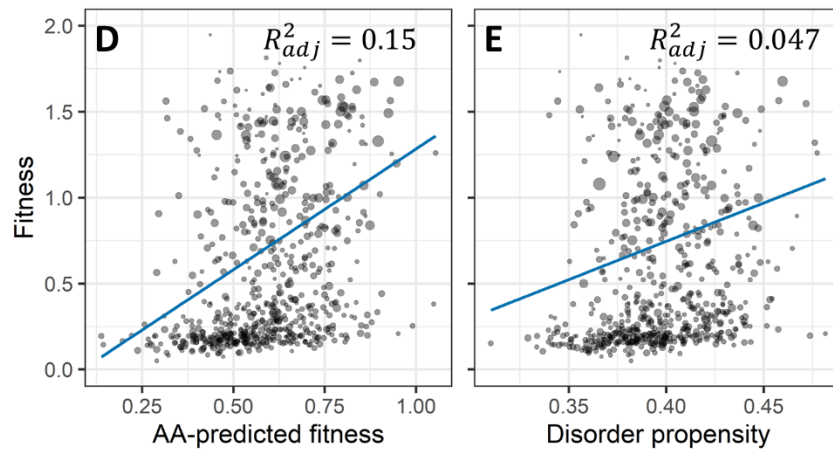
252 Most predictive power stems from amino acid frequencies rather than amino acid  
253 order

254 We estimated peptide disorder using several metrics that contain information both about  
255 amino acid frequencies and about their order: IUPred as an estimate of intrinsic structural  
256 disorder (Dosztányi, et al. 2005; Meszaros, et al. 2018), CamSol as an estimate of water  
257 solubility (Sormanni, et al. 2015), and Tango as an estimate of general aggregation propensity  
258 (Fernandez-Escamilla, et al. 2004; Linding, et al. 2004; Rousseau, et al. 2006). Fewer than 6% of  
259 the random peptides have a predicted transmembrane helix (Dataset S1) from TMHMM (Krogh,  
260 et al. 2001), so our choice of these predictors is guided by our assumption that the random  
261 peptides are predominantly located in the cytosol. Having a predicted transmembrane helix did  
262 not in itself predict random peptide fitness effects ( $P = 0.2$ , likelihood ratio test relative to  
263 mixed model with only the intercept as a fixed effect). In contrast, each of our cytosol-  
264 solubility-inspired metrics significantly predicted random peptide fitness (Fig. 1A – 1C), with  
265 effects in the predicted direction (more disorder and more solubility are good, more  
266 aggregation propensity is bad). Adjusted  $R^2$  values for IUPred, CamSol, and Tango are 0.027,  
267 0.029, 0.016, respectively. Another aggregation predictor, Waltz (Maurer-Stroh, et al. 2010),  
268 that specializes in  $\beta$  aggregates, was in the right direction but did not quite meet statistical  
269 significance ( $P = 0.06$ ).

## Metrics that include amino acid frequency + order information



## Frequency information only



270

271 Fig. 1. **Many metrics predict peptide fitness effects, but most predictive power comes from**  
272 **amino acid frequencies.** Three metrics that combine information on both amino acid  
273 frequencies and amino acid order ((A) IUPred, (B) CamSol, and (C) Tango), and two that contain  
274 only amino acid frequency information ((D) 19 custom weights on amino acid frequencies and  
275 (E) independently estimated disorder propensities used as weights on amino acid frequencies),  
276 each significantly predict peptide fitness on their own ( $P = 7 \times 10^{-4}$ , 0.003, 0.02,  $5 \times 10^{-6}$ , and  $9 \times$   
277  $10^{-7}$ , respectively, likelihood ratio test in mixed model compared to intercept-only model). Each  
278 point ( $n = 646$ ) shows a cluster of sequences with similar amino acid sequences (see Methods



279 for more details), and the area displayed for each point is proportional to summed weights  
280 across that cluster. Blue lines are fixed-effect weighted linear regressions of cluster fitness on  
281 the x-axis predictor, where clusters are collapsed to a single pseudo-datapoint by their  
282 weighted average and weights are sums within each cluster. Metrics that include both  
283 frequency and order information fail to outperform frequency-only based metrics, as shown by  
284 regression slopes (blue lines) and adjusted  $R^2$  values (top right of each figure panel). Adjusted  $R^2$   
285 is calculated as  $R_{adj}^2 = 1 - (1 - R^2) \frac{n-1}{n-p-1}$ , where  $n$  is the number of data points and  $p$  is the  
286 number of degrees of freedom in the predictor. Note that in part D the predictor (model-  
287 predicted fitness) is a composite of 19 degrees of freedom that have all been trained on the  
288 dataset, so care should be taken in comparing its blue regression line to that of the other  
289 panels, each of which has a predictor with only one degree of freedom – this problem does not  
290 apply to comparisons of adjusted  $R^2$  values. Seven clusters with fitness greater than 2 are not  
291 shown here for ease of visualization; a complete y-axis is shown in supplemental fig. 1. Log-  
292 transforming fitness would remove high fitness skew, but creates systematic heteroscedasticity,  
293 and so was not done (supplemental fig. 2). The lack of systematic heteroscedasticity can be  
294 seen here in the form of similar point size across fitness values.

295

296       Next we asked whether these sophisticated metrics offer additional predictive power  
297 beyond mere amino acid frequencies, in the light of prior work on young genes in which little  
298 additional predictive power was found (Wilson, et al. 2017). To do this, we fit a model of fitness  
299 predicted by amino acid frequencies, measured from counts of each amino acid in each

300 peptide's random region (Fig. 1D), and compared its performance to predictors that  
301 incorporate ordering information (Figs. 1A-C). The amino acid frequency-only model was a  
302 significant predictor of fitness ( $P = 4.5 \times 10^{-6}$ , likelihood ratio test compared to an intercept-only  
303 mixed model). It is also more biologically predictive than other metrics, with adjusted  $R^2 = 0.15$   
304 (adjusted to account for the number of predictors used) being far greater than the values of  
305 0.027, 0.029, and 0.016 found in Figs 1A-1C. Another, non-adjusted, way to look at biological  
306 effect size is the far steeper blue line in Fig. 1D than in Figs. 1A-1C. Statistically, when the  
307 frequencies of each of the twenty amino acids are used as predictors (Fig. 1D), then IUPred,  
308 CamSol, and Tango drop out of the model ( $P = 0.2, 0.2, \text{ and } 0.3$ , respectively, likelihood ratio  
309 test in mixed model, see Supplemental Table S1), suggesting that their predictive power in Figs.  
310 1A-1C came largely from being metrics of amino acid frequencies. These results are surprising:  
311 one might expect sophisticated metrics that incorporate both amino acid frequencies and order  
312 information to offer more predictive power and explain a greater range of fitness than simple  
313 amino acid frequencies, yet they fail to do so.

