

1 **Title:** Identifying drivers of parallel evolution: A regression model approach

2

3 **Authors:** Susan F. Bailey^{1,2,*}, Qianyun Guo¹, and Thomas Bataillon¹

4

5 **Author affiliations:**

6 ¹ Bioinformatics Research Centre, Aarhus University, C.F. Møllers Allé 8, DK-8000 Aarhus C,

7 Denmark.

8 ² Current affiliation: Department of Biology, Clarkson University, PO Box 5805, Potsdam, NY 13699-

9 5805

10

11 *Author for Correspondence: Susan F. Bailey, Department of Biology, Clarkson University, PO Box

12 5805, Potsdam, NY 13699-5805, phone: 315-268-6400 ext.4263, email: sbailey@clarkson.edu

13

14 **Running head:** Identifying drivers of parallel evolution

15

16 **Keywords:** parallel evolution, experimental evolution, Poisson regression, negative binomial

17 regression

18

19 **Data archival location:** Dryad, doi to be included later

20 **Abstract**

21 Parallel evolution, defined as identical changes arising in independent populations, is often attributed to
22 similar selective pressures favoring the fixation of identical genetic changes. However, some level of
23 parallel evolution is also expected if mutation rates are heterogeneous across regions of the genome.
24 Theory suggests that mutation and selection can have equal impacts on patterns of parallel evolution,
25 however empirical studies have yet to jointly quantify the relative importance of these two processes.
26 Here, we introduce several statistical models to examine the contributions of mutation and selection
27 heterogeneity to shaping parallel evolutionary changes at the gene-level. Using this framework we
28 analyze published data from forty experimentally evolved *Saccharomyces cerevisiae* populations. We
29 can partition the effects of a number of genomic variables into those affecting patterns of parallel
30 evolution via effects on the rate of arising mutations, and those affecting the retention versus loss of the
31 arising mutations (i.e. selection). Our results suggest that gene-to-gene heterogeneity in both mutation
32 and selection, associated with gene length, recombination rate, and number of protein domains drive
33 parallel evolution at both synonymous and nonsynonymous sites.

34 **Introduction**

35 Documenting patterns of parallel evolution during the adaptive divergence of populations or during
36 repeated bouts of adaptation in populations maintained in the lab is becoming increasingly feasible.
37 Beyond the fascination for the pattern of repeatable evolution, an outstanding open question is to
38 understand which underlying processes are driving the pattern of molecular evolution during
39 adaptation. Theory makes clear cut predictions: in the absence of selective interference between
40 beneficial mutations (the so called strong selection weak mutation, or SSWM, domain), heterogeneity
41 in mutation rates and selection coefficients between loci are expected to have equal influence on
42 patterns of parallel evolution (Chevin et al., 2010; Lenormand et al., 2016). So far very few empirical
43 studies have attempted to jointly quantify the relative importance of these two processes in shaping
44 patterns of parallel evolution in genetic data. One study has explored this indirectly by quantifying the
45 contribution of these two processes in shaping the parallel evolution of heritable traits that are assumed
46 to be associated with parallel genetic changes (Streisfeld and Rausher, 2011). Recent work by Bailey et
47 al., 2017 outlines an approach for quantifying the effects of mutation and selection heterogeneity in
48 driving parallel evolution in experimental evolution data, but this alternate approach can not identify
49 potential genomic drivers of that heterogeneity, as we do here. Other previous studies looking
50 explicitly at parallel genetic changes have focused on the impacts of either selection or mutation
51 separately.

52 A number of studies have examined gene-level mutation counts, looking for levels of parallel
53 evolution that exceed what one would expect in the absence of selection (Caballero et al., 2015; Marvig
54 et al., 2015; Woods et al., 2006), according to some null model, with an aim to identify genes that are
55 under selection. For example, (Caballero et al., 2015) calculated the probability of instances of gene-
56 level parallel evolution in whole genome sequences of *Pseudomonas aeruginosa* repeatedly sampled
57 over the course of a year from the sputum of a cystic fibrosis (CF) patient assuming uniform re-

58 sampling of ~150 mutation events across the approximately 6000 genes in the genome. The authors
59 were able to identify 19 different genes for which there was significant deviation from their null model,
60 and that pattern was interpreted as evidence for selection acting on these genes. However this study and
61 other similar approaches do not account for the possibility of heterogeneity in mutation rate from gene-
62 to-gene, a process that can generate false positives when using “abnormal” levels of parallel evolution
63 as a means to detect selected genes.

64 Others have compared instances of parallel and convergent evolution across species (see
65 Christin et al., 2010 for a review and examples). These studies also aim to identify genes under
66 selection by searching for genes that exhibit a higher than expected number of instances of parallel
67 evolution according to a specified null model for evolution. Many cross-species comparative studies
68 report instances of parallel molecular evolution and readily interpret these as being driven by positive
69 selection (e.g. Castoe et al., 2009; Feldman et al., 2012; Jost et al., 2008; Liu et al., 2014). However
70 (Zou and Zhang, 2015) show that in this type of analysis the choice of null model is crucial and suggest
71 that many previously reported instances of parallel evolution driven by selection could in fact have
72 resulted simply from mutation biases and mutational heterogeneity in the absence of selection.

73 In contrast to studies aimed at identifying selection, other work has focused on examining how
74 heterogeneity in mutation rate can effect the distribution of mutations across a genome, and so the
75 probability of parallel evolution. These studies focus exclusively on either those mutations that are
76 assumed to be to a first approximation neutral (e.g. synonymous mutations, Maddamsetti et al., 2015)
77 or mutations arising in the course of experiments where selection is minimal (e.g. mutations arising in a
78 mutation accumulation experiment, Ness et al., 2015). On the whole, these studies suggest substantial
79 gene-to-gene heterogeneity in mutation rate and this can arguably also generate differences in the
80 distribution of mutations across the genome (although studies differ in the factors identified that drive
81 that heterogeneity). However, it is not clear what the relative contribution of mutation rate

82 heterogeneity is when the mutations of interest also have the potential to be under varying degrees of
83 selection.

84 In this study we aim to explore the effects of both mutation and selection in generating observed
85 patterns of parallel evolution at the gene-level. To do this we propose a framework that explicitly
86 considers both selection and mutational heterogeneity. Using both Poisson and negative binomial
87 regression models, we analyze gene-level mutation count data obtained from whole genome
88 sequencing of a large set of yeast (*Saccharomyces cerevisiae*) experimental populations that were
89 adapted in parallel to a glucose environment in the lab (Lang et al., 2013). We find that the best
90 predictor of parallel mutations at the gene-level is simply the length of the gene, and along with this, a
91 few other genomic covariates – namely the number of protein domains and the rate of recombination –
92 also affect patterns of parallel evolution.