314 Our Fig. 1D model using the frequencies of the 20 amino acids involves 19 degrees of  
315 freedom, while the other metrics we examine involve only one. This makes it inappropriate to  
316 compare the slopes of the blue lines, although adjusted  $R^2$  values can still be compared, and the  
317 fact that the other metrics drop out of a combined model is also informative. We also  
318 investigated a one degree of freedom model of amino acid frequencies, in which relative  
319 weights were specified in advance by a disorder propensity metric that assigns each amino acid  
320 a score based on how frequently it is found in known disordered versus ordered proteins  
321 (Theillet, et al. 2013). Average disorder scores over each peptide's random region significantly

322 predicted random peptide fitness effects in a linear mixed model (Fig. 1E,  $P = 9 \times 10^{-7}$ , likelihood  
323 ratio test compared to an intercept-only model). The effect size on predicted fitness from the  
324 10% to the 90% quantiles of disorder propensity is 0.49 to 0.70, and the adjusted  $R^2$  for the  
325 disorder propensity model 0.047. For comparison to other predictors with a single degree of  
326 freedom, the largest effect size model that incorporates both amino acid frequency and order  
327 information was IUPred with an effect size from 0.51 to 0.69, and the best adjusted  $R^2$  model  
328 was CamSol with 0.029. This further suggests that predictive power resides with amino acid  
329 frequencies, not order information.

330 To understand whether order information has additional predictive power beyond that  
331 of amino acid frequencies, we next investigated a metric of ISD that is comprised of only order  
332 information. This can be calculated as the excess IUPred score of the real peptide in comparison  
333 to the average IUPred score of a set of hypothetical peptides in which the order of the amino  
334 acids has been randomly scrambled; this metric was previously found to be elevated in younger  
335 mouse genes (Wilson, et al. 2017). However, adding this  $\Delta$ ISD metric to our model with amino  
336 acid frequencies as predictors did not significantly improve the model ( $P = 0.2$ ). This further  
337 supports our conclusion that amino acid ordering plays only a minor role compared to amino  
338 acid frequencies in the fitness effects of the random peptides examined here.

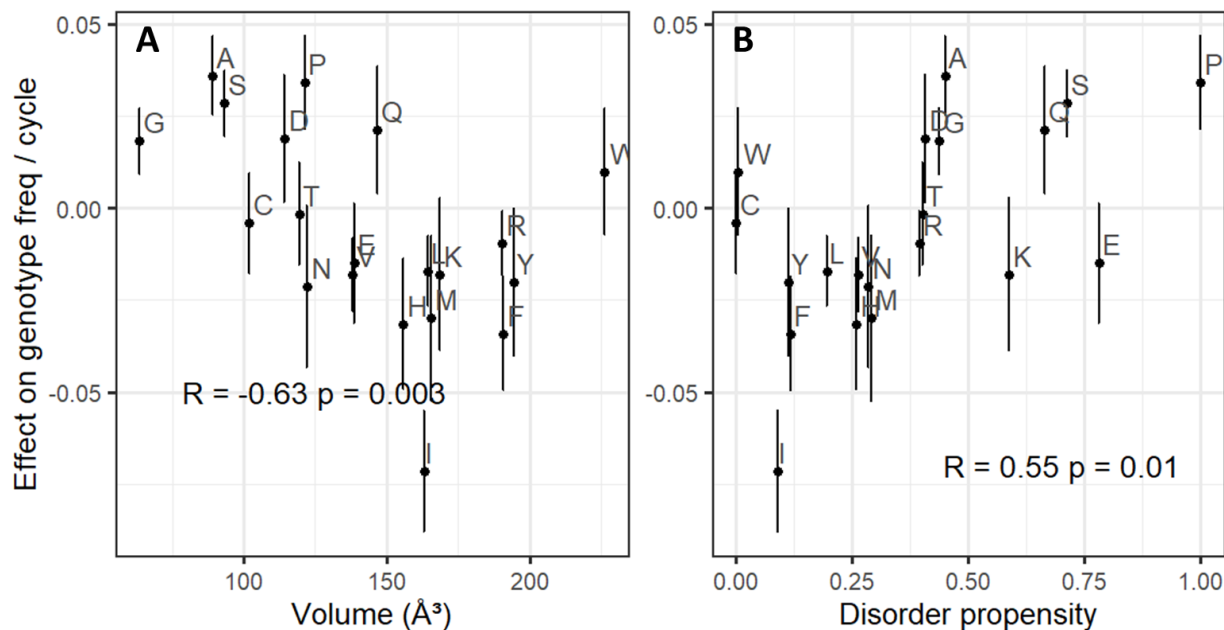
339

### 340 Small and disorder-promoting amino acids predict benign fitness effects

341 Next we quantify the statistical effect of each of the 20 amino acids on fitness. Naively, we  
342 could take the associated slope coefficient in a multiple regression model, which represents the

343 change in fitness when one amino acid is gained. But in a peptide of fixed length, one amino  
344 acid cannot be gained without another amino acid being lost. We therefore instead calculate  
345 the marginal fitness effect of each amino acid on fitness (see supplementary text and Table S2,  
346 displayed in fig. 2, y-axis), representing the effect of gaining that amino acid and losing a  
347 randomly selected alternative.

348 Amino acids with smaller volumes (Tsai, et al. 1999) and higher disorder propensities  
349 (Theillet, et al. 2013) tend to have higher marginal fitness effects (fig. 2A and 2B;  $P = 0.01$  for  
350 both disorder propensity and volume, likelihood ratio test for dropping either term from a  
351 weighted regression of marginal effect on both volume and disorder propensity). Volume and  
352 disorder propensity together explain over half the weighted variation in marginal fitness effect  
353 (weighted adjusted  $R^2 = 0.52$ ). Other properties of amino acids, such as stickiness (Levy, et al.  
354 2012), relative solvent accessibility (Tien, et al. 2013), amino acid cost in *E. coli* (Akashi and  
355 Gojobori 2002), and isoelectric point (Liu, et al. 2004) did not provide significant explanatory  
356 power on top of disorder propensity and volume (all  $P > 0.1$ , likelihood ratio test).



357

358 **Fig. 2. Amino acids that are small and are associated with disorder promote higher fitness.**

359 The y-axis shows each amino acid's marginal effect on fitness, which is the change in fitness  
360 when one amino acid of the focal type replaces one randomly chosen amino acid of a different  
361 type in a random peptide (see Supporting Information). Error bars are +/- one standard error. P-  
362 values and correlation coefficients come from weighted Pearson's correlations, where weights  
363 for marginal effects are calculated as  $1 / \text{s.e. (marginal fitness effect)}^2$ , and volume and disorder  
364 propensity are unweighted.

365

366 Tryptophan is an outlier for amino acid effects on fitness, with a slightly positive effect  
367 on fitness despite both its large volume and its underrepresentation in disordered regions (fig.  
368 2). Removing tryptophan from a weighted regression model of volume and disorder propensity  
369 predicting marginal effect increases the weighted adjusted  $R^2$  from 0.52 to 0.68. Tryptophan,  
370 encoded only by UGG, is nearly 60% more common among peptides with at least 5 sequence

371 reads than we expect from the 58% GC content of our dataset. Together with the confidence  
372 interval for its marginal fitness effect including 1, this provides further evidence that tryptophan  
373 is not harmful, making it a distinct outlier, for reasons that are not clear to us.

374 Isoleucine also stands out, as even more harmful than expected by its large size and  
375 order propensity. Isoleucine's harmful effects may be exacerbated by its role in amyloid  
376 formation. For example, familial amyloid cardiomyopathy is most commonly caused by a valine  
377 to isoleucine mutation (Jacobson, et al. 1997; Dubrey, et al. 2015), suggesting that isoleucine  
378 has potential to form dangerous amyloids where other hydrophobic amino acids do not.  
379 Isoleucine, valine, and leucine are all hydrophobic amino acids with a branched carbon, but only  
380 raised isoleucine levels are associated with a higher risk of Alzheimer's disease (Larsson and  
381 Markus 2017), further suggesting that isoleucine may be especially prone to amyloid formation.

382

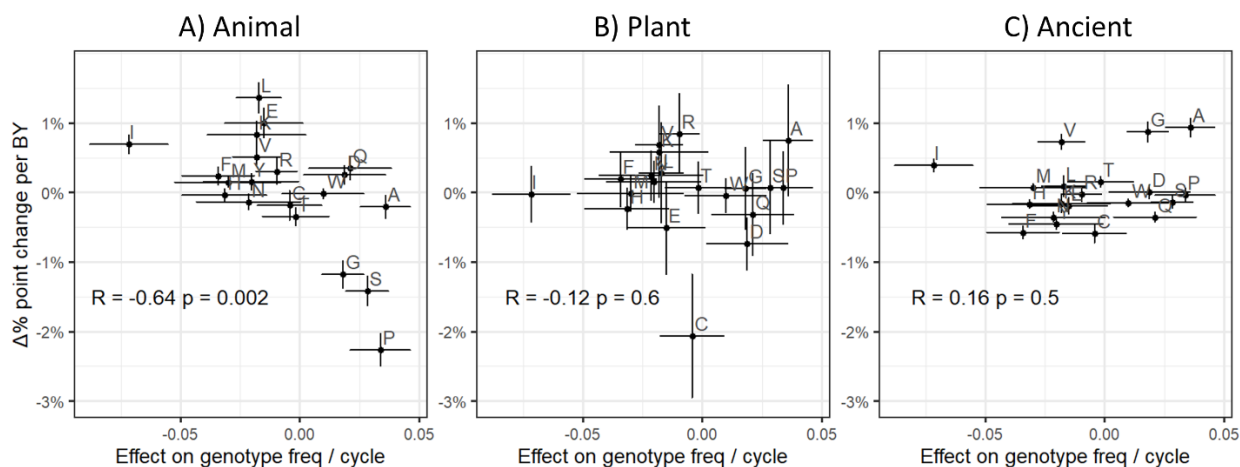
383 Young animal sequences are enriched for amino acids that increase fitness in  
384 random peptides

385 As discussed in the Introduction, young domains have higher predicted ISD than their  
386 older counterparts. One hypothesis to explain this observation is that in order to be successfully  
387 born *de novo*, a protein sequence is especially constrained to first do no harm (Wilson, et al.  
388 2017). However, the "phylostratigraphy" approach of assigning ages to genes is contentious.  
389 Detecting homologs is more difficult for fast-evolving sequences, which may be erroneously  
390 scored as young (Alba and Castresana 2007; Moyers and Zhang 2015, 2016). Disordered  
391 proteins tend to be fast evolving (Chen, et al. 2011), suggesting that highly disordered genes

392 could be misclassified as young because of their fast evolutionary rate. If the amino acid  
393 enrichments of higher fitness random peptides match the amino acid enrichments of young  
394 genes, this would be evidence that the *de novo* gene birth process, rather than homology  
395 detection bias alone, causes trends in protein properties as a function of apparent gene age.

396 To test this, we took the slopes of amino acid frequencies with protein domain age from  
397 James et al. (2021), as quantified across over 400 eukaryotic species. As predicted, amino acids  
398 that are good for random peptides are enriched among the youngest animal Pfams (fig. 3A).  
399 This prediction was not, however, supported for trends among recent plant domains (fig. 3B)  
400 nor among ancient (fig. 3C) domains older than 2.1 billion years. Plant and ancient trends  
401 reflect a *de novo* gene birth process that enriches for the most abundant amino acids in their  
402 respective lineages, such as cysteine, rather than for amino acids that promote ISD (James, et  
403 al. 2021). It is interesting that we find that ISD still predicts harmless in *E. coli*, even though  
404 we do not find evidence it shaped *de novo* gene birth in its distant ancestors. We also note that  
405 ISD does shape recent *de novo* gene birth in viruses (Willis and Masel 2018).

406



407

408 **Fig. 3. Purportedly young animal Pfams are enriched for amino acids that predict high fitness**  
409 **in random peptides.** The y-axis represents how the frequency of each amino acid depends on  
410 the age of the sequence in billion years (BY), estimated as a linear regression slope for non-  
411 transmembrane Pfam domains (James, et al. 2021). Frequency is in number of percentage  
412 points, e.g. a difference in glutamic acid content of 5% vs. 6% is a difference of one percentage  
413 point. The x-axis shows each amino acid's marginal effect on fitness, which is the change in  
414 fitness when one amino acid of the focal type replaces one randomly chosen amino acid of a  
415 different type in a random peptide (see Supporting Information). Error bars are +/- one  
416 standard error. Fitness effects predict A) animal, but not B) plant, or C) ancient (older than 2.1  
417 billion years) Pfam phylostratigraphy slopes. Correlation coefficients and P-values come from  
418 weighted Pearson correlations. Note that the P-value for animal phylostratigraphy slopes vs  
419 marginal effects survives a conservative Bonferroni correction ( $P = 0.002 < 0.05/3 = 0.017$ ).