93

94 **Models for identifying processes underlying parallel evolution**

95 We are interested in quantifying heterogeneity in mutation rate and selection, and how these in turn are
96 driving patterns of parallel evolution, and identifying genomic variables that predict how these
97 processes vary from gene-to-gene. To accomplish this, we need a framework that can explicitly
98 separate the effects of variation in mutation rate and variation in selection. We do this by examining
99 separately the observed synonymous and nonsynonymous mutations, making the assumption (which
100 we then check) that gene-to-gene variation in the rate at which synonymous mutations rise to
101 observable frequencies is driven solely by variation in the mutation rate per gene, while gene-to-gene
102 variation in the rate at which nonsynonymous mutations have arisen may be driven by heterogeneity in
103 both mutation and selection processes. We describe the number of mutations observed in gene i during
104 the course of an experiment, as $X_i = X_i^S + X_i^N$, where X_i^S and X_i^N denote the synonymous and
105 nonsynonymous mutation counts respectively. We assume these mutations are Poisson distributed with

106 rates λ_i^S and λ_i^N respectively. For synonymous mutations, this Poisson rate can be modeled as

107
$$\lambda_i^S = M_0 \mu_0 L_i \pi_0$$
 Eqn (1)

108 Here, M_0 is a parameter that absorbs both time and population size at which the evolution occurred and
109 that is constant across the genome, μ_0 is the per-nucleotide mutation rate that we assume (and check) is
110 constant across the genome, L_i is the length of gene i in nucleotides, and π_0 is the probability of a
111 synonymous mutation rising to an observable frequency in the population (we assume that synonymous
112 mutations are selectively neutral and so this probability is assumed to be constant across the genome).
113 For nonsynonymous mutations,

114
$$\lambda_i^N = \lambda_i^S \pi_i,$$
 Eqn (2)

115 where π_i is the probability of a nonsynonymous mutation in gene i rising to observable frequencies in
116 the population. This probability, π_i , is a function of the mean selection coefficient of gene i , s_i , and
117 under strong-selection-weak-mutation (SSWM) conditions, $\pi_i \propto s_i$ (Gillespie, 1984). The type of data
118 used and the underlying assumptions are summarized in Fig. 1.

119 Given these underlying assumptions about the processes giving rise to observable mutations in
120 the experimental sequence data, we can then use Poisson and negative binomial (NB) regression to
121 identify potential genomic variables that significantly explain variation in λ_i^N and λ_i^S , and thus
122 ultimately in the mutation and selection processes from gene-to-gene. The Poisson regression is used to
123 explore counts of rare events (i.e. the observed mutations) that have a fixed probability of being
124 observed, while for a NB regression, the rate of those rare events is itself a random variable that is
125 gamma-distributed. A NB regression incorporates an extra parameter beyond a Poisson rate, known as
126 the dispersion parameter (here denoted by θ), reflecting the amount of underlying variation in the rate
127 of observed mutations from gene-to-gene and governs the “extra” variance of the NB distribution
128 relative to a Poisson distribution with identical mean. If there is no heterogeneity among the rate of
129 observed mutations from gene-to-gene, the dispersion parameter θ goes to zero and we recover a

130 Poisson regression model. Therefore, the Poisson regression model is a special case of the NB
131 regression model, as $NB(\lambda_i, \theta)$ reduces to $Poisson(\lambda_i)$ at the limit of $\theta \rightarrow 0$ (see for instance Zuur,
132 2009). As a consequence, the Poisson and NB models are “nested” and their relative fit can be
133 compared using a likelihood ratio test when exploring the fit of both types of regression models in this
134 study.

135 More precisely, we use the models $X_i \sim Poisson(\lambda_i)$ or $X_i \sim NB(\lambda_i, \theta)$, fitting the following
136 regression:

$$137 \quad \log(\lambda) = constant + \alpha_1 \mathbf{A}_1 + \alpha_2 \mathbf{A}_2 + \dots + \alpha_j \mathbf{A}_j, \quad \text{Eqn (3)}$$

138 where $\lambda = (\lambda_1, \dots, \lambda_i, \dots, \lambda_n)$ are the Poisson rates for all n genes, $\mathbf{A}_1 \dots \mathbf{A}_j$ are the j potential genomic
139 explanatory variables, and $\alpha_1 \dots \alpha_j$, the estimated regression coefficients for those j variables. Thus, in
140 the case of the synonymous mutations, $constant = \log(M_0 \pi_0 \mu_0)$, $\mathbf{A}_1 = \log(L_i)$ setting $\alpha_1 = 1$. For
141 nonsynonymous mutations, $\alpha_2 \mathbf{A}_2 + \dots + \alpha_j \mathbf{A}_j = \log(\pi_i)$. Details of the implementation of these models
142 is provided below.

143 **Methods**

144 ***The data***

145 *Evolution experiment data* – We analyzed data obtained from whole genome re-sequencing of forty
146 populations of *S. cerevisiae* adapted in parallel to a glucose-limited environment in the lab (Lang et al.,
147 2013). In our analysis we include all detected genic mutations, i.e. all genic mutations that were able to
148 escape drift and so rise to frequencies of at least approximately 10% in the populations (mutations
149 below this frequency could not reliably be detected, see Lang et al., 2013). Mutations were grouped by
150 gene across all forty populations, and categorized as synonymous (SYN) or nonsynonymous (NS), i.e.
151 those that do not confer amino acid changes, and those that do, respectively.

152

153 *Comparative genomics data* – We used a set of orthologous gene alignments spanning four distinct

154 yeast species (*S. cerevisiae*, *S. paradoxus*, *S. bayanus*, and *S. mikatae*; available from
155 www.yeastgenome.org/download-data/genomics; Kellis et al., 2003; Cliften et al., 2003) to infer the
156 gene-to-gene heterogeneity of the substitution rates at synonymous sites and nonsynonymous sites,
157 hereafter dS and dN respectively. To do so, we first realigned the gene sequences using ClustalW
158 (Larkin et al., 2007) on the translated protein sequence data and then applied a number of filters to the
159 data with an aim at removing those gene alignments that might result in inaccurate codon substitution
160 model predictions. We removed alignments for those genes where sequences were not available from
161 all four species, alignments for which at least one sequence had <30% overlap with the one of the other
162 3 sequences, and alignments for which at least one sequence was <300 bps in length. We then used a
163 maximum likelihood codon based method (CodeML in the PAML software package; Yang, 2007) to
164 infer dS and dN , for each gene in our data set. We used a codon table model with a fixed tree topology
165 (a comparison of AICs among alternative codon based models indicated this was the most appropriate
166 model for the data set).