420

421

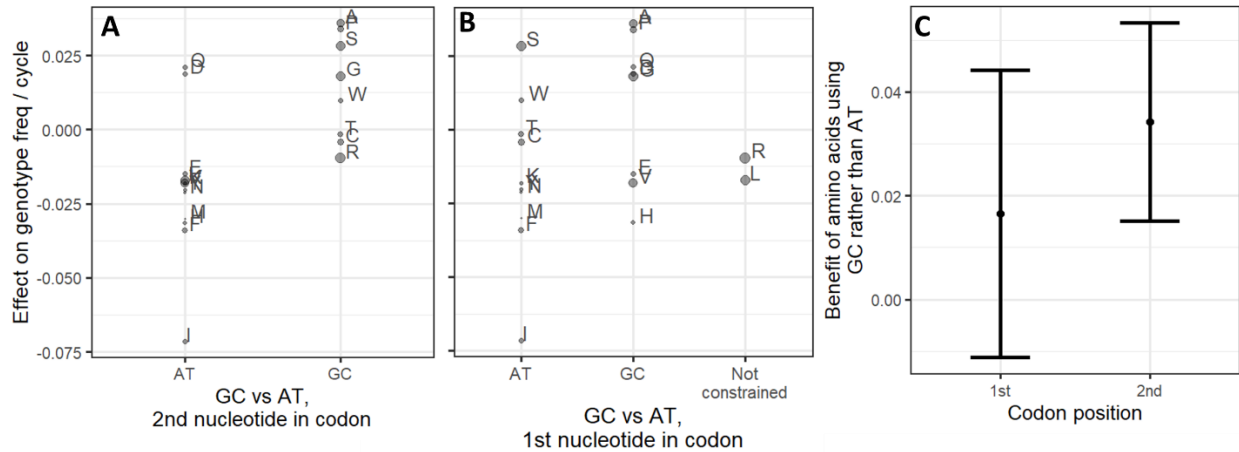
422 **Fitness is better predicted by amino acid frequencies than by GC content**

423 Long et al. (2018) proposed that selection acts directly on GC content, perhaps due to the three  
424 hydrogen bonds of G-C pairs. Amino acids encoded by Gs and Cs tend to promote higher ISD  
425 (Angyan, et al. 2012), making it difficult to distinguish between selection for high GC content  
426 and selection for disorder-promoting amino acids. To attempt to distinguish between the two,  
427 we compare amino acids that always have G or C to those that always have A or T, at both the  
428 first and second nucleotide positions in the codon. If selection were for GC nucleotides, we



429 would expect GC to predict high marginal amino acid fitness effects at both positions. But if  
430 results are dramatically different at the two positions, this would show that it is selection on  
431 amino acid content that drives GC as a correlated trait. Results are statistically significant in the  
432 predicted direction at the second position (fig. 4A,  $P = 0.001$ , weighted Welch's t-test), and in  
433 the predicted direction but not statistically significant at the first (fig. 4B,  $P = 0.2$ ). The effect  
434 size of GC content on fitness could not be statistically distinguished between the first and  
435 second position (fig. 4C), with wide and hence inconclusive error bars.

436 Linear models are compatible with partially independent contributions of both amino  
437 acid frequencies and GC content to harm avoidance. GC content is a statistically significant  
438 predictor of fitness by itself ( $P = 6 \times 10^{-11}$ , likelihood ratio test for nested fixed-effect models  
439 relative to intercept-only model). However, the weighted adjusted  $R^2$  of 0.06 for GC content is  
440 much lower than the weighted adjusted  $R^2$  of 0.15 ( $P = 10^{-18}$ ) for full amino acid frequency  
441 information, suggesting it explains less of the variation than amino acid frequencies. Adding GC  
442 content to the amino acid frequencies-only model offers a modest improvement ( $P = 0.004$ ,  
443 weighted adjusted  $R^2$  values improves from 0.15 to 0.16), while adding amino acid frequencies  
444 to a GC content only model offers a notably larger improvement ( $P = 10^{-11}$ , weighted adjusted  
445  $R^2$  improves from 0.06 to 0.16). These weighted adjusted  $R^2$  values suggest that while there  
446 may be some direct selection on GC content, the effect of amino acid frequencies appears to be  
447 well beyond what can be explained by GC content.



448

449 **Fig. 4. Amino acids that are constrained to use Gs and Cs tend to have higher marginal effects**

450 **on fitness than those constrained to use As and Ts.** The difference is significant for constraints

451 at the second nucleotide position of a codon (A) ( $P = 0.001$ , weighted Welch's t-test), but not at

452 the first (B) ( $P = 0.2$ ). Point area is proportional to weight, which is calculated as  $1 /$

453  $s.e.(\text{marginal fitness effect})^2$ , as described in Supporting Information. The y-axis is the same as

454 the fig. 2 y-axis and fig. 3 x-axis. C) The mean advantage of amino acids constrained to use GC

455 rather than constrained to use AT is not distinguishable in size between the first and second

456 codon positions. Y-axis gives the difference in the two weighted means of marginal fitness

457 effects from A) and B). Error bars represent 95% confidence intervals on the difference

458 between the means (calculated as  $\text{difference} \pm t_{\text{crit}} \times se$ ), where  $t_{\text{crit}} \approx 2.1$  is the critical value of

459 the t-statistic with the appropriate degrees of freedom. Weighted Welch's t-test statistic and

460 the corresponding standard error of the difference in means were calculated using the

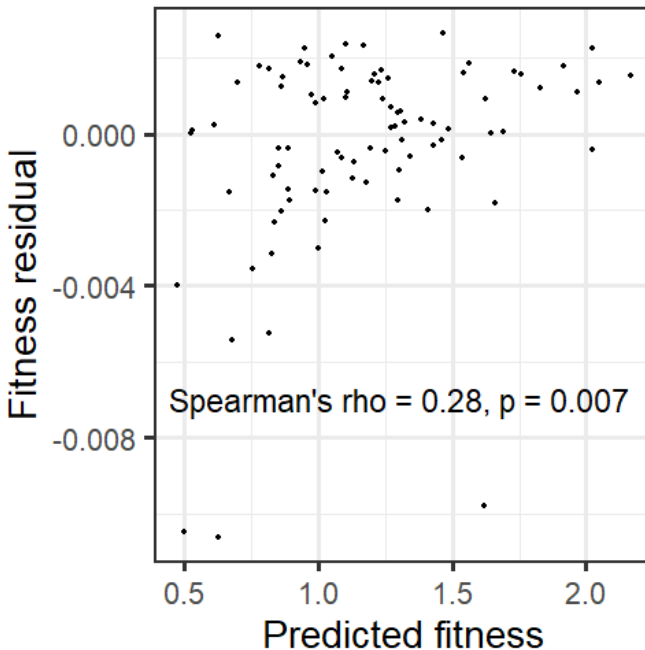
461 "wtd.t.test" function from the "weights" R package, version 1.0.1.