167

168 *Additional genomic data* – We included eight additional genomic variables in our analysis that we
169 expected could have the potential to effect the probability of a gene to harbor mutations. Our collection
170 of variables is not meant to be exhaustive, but simply meant to illustrate the potential for additional
171 genomic information to improve our predictions of which genes bear mutations across the genome. For
172 each gene we consider: gene length, % GC content, multi-functionality, degree of protein-protein
173 interaction (PPI), codon adaptation index (CAI), number of domains, level of expression, local
174 recombination rate, and essential genes. We expect some of these variables may capture heterogeneity
175 in the per-gene mutation rates, for example: gene length, which likely captures variation in a gene's
176 mutational target size, and local recombination rate, which has been shown to be associated with
177 mutability in yeast (Holbeck and Strathern, 1997; Strathern et al., 1995). We expect other variables may

178 capture heterogeneity in selection from gene-to-gene, for example: multi-functionality and PPI, which
179 may characterize aspects of how pleiotropic a given gene is and so the level of evolutionary constraint
180 it is under. We expect still other genomic variables may capture heterogeneity in both mutation and
181 selection. For example, level of expression of a gene may be correlated with gene-to-gene variation in
182 selection as highly expressed genes have been shown to be more highly conserved, both specifically in
183 yeast (Drummond et al., 2005; Pál et al., 2001) and as a more general phenomenon across species
184 (Drummond and Wilke, 2008). On the other hand, level of expression of a gene has also been shown to
185 be positively correlated with mutability (Ness et al., 2015). Descriptions of the variables used in this
186 study and sources from which the data were obtained are provided in Table 1.

187 A data set integrating the mutation counts originally made available by Lang et al., 2013 (from
188 their Supplementary Table 1) and all the genomic covariates that we aggregated for this study, as well
189 as the gene alignments used for estimating dN and dS are available on Dryad (doi will be inserted here).

190

191 ***Regression models***

192 *Regression models and explanatory variables tested* – We used the Poisson and negative binomial
193 regression models described in the “Models” section above to examine how much of the variation in
194 our explanatory variables could account for patterns of variation in mutation counts per gene. We used
195 the 'glm' and 'glm.nb' functions in R (R Development Core Team, 2014) to implement these models.
196 We fit a series of models to synonymous and nonsynonymous mutation count data separately. To start,
197 we fit the synonymous mutations (model M_S), testing our assumptions that rate of observed mutations
198 per gene is proportional to number of nucleotide sites in the gene (L_i), and the per nucleotide mutations
199 does not vary significantly across the genome – i.e. a model assuming μ_0 is a fixed parameter (Poisson
200 regression) fits the data better than a model where μ_0 for each gene is drawn from a gamma distribution
201 (NB regression).

202 After these assumptions were confirmed, we moved on to fit the nonsynonymous mutation data
203 (M_N), testing the 11 genomic variables listed in Table 1. We then examined an alternate model (M_N^{PC}),
204 fitting the nonsynonymous mutations using the principal components of the 11 genomic variables in
205 place of the raw variables. The reason we explore this model is that many genomic variables tend to be
206 correlated (for correlations between the particular variables used in this study, see supplementary Table
207 S1), and one approach to reducing potential problems with co-linearity is to transform the raw variables
208 into their principal components and use the resulting uncorrelated composite variables for the
209 regression analysis. We performed a principal component analysis on 11 genomic variables using the
210 'prcomp' function in R to obtain 11 principal components (PCs).

211

212 *Model selection and significance of variables* – For each variable and parameter of interest we tested
213 significance by comparing versions of the models with and without that variable or parameter of
214 interest through a likelihood-ratio test (LRT). Significance testing for LRTs was done using
215 permutation tests instead of relying on asymptotic distribution of the LRTs, approximating the null
216 distribution and obtaining P-values by calculating the frequency of permutations where the model fit
217 resulted in a likelihood-ratio greater than or equal to the observed value. Variables found to
218 significantly improve model fit were retained in the final “best” model. We choose to test significance
219 using permutations given that asymptotic results on the distribution of the likelihood ratio test may
220 break down as the reduced model – the Poisson regression – lies at the boundaries of the parameter
221 space for θ , included in the NB regression (see for instance Self and Liang, 1987). In practice, 1000
222 permutations were used to approximate the null and obtain p-values on each variable (more
223 permutations might be required if needed to approximate p-values that are much smaller than 10^{-3}).

224 The two nonsynonymous mutation models M_N and M_N^{PC} were compared with each other using
225 Akaike information criterion (AIC; Akaike, 1973), and the proportion of variation explained (pseudo-

226 R^2) was estimated as the R^2 obtained from a linear regression (using 'lm' in R) between the observed
227 and predicted mutation counts for a given model. Note that this statistic is not used for any formal
228 goodness-of-fit but as an illustrative way to report how much of the whole variation is accounted for by
229 any model we fit to the mutation count data.

230 All statistical analyses, including the permutation tests, were scripted in R (R Development
231 Core Team, 2014) and an example script for implementing our model framework and hypothesis testing
232 is available on Dryad (doi will be inserted here).

233

234 **Results**

235 ***The data***

236 *Mutation counts data* – We used experimental data comprising all mutations detected at a frequency
237 over 10% in the forty evolved *S. cerevisiae* populations described in Lang et al., 2013. After removing
238 those genes for which we had incomplete or unreliable data (see Methods), we were left with 2891
239 genes out of a total of 6603. The filtered data set contained 357 nonsynonymous mutations distributed
240 across 267 genes, and 58 synonymous mutations distributed across 57 genes. The genes removed by
241 our filtering rules had disproportionately more mutations compared to those genes that were retained in
242 the data set ($\chi^2 = 50.57$, $df = 1$, $P < 0.001$). This is not unexpected as highly divergent genes are more
243 likely to be filtered out due to alignment issues, and it is not surprising that highly divergent genes
244 would tend see more mutations than average, whether it be as a result of mutation and / or selection
245 mechanisms. This bias in the filtering means that our results are likely conservative in terms of
246 detecting significant relationships between long-term (from comparative genomics data) and short-term
247 (from experimental evolution data) measures of divergence.

248

249 *Genomic variables* – We used codon substitution models comparing four yeast species to estimate dS

250 and dN/dS for each gene. Estimates for dS ranged widely, from 0.21 to 68.7, however the vast majority
251 of dS estimates (~95%) were less than 4. Estimates for dN/dS ranged from 0.00010 to 0.43, and these
252 values are weakly negatively correlated with dS ($r = -0.043$, $P = 0.021$). We collated and/ or calculated
253 nine other genomic variables with the potential to effect the mutation and selection processes in this
254 system and estimated correlation coefficients between all pairs of explanatory variables used in this
255 study (Table S1). While the correlations between these variables tend to be quite weak, many are, in
256 fact, significant due to the large number of observations in the data set.

257

258 ***Mutation counts analysis***

259 *Synonymous mutations* – We used regression models to test our assumption that gene-level mutation
260 rate can be adequately described as simply being directly proportional to gene length. Restricting the
261 data to the synonymous mutations, we compared Poisson regression models with and without gene
262 length included as an explanatory variable (M_{S0} : $\lambda_s = constant$ and M_{S1} : $\lambda_s = constant * (L_i)^\alpha$,
263 respectively), and a Poisson regression model where rate is restricted to be directly proportional to gene
264 length (i.e. M_{S2} : $\lambda_s = constant * L_i$). We also compared with negative binomial versions of these model
265 to look for the possibility of additional unexplained variation in the rate λ . The results of these
266 comparisons are shown in Table 2. Model M_{S2} was the best model according to a comparison of AICs,
267 confirming our assumptions. The fits of these models to the distribution of synonymous mutation
268 counts per gene are visualized in Fig. 2A.