462

463 The same amino acids predict benign fitness effects in random N-terminal tags

464 The degree to which benign effects are due to low expression of a random peptide, vs. benign  
465 effects of the peptide once expressed, remains unclear. We therefore tested the ability of our  
466 amino-acid-frequencies-only model, trained on the data of Neme et al. (2017), to predict  
467 residual fitness effects in a dataset that controls for peptide expression level. Goodman et al.  
468 (2013) tagged the N-prime end of green fluorescent protein (GFP) with 137 different short  
469 random sequences (11 amino acids long), allowing random peptide expression level to be  
470 measured via fluorescence. Frumkin et al. (2017) measured the fitness effects of these random  
471 peptide-tagged GFPs in *E. coli* using FitSeq (Li, et al. 2018). For 89 of them, Frumkin et al. (2017)  
472 were able to calculate a “fitness residual” based on the deviation from the fitness expected  
473 from the level of GFP expression. Note that while this fitness residual controls for expression  
474 level, it still contains the cost of inefficient expression in addition to the fitness effect of the  
475 peptide itself. Frumkin et al. (2017) found that low fitness residuals were associated with  
476 hydrophobic and expensive-to-synthesize amino acids. These findings are consistent with our  
477 own estimates of direct peptide effects, as hydrophobic amino acids tend to be order-prone  
478 (Linding, et al. 2004; Angyan, et al. 2012), and amino acid volume is highly correlated with  
479 synthesis cost in *E. coli* (Pearson’s correlation coefficient = 0.85,  $P = 2 \times 10^{-6}$ , cost for amino acid  
480 synthesis in *E. coli* taken from (Akashi and Gojobori 2002)). Indeed, predicted fitness values for  
481 Frumkin et al.’s (2017) N-terminal tags were significantly correlated with their actual fitness  
482 residuals (fig. 5). The consistency between our results and the findings of Frumkin et al. (2017),  
483 who control for peptide expression level, provides an external validation of our results and

484 suggests that our findings are unlikely to be due to differences in peptide expression levels  
485 alone.



486

487 **Fig. 5. Fitness predictions trained on the random peptides of Neme et al. (2017) also work for**  
488 **short random tags attached to the N-terminus of GFP.** Predicted fitness comes from our amino  
489 acid frequencies-only mixed model. “Fitness residuals” of N-terminal tags are from Frumkin et  
490 al. (2017), and represent the difference between the fitness of the construct and the expected  
491 fitness from expression level.  $n = 89$ .

## 492 Discussion

493 We found that, while many metrics of peptide properties have some ability to predict the  
494 fitness effects of random peptides expressed in *E. coli*, most predictive power stems from  
495 amino acid frequencies. Simply knowing how many of which amino acids are present in these

496 random peptides can account for 15% of the variance in fitness among lineages, and adding  
497 more predictors to account for amino acid order fails to add more predictive power. This  
498 indicates both the success of our statistical method for minimizing the noise in our fitness  
499 estimates, and that mere amino acid frequencies without amino acid order can be informative  
500 of peptide fitness effects. Amino acids that are small and promote disorder predict high fitness  
501 in *E. coli*, and align with those that are enriched in young protein domains in animals.

502         Most studies of random peptides have focused on finding peptides that have specific  
503 binding or function (e.g. Kaiser, et al. 1987; Keefe and Szostak 2001; Frulloni, et al. 2009). Some  
504 were motivated as proof-of-concept that random peptides can exhibit properties of native  
505 proteins, such as folding (Davidson and Sauer 1994; Chiarabelli, et al. 2006; LaBean, et al. 2011)  
506 and being soluble (Priyambada, et al. 1996). Others focus on how to increase the percentage of  
507 native-like random peptides, e.g. by showing that more hydrophilic random peptide libraries  
508 have a higher percentage of stable and soluble peptides (Davidson, et al. 1995). Our work has a  
509 different intent, identifying properties that make a peptide less likely to be harmful. Neme et  
510 al.'s (2017) experiment was suitable for this purpose because it used a large library of peptides  
511 with diverse properties, competed lineages growing under permissive conditions, and  
512 measured relative growth rates (i.e. fitness). In contrast, a study design such as that of Knopp et  
513 al. (2019), who selected random peptides that rescue viability in the presence of antibiotics, is  
514 less suitable for our purposes because so few peptides, including harm-avoiding peptides, are  
515 viable. Neme et al.'s (Neme, et al. 2017) study was also convenient because all peptides were  
516 the same length – 65 amino acids with 50 amino acids of random sequence – allowing us to  
517 neglect length in our analysis.

518           Having a higher proportion of random peptides do no harm is expected to increase the  
519 success rate of future screens for peptide with specific properties. Nucleotide sequences with  
520 high %GC content tend to encode peptides with more benign fitness effects, suggesting that  
521 higher %GC should be used in future random peptide libraries. However, very high GC content  
522 will yield low complexity sequences, which our predictor has not been trained on. The marginal  
523 fitness effects of each amino acid might be different in this very different context.

524           While the library used by Neme et al. (2017) was designed to have equal frequencies of  
525 each nucleotide in the random region, and thus 50% GC content, the over two million random  
526 peptides that had at least one sequencing read had a GC content of ~59% in their random  
527 portion. The mean GC content of the peptide clusters we analyzed (see Materials and Methods)  
528 was similar, at ~58%, with higher fitness peptides within this group having still higher %GC, as  
529 discussed in the Results. The enrichment from 50% GC to ~59% GC might be because many  
530 lower GC content sequences were so harmful that lineages that carried them went extinct prior  
531 to detection via sequencing. Note that it might also reflect a bias toward GC in sequencing  
532 methods (Benjamini and Speed 2012; Choudhari and Grigoriev 2017) – a bias that affects all  
533 time points equally and so should not affect our fitness estimates.

534           Long et al. (2018) proposed that there is direct selection for high GC content, as  
535 evidenced in part by a preference for amino acids with G or C at the second position of codons,  
536 in excess of that predicted from mutation accumulation experiments. Our findings cannot  
537 exclude this hypothesis, but show stronger selection on amino acid frequencies, selection that  
538 is capable of driving increased GC content in coding regions as a correlated trait. In intergenic  
539 regions, elevated %GC is likely driven mostly by GC-biased gene conversion. However, elevated

540 GC content could also be due, at least in part, to selection on peptides from non-coding regions  
541 translated by error (Rajon and Masel 2011; Wilson and Masel 2011). Selection on translation  
542 errors is for example strong enough to shape non-coding sequences beyond stop codons in  
543 *Saccharomyces cerevisiae* (Kosinski and Masel 2020).

544 Fitness effects in Neme et al. (2017) might not be directly caused by peptide properties  
545 alone but instead by the effect of both nucleotide and peptide properties on expression (Knopp  
546 and Andersson 2018), with lower expression being less harmful. For example, auto-  
547 downregulation at the mRNA level can cause significant difference in expression among  
548 peptides, despite identical promoters. However, the properties we find to be predictive, such as  
549 disorder and amino acid size, are not *a priori* related to auto-downregulation of mRNA in wild-  
550 type *E. coli*, making the latter an unlikely explanation for our findings.

551 While driven by amino acid frequencies, our findings are still consistent with the  
552 hypothesis that peptides with low structural disorder tend to be harmful. Disorder-promoting  
553 amino acids may help a peptide remain soluble even if unfolded. Small amino acids also tend to  
554 be benign, perhaps because they are hydrophobic enough to promote some amount of folding  
555 but flexible enough to avoid too much hydrophobic residue exposure.

556 Our findings suggest that the easiest way to avoid harm is through disorder and small  
557 size, but do not rule out other strategies that rely on capacity for folding. Indeed, BCS4, a *de*  
558 *novo* evolved protein in *Saccharomyces cerevisiae*, has a hydrophobic core and is capable of  
559 folding (Bungard, et al. 2017). Vakirlis et al. (2020) found that *de novo* proteins can emerge as

560 transmembrane proteins, which need to be lipid soluble, presumably requiring different harm-  
561 avoidance strategies than peptides that are located in the cytosol.

562         The correlation between the extent to which an amino acid is enriched in young animal  
563 protein domains and its marginal fitness effect in random peptides in *E. coli* is intriguing, and  
564 consistent with a body of literature that *de novo* gene birth favors protein disorder. What is  
565 more, our ability to externally validate animal phylostratigraphy slopes against random  
566 peptides in *E. coli* provides additional support that these slopes represent more than mere bias,  
567 in contrast to suggests that all patterns are due to homology detection bias (Alba and  
568 Castresana 2007; Moyers and Zhang 2015, 2016). That is, if phylostratigraphy trends were due  
569 to an artifact such as homology detection bias, such an artifact would be unlikely to bias our  
570 random peptide analysis in the same direction.

571         Plants have different trends in amino acid frequencies as a function of sequence age  
572 than animals do, with young genes seeming to prefer readily available amino acids, rather than  
573 amino acids that promote ISD (James, et al. 2021). This could be because: 1) plants are less  
574 susceptible to harm from random peptides, 2) other properties, such as amino acid availability,  
575 drive the emergence of *de novo* genes in plants, or 3) the plant data lack the resolution needed  
576 to identify a correlation with the properties studied here. We do not have the ability to  
577 differentiate between these three possibilities here.

578         Nevertheless, our finding of consistency between what is benign in *E. coli* and what is  
579 benign in animals suggests the possibility of a deep concordance in what makes a peptide  
580 harmful between two apparently disparate branches of life. The forces that drive protein birth



581 therefore appear to share a key similarity between bacteria and Animalia. Monod once  
582 suggested that what is true in *E. coli* must also be true in elephants; our work suggests that this  
583 may apply to the properties that tend to make peptides less harmful. To modify Monod's  
584 famous quote, what is harmful in *E. coli* is also harmful in elephants, but not necessarily in  
585 eucalyptus.

586         A major idea in our understanding of proteins is that form – that is, the fold that is  
587 determined by the exact sequence of amino acids – determines function and thus fitness.  
588 However, for these random peptides in *E. coli*, the amino acid content but not the sequence in  
589 which they occur is the main determinant of benign vs harmful effects. Random peptides likely  
590 exist as a diverse ensemble of structural states, but the same is increasingly acknowledged to  
591 be true of functional proteins. While the ordering of amino acids in functional proteins no  
592 doubt plays a role, perhaps mere amino acid frequencies are also more important than once  
593 thought in this context too, especially in structurally disordered protein regions.

## 594 [Acknowledgements](#)

595 This work was supported by the John Templeton Foundation (39667, 60814) and the National  
596 Institutes of Health (GM-104040, T32GM-008659, T32GM-084905). We thank Rafik Neme and  
597 Diethard Tautz for sharing their data with us and for graciously answering all our questions  
598 regarding their analyses, Dvir Schirman and Tzachi Pilpel for sharing their data with us, Joe  
599 Watkins for helpful discussions about our likelihood estimation procedure, and Catherine  
600 Weibel for driving the GC content analysis forward.

601 Works cited

- 602 Akashi H, Gojobori T. 2002. Metabolic efficiency and amino acid composition in the proteomes  
603 of Escherichia coli and Bacillus subtilis. Proceedings of the National Academy of Sciences  
604 of the United States of America 99:3695-3700.
- 605 Alba MM, Castresana J. 2007. On homology searches by protein Blast and the characterization  
606 of the age of genes. BMC Evolutionary Biology 7:53.
- 607 Angyan AF, Perczel A, Gaspari Z. 2012. Estimating intrinsic structural preferences of *de novo*  
608 emerging random-sequence proteins: Is aggregation the main bottleneck? FEBS Letters  
609 586:2468-2472.
- 610 Basile W, Sachenkova O, Light S, Elofsson A. 2017. High GC content causes orphan proteins to  
611 be intrinsically disordered. PLoS computational biology 13:e1005375.
- 612 Bates D, Maechler M, Bolker B, Walker S. 2015. Fitting linear mixed-effects models using lme4.  
613 Journal of Statistical Software 67:1-48.
- 614 Benjamini Y, Speed TP. 2012. Summarizing and correcting the GC content bias in high-  
615 throughput sequencing. Nucleic Acids Research 40:e72.
- 616 Bungard D, Copple JS, Yan J, Chhun JJ, Kumirov VK, Foy SG, Masel J, Wysocki VH, Cordes MHJ.  
617 2017. Foldability of a natural *de novo* evolved protein. Structure 25:1687-1696.
- 618 Carvunis A-R, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, Charlotiaux B,  
619 Hidalgo CA, Barbette J, Santhanam B. 2012. Proto-genes and *de novo* gene birth. Nature  
620 487:370-374.