269

270 *Nonsynonymous mutations* – We fit regression models to the nonsynonymous mutation data, including
271 eleven genomic variables, trying to identify which of those variables could significantly explain
272 variation in the number of observed mutations per gene. We found that gene length (L), number of
273 domains in the encoded protein ($num.dom$), and recombination rate (r) were significant in our model

274 (see model $M_N.NB$ in Table 3).

275 When we fit regression models using the principal components of the genomic variables in
276 place of the raw variables, we found that only a single principal component, PC10, was significant in
277 the model (see model $M_N.NB_{PC}$ in Table 3). PC10 is fairly evenly loaded with a number genomic
278 variables (see Fig. 3), however the three significant genomic variables from $M_N.NB$ (L , *num.dom*, and
279 r) are among the variables more heavily loaded on PC10, so the two models seem to be roughly in
280 agreement. A comparison of Poisson and negative binomial regression models, as well as models
281 including the raw genomic variables versus the transformed principal component variables, suggests
282 that the best model for these nonsynonymous mutation count data is a negative binomial regression
283 using the raw genomic variables (see AIC values in Table 4). The fits of these models to the distribution
284 of nonsynonymous mutation counts per gene are visualized in Fig. 2B.

285

286 **Discussion**

287 Here we present a modeling framework to infer what genomic variables may underlie gene to gene
288 variation in mutation rate and intensity of selection. We use these models to provide evidence that
289 parallel evolution at both nonsynonymous and synonymous sites is driven by non trivial amounts of
290 gene-to-gene heterogeneity in the mutation and selection processes. Using our modeling approach, we
291 identified a number of genomic variables that can significantly predict the distribution of mutations
292 observed across genes in experimentally evolved populations of *S. cerevisiae* (Lang et al., 2013). We
293 are also able to classify genomic variables into those that have affected mutation counts 1) through
294 their effect on the mutation rate (variables that significantly predict synonymous mutations), and/ or 2)
295 through their effect on the probability of a mutation being either observed/ lost due to selection
296 (variables that significantly predict nonsynonymous mutations). Out of all the variables tested, we
297 found that *gene length* explained the most variation in both synonymous and nonsynonymous mutation

298 counts per gene – plainly speaking, longer genes accumulate more mutations. However, *number of*
299 *domains* and *recombination* also had significant effects. Below we discuss in detail these genomic
300 variables and their potential contributions to the probability of parallel evolution via the processes of
301 mutation and selection.

302

303 *Longer genes harbor more mutations* – By far, the variable having the largest effect on variation in the
304 number of synonymous and nonsynonymous mutations observed was *gene length*. More specifically,
305 gene length positively affected the rate of mutation at the gene-level, meaning genes comprising more
306 nucleotides were more likely to harbor mutations. This result is not surprising and is in agreement with
307 recent analysis of synonymous mutation counts from Lenksi's long term evolution experiment with *E.*
308 *coli* (Maddamsetti et al., 2015).

309

310 *Long-term divergence does not predict short-term mutation counts* – Our model for synonymous
311 mutation counts suggests that divergence estimates from long-term evolutionary comparisons at the
312 species level do not provide insight into expected mutation counts on the shorter time scale of evolution
313 in the lab, also in agreement with recent analysis of *E. coli* data (Maddamsetti et al., 2015).
314 Maddamsetti et al found that their proxy for long-term per gene mutation rate, θ_s (a measure of within-
315 species nucleotide diversity), did not explain gene-to-gene variation in synonymous mutation counts in
316 their data. The authors argued that horizontal gene transfer (HGT) is therefore likely a more important
317 process driving gene-to-gene variation in long-term divergence between naturally occurring *E. coli*
318 strains, and since HGT did not occur in their evolution experiment, it is not surprising that the
319 experiment's synonymous mutation counts did not correlate with θ_s . However, rates of HGT tend to be
320 higher in bacteria, and in particular *E. coli*, as compared to yeast and other eukaryotes (e.g. Boto 2010).
321 Furthermore, a recent mutation accumulation experiment with the eukaryote *Chlamydomonas*

322 *reinhardtii* showed a positive correlation between a proxy for long-term mutation rate (θ_s) and per site
323 mutability (Ness et al., 2015). Thus, it is somewhat surprising that we do not see a significant
324 relationship between dS and dN/dS and counts of synonymous and nonsynonymous mutations
325 respectively in our examination of the *S. cerevisiae* data used in this study. One possibility might also
326 be that dS and dN/dS are noisy to estimate at the gene level and that tends to downplay their predictive
327 power in our analysis of counts in evolve and re-sequence experiment.

328

329 *Nonsynonymous mutation counts show evidence of selection heterogeneity* – As expected (Lenormand
330 et al., 2016), we report strong evidence that the distribution of nonsynonymous mutations across the
331 genome was driven in part by gene-to-gene heterogeneity in selection. Of those genomic variables
332 tested, we found three that were significant predictors of nonsynonymous mutation counts, suggesting
333 that those variables may drive or are correlated with processes that modulate the intensity of selection
334 across genes. The significant variables were *gene length*, *recombination rate*, and *number of protein*
335 *domains*.

336 We found that *gene length* predicts nonsynonymous mutation count via selection, over and
337 above its effects on per gene mutation rate – as estimated from models aimed at explaining the
338 synonymous mutation count only. While one might not expect *gene length* to have direct effects on
339 selection, we suggest that *gene length* may show a significant effect here because it is correlated with
340 other attributes of the genome that could have important effects on selection, for example *gene*
341 *expression levels* and *multifunctionality*. Because of these correlations, it could be that *gene length* acts
342 as a kind of summary variable for these covariates and other unidentified factors we have not captured
343 in these models. Further evidence that *gene length* acts as a summary variable comes from the M3
344 results (summarized in Table 3), where we see that *gene length* is no longer significant when other
345 summary variables – the principal components – are included in the model.

346 In contrast to the positive relationship between *gene length* and number of nonsynonymous
347 mutations, we also found that the *number of protein domains* that a gene codes for (a variable that is
348 positively correlated with gene length; Table S1) actually negatively predicts the number of
349 nonsynonymous mutations. In other words, the more domains in the encoded protein of a gene, the
350 fewer mutations that gene is expected to incur in the course of the yeast evolution experiment analyzed
351 here. The mechanism behind this effect is not clear, but certainly protein structure has previously been
352 reported to have significant impacts on evolutionary rates in yeast (Bloom et al., 2006) and one can
353 also posit that genes encoding proteins with multiple domains and thereby involved in more numerous
354 interactions are – all else being equal – more severely constrained by purifying selection. It is
355 interesting that this effect can be observed in the course of relatively short time span (relative to
356 between species divergence times) through the relative paucity of nonsynonymous mutations in these
357 genes.