- 621 Chen SCC, Chuang TJ, Li WH. 2011. The Relationships Among MicroRNA Regulation, Intrinsically  
622 Disordered Regions, and Other Indicators of Protein Evolutionary Rate. *Molecular*  
623 *Biology and Evolution* 28:2513-2520.
- 624 Chiarabelli C, Vrijbloed JW, De Lucrezia D, Thomas RM, Stano P, Polticelli F, Ottone T, Papa E,  
625 Luisi PL. 2006. Investigation of de novo totally random biosequences Part II: On the  
626 folding frequency in a totally random library of de novo proteins obtained by phage  
627 display. *Chemistry & Biodiversity* 3:840-859.
- 628 Chiti F, Dobson CM. 2017. Protein misfolding, amyloid formation, and human disease: A  
629 summary of progress over the last decade. *Annual Review of Biochemistry* 86:27-68.
- 630 Choudhari S, Grigoriev A. 2017. Phylogenetic heatmaps highlight composition biases in  
631 sequenced reads. *Microorganisms* 5:4.
- 632 Davidson AR, Lumb KJ, Sauer RT. 1995. Cooperatively folded proteins in random sequence  
633 libraries. *Nature Structural Biology* 2:856-864.
- 634 Davidson AR, Sauer RT. 1994. Folded proteins occur frequently in libraries of random amino-  
635 acid sequences. *Proceedings of the National Academy of Sciences of the United States*  
636 *of America* 91:2146-2150.
- 637 Dosztányi Z, Csizmok V, Tompa P, Simon I. 2005. The pairwise energy content estimated from  
638 amino acid composition discriminates between folded and intrinsically unstructured  
639 proteins. *Journal of Molecular Biology* 347:827-839.
- 640 Dubrey S, Ackermann E, Gillmore J. 2015. The transthyretin amyloidoses: advances in therapy.  
641 *Postgraduate Medical Journal* 91:439-448.

- 642 Fernandez-Escamilla A-M, Rousseau F, Schymkowitz J, Serrano L. 2004. Prediction of sequence-  
643 dependent and mutational effects on the aggregation of peptides and proteins. *Nature*  
644 *biotechnology* 22:1302-1306.
- 645 Foy SG, Wilson BA, Bertram J, Cordes MHJ, Masel J. 2019. A shift in aggregation avoidance  
646 strategy marks a long-term direction to protein evolution. *Genetics* 211:1345-1355.
- 647 Frulloni L, Lunardi C, Simone R, Dolcino M, Scattolini C, Falconi M, Benini L, Vantini I, Corrocher  
648 R, Puccetti A. 2009. Identification of a Novel Antibody Associated with Autoimmune  
649 Pancreatitis. *New England Journal of Medicine* 361:2135-2142.
- 650 Frumkin I, Schirman D, Rotman A, Li F, Zahavi L, Mordret E, Asraf O, Wu S, Levy SF, Pilpel Y.  
651 2017. Gene Architectures that Minimize Cost of Gene Expression. *Molecular Cell* 65:142-  
652 153.
- 653 Goodman DB, Church GM, Kosuri S. 2013. Causes and effects of N-terminal codon bias in  
654 bacterial genes. *Science* 342:475-479.
- 655 Heames B, Schmitz J, Bornberg-Bauer E. 2020. A continuum of evolving de novo genes drives  
656 protein-coding novelty in *Drosophila*. *Journal of Molecular Evolution* 38:382-398.
- 657 Jacobson DR, Pastore RD, Yaghoubian R, Kane I, Gallo G, Buck FS, Buxbaum JN. 1997. Variant-  
658 sequence transthyretin (isoleucine 122) in late-onset cardiac amyloidosis in black  
659 Americans. *New England Journal of Medicine* 336:466-473.
- 660 James JE, Willis SM, Nelson PG, Weibel C, Kosinski LJ, Masel J. 2021. Universal and taxon-  
661 specific trends in protein sequences as a function of age. *eLife* 10:e57347.
- 662 Kaiser CA, Preuss D, Grisafi P, Botstein D. 1987. Many random sequences functionally replace  
663 the secretion signal sequence of yeast invertase. *Science* 235:312-317.

- 664 Keefe AD, Szostak JW. 2001. Functional proteins from a random-sequence library. *Nature*  
665 410:715-718.
- 666 Knopp M, Andersson DI. 2018. No beneficial fitness effects of random peptides. *Nature Ecology*  
667 & *Evolution* 2:1046-1047.
- 668 Knopp M, Gudmundsdottir JS, Nilsson T, Konig F, Warsi O, Rajer F, Adelroth P, Andersson DI.  
669 2019. De novo emergence of peptides that confer antibiotic resistance. *Mbio*  
670 10:e00837-00819.
- 671 Kosinski LJ, Masel J. 2020. Readthrough errors purge cryptic sequences, facilitating the birth of  
672 coding sequence. *Molecular Biology and Evolution* 37:1761–1774.
- 673 Krogh A, Larsson B, von Heijne G, Sonnhammer ELL. 2001. Predicting transmembrane protein  
674 topology with a hidden Markov model: Application to complete genomes. *Journal of*  
675 *Molecular Biology* 305:567-580.
- 676 LaBean TH, Butt TR, Kauffman SA, Schultes EA. 2011. Protein folding absent selection. *Genes*  
677 2:608-626.
- 678 Larsson SC, Markus HS. 2017. Branched-chain amino acids and Alzheimer's disease: a  
679 Mendelian randomization analysis. *Scientific Reports* 7:13604.
- 680 Levy ED, De S, Teichmann SA. 2012. Cellular crowding imposes global constraints on the  
681 chemistry and evolution of proteomes. *Proceedings of the National Academy of*  
682 *Sciences of the United States of America* 109:20461-20466.
- 683 Levy SF, Blundell JR, Venkataram S, Petrov DA, Fisher DS, Sherlock G. 2015. Quantitative  
684 evolutionary dynamics using high-resolution lineage tracking. *Nature* 519:181-186.