358 Our analysis also showed that *recombination rate* is a significant predictor of the observed
359 number of nonsynonymous mutations observed in a given gene in these data. Genes with higher
360 recombination rates are more likely to bear nonsynonymous mutations. We expect recombination rate
361 to be correlated with mutation, as previous studies in yeast have shown that recombinational repair of
362 double strand breaks in substantially increases the frequency of nearby point mutations in nearby
363 intervals (e.g. Holbeck and Strathern, 1997; Strathern et al., 1995). However, it is not clear how high
364 recombination rates might drive, or be correlated with other processes that drive, selection – as our
365 models suggest is the case for this data set. Another non exclusive possibility might be the fact that
366 biased gene conversion might vary from gene to gene and also – like selection - affect the probability
367 of detecting variants in evolve and re-sequence experiments

368

369 *Factors driving mutation and selection are complex* – It is difficult to obtain any additional insights

370 from models that include principal components of the genomic covariate data, however there is at least
371 some level of agreement between those variables that are significant (i.e. length, recombination, and
372 number of domains) and ones that are heavily weighted in PC10 – the principal component that was
373 found to be significant (see Fig. 3). The local properties of the genome do appear to drive some
374 heterogeneity in the selection processes, and in turn, shape the patterns of parallel evolution, however
375 individual effects that can be ascribed to individual variables are not easy to parse out.

376 Finally we want to stress that while we were able to identify a number of factors affecting the
377 count of mutations observed in this evolution experiment data set, the total explained variance is still
378 low: 1 % and 16.0 % in the synonymous and nonsynonymous models respectively (calculated from
379 pseudo- r^2 estimates of the “best” models, see methods). While the models do capture the general
380 distribution of mutation counts (Fig. 2) and so the degree of parallel evolution, accurately predicting on
381 which genes those mutations will fall is still very difficult. This is not surprising given the amount of
382 stochasticity involved in both the origin of new mutations and their evolutionary fate through drift and
383 selection. A clearer picture might emerge when using our modeling approach in a meta-analysis
384 approach where several evolve and re-sequence experiments are considered together (see Bailey et al.,
385 2017 for a similar approach on summary statistics of the amount of parallel evolution at the gene level
386 across a wide range of experimental studies in yeast and bacteria)

387 While we do find a number of genomic variables that significantly affect the distribution of
388 mutations across the genome, it is noteworthy that these models are still unable to capture the more
389 extreme patterns of parallel evolution observed in this data set. For example, one gene (IRA1) saw
390 mutations in over 50% of the populations sequenced in this experimental data set (discussed in more
391 detail in Lang et al., 2013). Such a mutation count is completely out of the range of likely outcomes
392 predicted by our models. Some of this discrepancy may be because of the simplifying assumptions
393 made about the process of selection. Our framework models the process of mutation and its

394 heterogeneity but while we account for the fact that newly arising mutations may have different
395 probabilities of reaching an observable frequency, the modeling of that process could be made more
396 precise by incorporating an explicit underlying distribution of fitness effects of new mutations at each
397 gene. Incorporating a selection process that allows for different amounts of both positive and negative
398 selection, as well as further details about the selection pressures in the particular environment of
399 interest – something we do not consider at all in this study – would likely improve prediction for some
400 of these more extreme events.

401

402 *Advantages of this regression framework* – Relying on the assumption that synonymous mutations are
403 selectively neutral (which does appear to be the case for these data), the regression models we use in
404 this study allow us to distinguish between genomic variables influencing the observed distribution of
405 mutations across a genome through their potential effects on both gene-to-gene heterogeneity in
406 mutation rate and gene-to-gene heterogeneity in selection. The great advantage of this is that it allows
407 us to begin to break down the importance of these two processes in shaping patterns of parallel
408 evolution we see, and move closer the goal of predicting which genes will be involved in evolution
409 when organisms adapt to new environments. It will be interesting to apply this model framework to
410 other data sets of this type, as they become available, to see how general these patterns are across
411 different organisms and selection environments (Bailey and Bataillon, 2016).

412

413 **Acknowledgments**

414 This work was supported by the European Research Council under the European Union's Seventh
415 Framework Program [FP7/20072013, ERC grant number 311341 to T.B.].

416 References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Proceeding of the Second International Symposium on Information Theory*, (Budapest: Akademiai Kiado), pp. 267–281.
- Bailey, S.F., and Bataillon, T. (2016). Can the experimental evolution programme help us elucidate the genetic basis of adaptation in nature? *Mol. Ecol.* 25, 203–218.
- Bailey, S.F., Blanquart, F., Bataillon, T., and Kassen, R. (2017). What drives parallel evolution? *BioEssays* 39, 1–9.
- Bloom, J.D., Drummond, D.A., Arnold, F.H., and Wilke, C.O. (2006). Structural Determinants of the Rate of Protein Evolution in Yeast. *Mol. Biol. Evol.* 23, 1751–1761.
- Caballero, J.D., Clark, S.T., Coburn, B., Zhang, Y., Wang, P.W., Donaldson, S.L., Tullis, D.E., Yau, Y.C.W., Waters, V.J., Hwang, D.M., et al. (2015). Selective sweeps and parallel pathoadaptation drive *Pseudomonas aeruginosa* evolution in the cystic fibrosis lung. *mBio* 6, e00981-15.
- Castoe, T.A., de Koning, A.J., Kim, H.-M., Gu, W., Noonan, B.P., Naylor, G., Jiang, Z.J., Parkinson, C.L., and Pollock, D.D. (2009). Evidence for an ancient adaptive episode of convergent molecular evolution. *Proc. Natl. Acad. Sci.* 106, 8986–8991.
- Cherry, J.M., Hong, E.L., Amundsen, C., Balakrishnan, R., Binkley, G., Chan, E.T., Christie, K.R., Costanzo, M.C., Dwight, S.S., Engel, S.R., et al. (2012). *Saccharomyces* Genome Database: the genomics resource of budding yeast. *Nucleic Acids Res.* 40, D700–D705.
- Chevin, L.-M., Martin, G., and Lenormand, T. (2010). Fisher’s model and the genomics of adaptation: restricted pleiotropy, heterogenous mutation, and parallel evolution. *Evolution* 64, 3213–3231.
- Christin, P.-A., Weinreich, D.M., and Besnard, G. (2010). Causes and evolutionary significance of genetic convergence. *Trends Genet.* 26, 400–405.
- Cliften, P., Sudarsanam, P., Desikan, A., Fulton, L., Fulton, B., Majors, J., Waterston, R., Cohen, B.A., and Johnston, M. (2003). Finding functional features in *Saccharomyces* genomes by phylogenetic footprinting. *Science* 301, 71–76.
- Drummond, D.A., and Wilke, C.O. (2008). Mistranslation-induced protein misfolding as a dominant constraint on coding-sequence evolution. *Cell* 134, 341–352.
- Drummond, D.A., Bloom, J.D., Adami, C., Wilke, C.O., and Arnold, F.H. (2005). Why highly expressed proteins evolve slowly. *Proc. Natl. Acad. Sci. U. S. A.* 102, 14338–14343.
- Feldman, C.R., Brodie, E.D., Brodie, E.D., and Pfrender, M.E. (2012). Constraint shapes convergence in tetrodotoxin-resistant sodium channels of snakes. *Proc. Natl. Acad. Sci.* 109, 4556–4561.
- Gillespie, J.H. (1984). Molecular evolution over the mutational landscape. *Evolution* 38, 1116–1129.