- 685 Li F, Salit ML, Levy SF. 2018. Unbiased fitness estimation of pooled barcode or amplicon  
686 sequencing studies. *Cell Systems* 7:521-525.
- 687 Linding R, Schymkowitz J, Rousseau F, Diella F, Serrano L. 2004. A comparative study of the  
688 relationship between protein structure and  $\beta$ -aggregation in globular and intrinsically  
689 disordered proteins. *Journal of Molecular Biology* 342:345-353.
- 690 Liu HX, Zhang RS, Yao XJ, Liu MC, Hu ZD, Fan BT. 2004. Prediction of the isoelectric point of an  
691 amino acid based on GA-PLS and SVMs. *Journal of Chemical Information and Computer  
692 Sciences* 44:161-167.
- 693 Long H, Sung W, Kucukyildirim S, Williams E, Miller SF, Guo W, Patterson C, Gregory C, Strauss  
694 C, Stone C, et al. 2018. Evolutionary determinants of genome-wide nucleotide  
695 composition. *Nature Ecology & Evolution* 2:237-240.
- 696 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for  
697 RNA-seq data with DESeq2. *Genome Biology* 15:550.
- 698 Maurer-Stroh S, Debulpaep M, Kuemmerer N, de la Paz ML, Martins IC, Reumers J, Morris KL,  
699 Copland A, Serpell L, Serrano L, et al. 2010. Exploring the sequence determinants of  
700 amyloid structure using position-specific scoring matrices. *Nature Methods* 7:237-242.
- 701 McLysaght A, Guerzoni D. 2015. New genes from non-coding sequence: the role of de novo  
702 protein-coding genes in eukaryotic evolutionary innovation. *Philosophical Transactions  
703 of the Royal Society B-Biological Sciences* 370:20140332.
- 704 Meszaros B, Erdos G, Dosztanyi Z. 2018. IUPred2A: context-dependent prediction of protein  
705 disorder as a function of redox state and protein binding. *Nucleic Acids Research*  
706 46:W329-W337.

- 707 Moyers BA, Zhang JZ. 2016. Evaluating Phylostratigraphic Evidence for Widespread De Novo  
708 Gene Birth in Genome Evolution. *Molecular Biology and Evolution* 33:1245-1256.
- 709 Moyers BA, Zhang JZ. 2015. Phylostratigraphic Bias Creates Spurious Patterns of Genome  
710 Evolution. *Molecular Biology and Evolution* 32:258-267.
- 711 Neme R, Amador C, Yildirim B, McConnell E, Tautz D. 2017. Random sequences are an abundant  
712 source of bioactive RNAs or peptides. *Nature Ecology & Evolution* 1:0127.
- 713 Prijambada ID, Yomo T, Tanaka F, Kawama T, Yamamoto K, Hasegawa A, Shima Y, Negoro S,  
714 Urabe I. 1996. Solubility of artificial proteins with random sequences. *FEBS Letters*  
715 382:21-25.
- 716 R Core Team. 2019. R: A language and environment for statistical computing: R Foundation for  
717 Statistical Computing.
- 718 Rajon E, Masel J. 2011. Evolution of molecular error rates and the consequences for  
719 evolvability. *Proceedings of the National Academy of Sciences of the United States of*  
720 *America* 108:1082-1087.
- 721 Rousseau F, Schymkowitz J, Serrano L. 2006. Protein aggregation and amyloidosis: confusion of  
722 the kinds? *Current Opinion in Structural Biology* 16:118-126.
- 723 Sormanni P, Aprile FA, Vendruscolo M. 2015. The CamSol Method of Rational Design of Protein  
724 Mutants with Enhanced Solubility. *Journal of Molecular Biology* 427:478-490.
- 725 Theillet F-X, Kalmar L, Tompa P, Han K-H, Selenko P, Dunker AK, Daughdrill GW, Uversky VN.  
726 2013. The alphabet of intrinsic disorder: I. Act like a Pro: On the abundance and roles of  
727 proline residues in intrinsically disordered proteins. *Intrinsically Disordered Proteins*  
728 1:e24360.

- 729 Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. 2013. Maximum allowed solvent  
730 accessibilities of residues in proteins. *PLoS One* 8:e80635.
- 731 Tretyachenko V, Vymetal J, Bednarova L, Kopecky V, Hofbauerova K, Jindrova H, Hubalek M,  
732 Soucek R, Konvalinka J, Vondrasek J, et al. 2017. Random protein sequences can form  
733 defined secondary structures and are well-tolerated in vivo. *Scientific Reports* 7:15449.
- 734 Tsai J, Taylor R, Chothia C, Gerstein M. 1999. The packing density in proteins: Standard radii and  
735 volumes. *Journal of Molecular Biology* 290:253-266.
- 736 Vakirlis N, Acar O, Hsu B, Coelho NC, Van Oss SB, Wacholder A, Medetgul-Ernar K, Bowman RW,  
737 Hines CP, Iannotta J, et al. 2020. De novo emergence of adaptive membrane proteins  
738 from thymine-rich genomic sequences. *Nature Communications* 11:781.
- 739 Vakirlis N, Carvunis A-R, McLysaght A. 2020. Synteny-based analyses indicate that sequence  
740 divergence is not the main source of orphan genes. *eLife* 9:e53500.
- 741 Vakirlis N, Hebert AS, Opuente DA, Achaz G, Hittinger CT, Fischer G, Coon JJ, Lafontaine I. 2018.  
742 A Molecular Portrait of De Novo Genes in Yeasts. *Molecular Biology and Evolution*  
743 35:631–645.
- 744 Van Oss SB, Carvunis AR. 2019. *De novo* gene birth. *PLoS Genetics* 15:e1008160.
- 745 Vecchi G, Sormanni P, Mannini B, Vandelli A, Tartaglia GG, Dobson CM, Hartl FU, Vendruscolo  
746 M. 2020. Proteome-wide observation of the phenomenon of life on the edge of  
747 solubility. *Proceedings of the National Academy of Sciences of the United States of*  
748 *America* 117:1015-1020.
- 749 Weisman CM, Eddy SR. 2017. Gene Evolution: Getting Something from Nothing. *Current Biology*  
750 27:R661-R663.



- 751 Wickham H. 2016. ggplot2: Elegant graphics for data analysis. New York: Springer-Verlag.
- 752 Willis S, Masel J. 2018. Gene birth contributes to structural disorder encoded by overlapping  
753 genes. *Genetics* 210:303-313.
- 754 Wilson BA, Foy SG, Neme R, Masel J. 2017. Young genes are highly disordered as predicted by  
755 the preadaptation hypothesis of *de novo* gene birth. *Nature Ecology & Evolution* 1:0146.
- 756 Wilson BA, Masel J. 2011. Putatively noncoding transcripts show extensive association with  
757 ribosomes. *Genome Biology and Evolution* 3:1245-1252.
- 758