- Holbeck, S.L., and Strathern, J.N. (1997). A role for REV3 in mutagenesis during double-strand break repair in *Saccharomyces cerevisiae*. *Genetics* 147, 1017–1024.
- Holstege, F.C., Jennings, E.G., Wyrick, J.J., Lee, T.I., Hengartner, C.J., Green, M.R., Golub, T.R., Lander, E.S., and Young, R.A. (1998). Dissecting the regulatory circuitry of a eukaryotic genome. *Cell* 95, 717–728.
- Illingworth, C.J.R., Parts, L., Bergström, A., Liti, G., and Mustonen, V. (2013). Inferring Genome-Wide Recombination Landscapes from Advanced Intercross Lines: Application to Yeast Crosses. *PLoS ONE* 8, e62266.
- Jost, M.C., Hillis, D.M., Lu, Y., Kyle, J.W., Fozzard, H.A., and Zakon, H.H. (2008). Toxin-resistant sodium channels: Parallel adaptive evolution across a complete gene family. *Mol. Biol. Evol.* 25, 1016–1024.
- Kellis, M., Patterson, N., Endrizzi, M., Birren, B., and Lander, E.S. (2003). Sequencing and comparison of yeast species to identify genes and regulatory elements. *Nature* 423, 241–254.
- Koch, E.N., Costanzo, M., Bellay, J., Deshpande, R., Chatfield-Reed, K., Chua, G., D’Urso, G., Andrews, B.J., Boone, C., Myers, C.L., et al. (2012). Conserved rules govern genetic interaction degree across species. *Genome Biol* 13, R57.
- Lang, G.I., Rice, D.P., Hickman, M.J., Sodergren, E., Weinstock, G.M., Botstein, D., and Desai, M.M. (2013). Pervasive genetic hitchhiking and clonal interference in forty evolving yeast populations. *Nature* 500, 571–574.
- Larkin, M.A., Blackshields, G., Brown, N.P., Chenna, R., McGettigan, P.A., McWilliam, H., Valentin, F., Wallace, I.M., Wilm, A., Lopez, R., et al. (2007). Clustal W and Clustal X version 2.0. *Bioinformatics* 23, 2947–2948.
- Lenormand, T., Chevin, L.M., and Bataillon, T. (2016). Parallel evolution: what does it (not) tell us and why is it (still) interesting. In *Chance in Evolution*, (Chicago, Illinois: Chicago University Press), p.
- Liu, Z., Qi, F.-Y., Zhou, X., Ren, H.-Q., and Shi, P. (2014). Parallel sites implicate functional convergence of the hearing gene prestin among echolocating mammals. *Mol. Biol. Evol.* 31, 2415–2424.
- Maddamsetti, R., Hatcher, P.J., Cruveiller, S., Médigue, C., Barrick, J.E., and Lenski, R.E. (2015). Synonymous genetic variation in natural isolates of *Escherichia coli* does not predict where synonymous substitutions occur in a long-term experiment. *Mol. Biol. Evol.* msv161.
- Marvig, R.L., Sommer, L.M., Molin, S., and Johansen, H.K. (2015). Convergent evolution and adaptation of *Pseudomonas aeruginosa* within patients with cystic fibrosis. *Nat. Genet.* 47, 57–64.
- McVean, G.A.T., Myers, S.R., Hunt, S., Deloukas, P., Bentley, D.R., and Donnelly, P. (2004). The Fine-Scale Structure of Recombination Rate Variation in the Human Genome. *Science* 304, 581–584.
- Ness, R.W., Morgan, A.D., Vasanthakrishnan, R.B., Colegrave, N., and Keightley, P.D. (2015).

- Extensive de novo mutation rate variation between individuals and across the genome of *Chlamydomonas reinhardtii*. *Genome Res.* 25, 1739–1749.
- Pál, C., Papp, B., and Hurst, L.D. (2001). Highly expressed genes in yeast evolve slowly. *Genetics* 158, 927–931.
- Punta, M., Coghill, P.C., Eberhardt, R.Y., Mistry, J., Tate, J., Boursnell, C., Pang, N., Forslund, K., Ceric, G., Clements, J., et al. (2012). The Pfam protein families database. *Nucleic Acids Res.* 40, D290–D301.
- R Development Core Team (2014). R: a language and environment for statistical computing (R Foundation for Statistical Computing, Vienna, Austria).
- Self, S.G., and Liang, K.-Y. (1987). Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *J. Am. Stat. Assoc.* 82, 605–610.
- Sharp, P.M., and Li, W.-H. (1987). The Codon Adaptation Index - A measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.* 15, 1281–1295.
- Stark, C., Breitkreutz, B.-J., Reguly, T., Boucher, L., Breitkreutz, A., and Tyers, M. (2006). BioGRID: a general repository for interaction datasets. *Nucleic Acids Res.* 34, D535–539.
- Strathern, J.N., Shafer, B.K., and McGill, C.B. (1995). DNA Synthesis Errors Associated with Double-Strand-Break Repair. *Genetics* 140, 965–972.
- Streisfeld, M.A., and Rausher, M.D. (2011). Population genetics, pleiotropy, and the preferential fixation of mutations during adaptive evolution. *Evolution* 65, 629–642.
- Winzeler, E.A., Shoemaker, D.D., Astromoff, A., Liang, H., Anderson, K., Andre, B., Bangham, R., Benito, R., Boeke, J.D., Bussey, H., et al. (1999). Functional Characterization of the *S. cerevisiae* Genome by Gene Deletion and Parallel Analysis. *Science* 285, 901–906.
- Woods, R., Schneider, D., Winkworth, C.L., Riley, M.A., and Lenski, R.E. (2006). Tests of parallel molecular evolution in a long-term experiment with *Escherichia coli*. *Proc. Natl. Acad. Sci.* 103, 9107–9112.
- Yang, Z. (2007). PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol. Biol. Evol.* 24, 1586–1591.
- Zou, Z., and Zhang, J. (2015). Are convergent and parallel amino acid substitutions in protein evolution more prevalent than neutral expectations? *Mol. Biol. Evol.* 32, 2085–2096.
- Zuur, A.F. (2009). Mixed effects models and extensions in ecology with R (Springer).

FIGURES

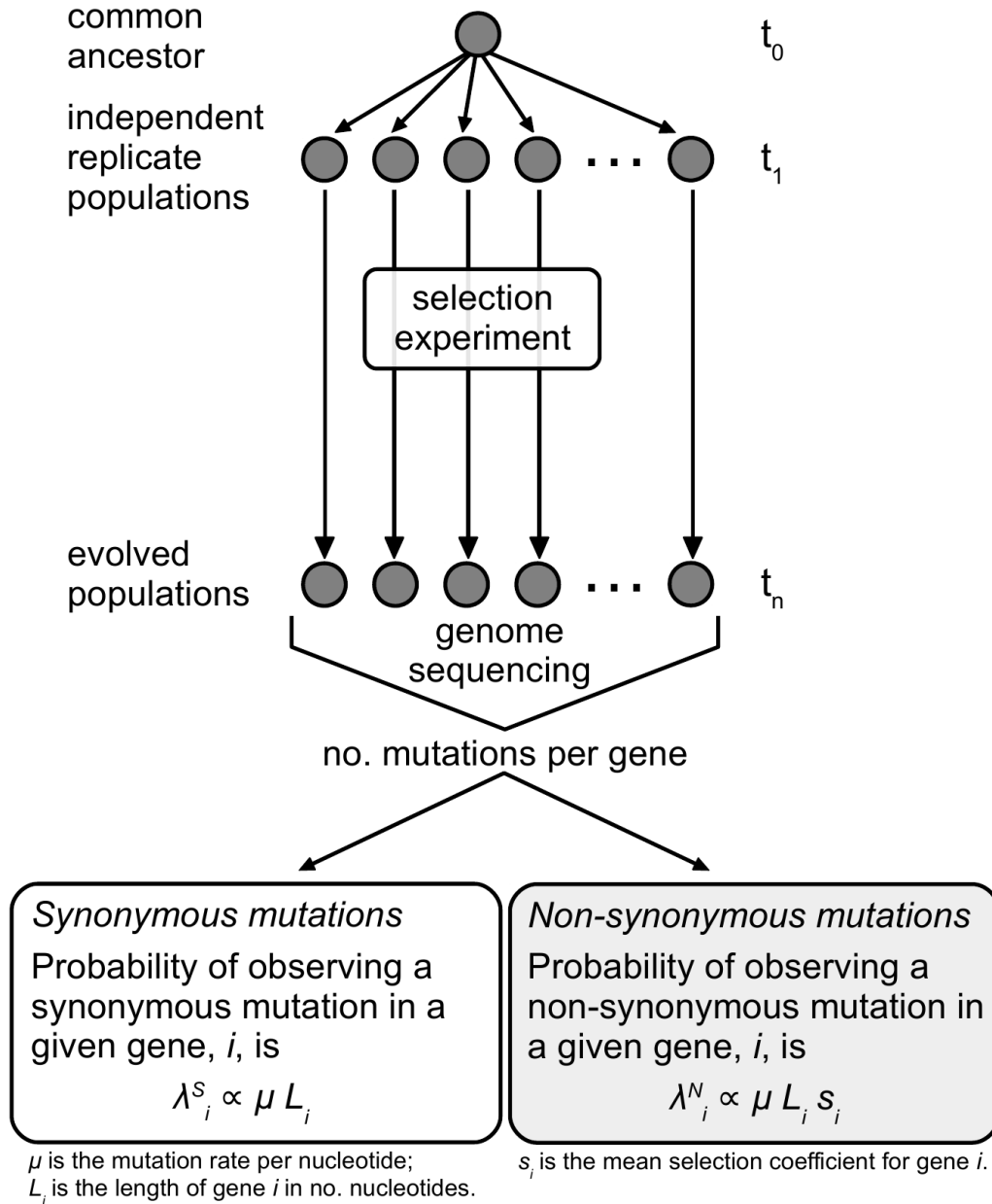
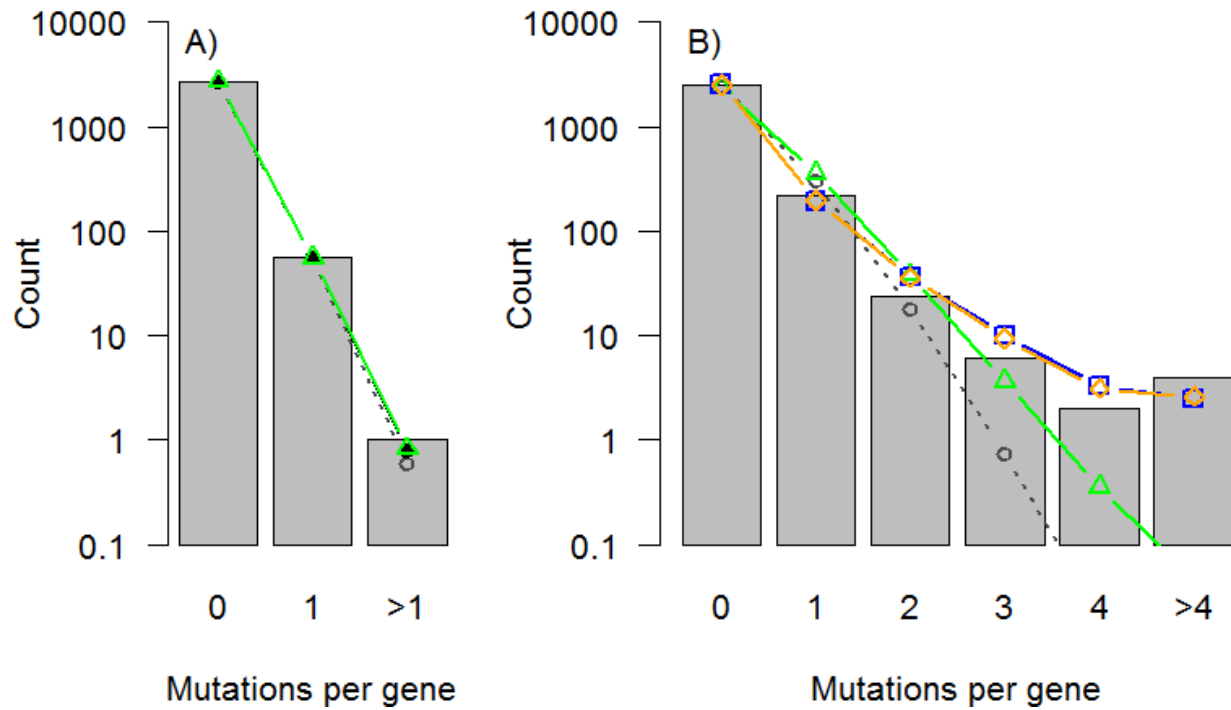


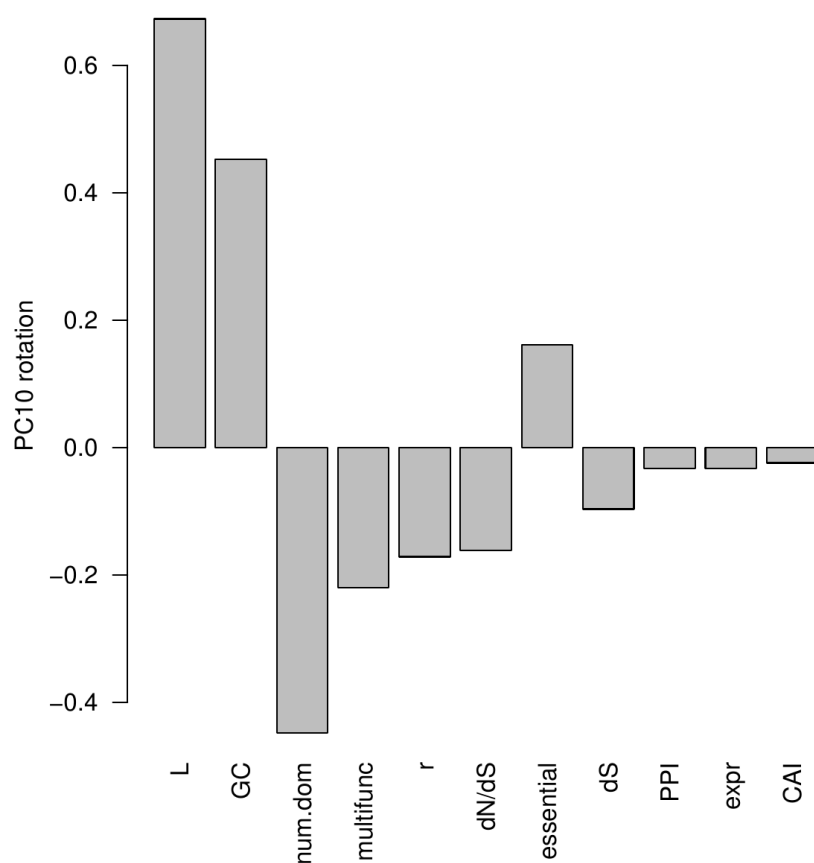
Figure 1: Schematic showing how the mutation counts data are generated and general assumptions underlying these data.



418

419

420 **Figure 2:** Distribution of A) synonymous and B) nonsynonymous mutations per gene and predicted
421 model distributions from M0.P (grey circles), M1.P (black points), M2.P (green triangles), and M_N.NB
422 (blue squares), and M_N.NB_{PC} (orange diamonds).



424 **Figure 3:** Loadings of the 11 genomic variables on PC10 – the only principal component that
425 significantly explains variation in nonsynonymous mutation counts. Genomic variables are ordered
426 from largest to smallest in terms of the absolute value of their loading.

TABLES

Table 1: Genomic variables used in this study.

Variable name	Description	Reference
<i>dS</i>	Number of synonymous substitutions per synonymous site, estimated from gene alignments of <i>S. cerevisiae</i> , <i>S. paradoxus</i> , <i>S. bayanus</i> , and <i>S. mikatae</i> (Cliften et al., 2003; Kellis et al., 2003).	Estimated for this study.
<i>dN/dS</i>	Number of nonsynonymous substitutions per nonsynonymous site, estimated from <i>S. cerevisiae</i> , <i>S. paradoxus</i> , <i>S. bayanus</i> , and <i>S. mikatae</i> (Cliften et al., 2003; Kellis et al., 2003).	Estimated for this study.
Gene length (<i>L</i>)	The number of nucleotides.	(Cherry et al., 2012)
% GC content (<i>GC</i>)	Percentage of nucleotides in the gene sequence that are either guanine or cytosine.	(Cherry et al., 2012)
Multi-functionality (<i>multifunc</i>)	Number of different GO slim categories assigned to a gene.	(Cherry et al., 2012)
Degree of protein-protein interaction (<i>PPI</i>)	The number of physical interactions reported by BioGRID (Stark et al., 2006).	(Koch et al., 2012)
Codon adaptation index (<i>CAI</i>)	A measure of bias in the usage of synonymous codons, based on a comparison between codon frequencies in the gene and frequencies observed in a set of highly expressed genes (Sharp and Li, 1987).	(Koch et al., 2012)
Number of domains (<i>num.dom</i>)	The number of regions that Pfam (Punta et al., 2012) has identified as domains in the protein sequence of each gene.	(Koch et al., 2012)
Level of expression (<i>expr</i>)	A measure of mRNA level for each gene when grown in standard lab conditions.	(Holstege et al., 1998)
Local recombination rate (<i>r</i>)	Mean recombination rate for a given gene calculated from recombination rate estimate at 0.5 kb intervals using <i>LDhat</i> (McVean et al., 2004).	(Illingworth et al., 2013)
Essential genes (<i>essential</i>)	A true/ false indicator variable denoting whether or not a gene is essential, based growth assays of deletion strains.	(Winzeler et al., 1999)

428 **Table 2:** 'M_S' models testing assumptions with the synonymous mutation data. Log-likelihoods, and
 429 AIC values are provided. The best model as determined by the lowest AIC with the fewest parameters
 430 is highlighted in grey.

431	432		log-lik.	No. param.	AIC
433	M _S 0.P:	Pois($\lambda_S = constant$)	-283.0	1	568.2
434	M _S 0.NB:	NB($\lambda_S = constant, \theta_S$)	-283.0	1	569.9
435					
436	M _S 1.P:	Pois($\lambda_S = constant * L_i^\alpha$)	-273.9	2	551.8
437	M _S 1.NB:	NB($\lambda_S = constant * L_i^\alpha, \theta_S$)	-273.9	3	553.8
438					
439	M _S 2.P:	Pois($\lambda_S = constant * L_i$)	-274.0	1	549.9
440	M _S 2.NB:	NB($\lambda_S = constant * L_i, \theta_S$)	-274.0	2	551.9

441
 442
 443 **Table 3:** 'M_N' models parameter estimates (*constant*, $\alpha 1$, $\alpha 2$, etc) and P-values for those estimates. Only
 444 those variables that significantly improved model fit are included.

445	446 M _N .NB: NB($\lambda_N = constant * L_i * L_i^{\alpha 1} * num.dom_i^{\alpha 2} * r_i^{\alpha 3}, \theta_N$)			
447		Estimate	P-value	
448	$\lambda_N \sim L$	$\alpha 1 = 0.4431$	0.001	
449	<i>num.dom</i>	$\alpha 2 = -0.4638$	0.004	
450	<i>r</i>	$\alpha 3 = 0.1033$	0.041	
451	<i>constant</i>	$8.015 * 10^{-6}$	<0.001	
452	θ_N	0.3877	<0.001	
453				
454	M _N .NB _{PC} : NB($\lambda_N = constant * L_i * \exp(PC10_i)^{\alpha 1}, \theta_N$)			
455		Estimate	P-value	
456	$\lambda_N \sim \exp(PC10)$	$\alpha 1 = 0.2984$	<0.001	
457	<i>constant</i>	$8.846 * 10^{-5}$	<0.001	
458	θ_N	0.3988	<0.001	

459
 460
 461 **Table 4:** Log-likelihoods, and AIC values for the 'M_N' models. The best model as determined by the
 462 lowest AIC with the fewest parameters is highlighted in grey.

463	464	log-lik.	No. param.	AIC
465	M _N .P	-1029.3	4	2066.5
466	M _N .P _{PC}	-1021.3	2	2046.6
467	M _N .NB	-953.4	5	1916.9
468	M _N .NB _{PC}	-956.8	3	1919.5