

<sup>1</sup> Unlocking a signal of introgression from codons in *Lachancea*  
<sup>2</sup> *kluveri* using a mutation-selection model

<sup>3</sup> Cedric Landerer \*

Department of Ecology & Evolutionary Biology, University of Tennessee, Knoxville TN 37996

National Institute for Mathematical and Biological Synthesis, Knoxville, TN 37996

Brian C. O'Meara

Department of Ecology & Evolutionary Biology, University of Tennessee, Knoxville TN 37996

National Institute for Mathematical and Biological Synthesis, Knoxville, TN 37996

Russell Zaretzki

Department of Business Analytics and Statistics, University of Tennessee, Knoxville TN 37996

National Institute for Mathematical and Biological Synthesis, Knoxville, TN 37996

Michael A. Gilchrist

Department of Ecology & Evolutionary Biology, University of Tennessee, Knoxville TN 37996

National Institute for Mathematical and Biological Synthesis, Knoxville, TN 37996

<sup>4</sup> June 14, 2019

---

\*Corresponding author: cedric.landerer@gmail.com

## Abstract

For decades, codon usage has been used as a measure of adaptation for translational efficiency of a gene's coding sequence. These patterns of codon usage reflect both the selective and mutational environment in which the coding sequences evolved. Over this same period, gene transfer between lineages has become widely recognized as an important biological phenomenon. Nevertheless, most studies of codon usage implicitly assume that all genes within a genome evolved under the same selective and mutational environment, an assumption violated when introgression occurs. In order to better understand the effects of introgression on codon usage patterns and vice versa, we examine the patterns of codon usage in *Lachancea kluyveri*, a yeast which has experienced a large introgression. We quantify the effects of mutation bias and selection for translation efficiency on the codon usage pattern of the endogenous and introgressed exogenous genes using a Bayesian mixture model, ROC SEMPPR, which is built on mechanistic assumptions of protein synthesis and grounded in population genetics. We find substantial differences in codon usage between the endogenous and exogenous genes, and show that these differences can be largely attributed to a shift in mutation bias favoring A/T ending codons in the endogenous genes to C/G ending codons in the exogenous genes. Recognizing the two different signatures of mutation bias and selection improves our ability to predict protein synthesis rate by 42% and allowed us to accurately assess endogenous codon preferences. In addition, using our estimates of mutation bias and selection, we identify *Eremothecium gossypii* as the closest relative to the exogenous genes, providing an alternative hypothesis about the origin of the exogenous genes, estimate the introgression occurred  $\sim 6 \times 10^8$  generation ago, and estimate its historic and current selection against mismatched codon usage. Together, our work illustrates the advantage of mechanistic, population genetic models like ROC SEMPPR and the quantitative estimates they provide when analyzing sequence data.

## 27 Introduction

28 Synonymous codon usage patterns varies within a genome and between taxa, reflecting differences in  
29 mutation bias, selection, and genetic drift. The signature of mutation bias is largely determined by the  
30 organism's internal or cellular environment, such as their DNA repair genes or UV exposure. While  
31 this mutation bias is an omnipresent evolutionary force, its impact can be obscured or amplified by  
32 selection. In contrast, the signature of selection on codon usage is largely determined by an organism's  
33 cellular environment alone, such as its tRNA species, their copy number, and their post-transcriptional  
34 modifications. The strength of selection on the codon usage of an individual gene is largely determined  
35 by its expression and synthesis rate which, in turn, is largely determined by the organism's external  
36 environment. In general, the strength of selection on codon usage increases with its expression level  
37 (Gouy and Gautier, 1982; Ikemura, 1985; Bulmer, 1990), specifically its protein synthesis rate (Gilchrist,  
38 2007). Thus as protein synthesis increases, codon usage shifts from a process dominated by mutation to  
39 a process dominated by selection. The overall efficacy of selection on codon usage is a function of the  
40 organism's effective population size  $N_e$  which, in turn, is largely determined by its external environment.  
41 ROC SEMPPR allows us disentangle the evolutionary forces responsible for the patterns of codon usage  
42 bias (CUB) encoded in an species' genome, by explicitly modeling the combined evolutionary forces of  
43 mutation, selection, and drift (Gilchrist, 2007; Shah and Gilchrist, 2011; Wallace *et al.*, 2013; Gilchrist  
44 *et al.*, 2015). In turn, these evolutionary forces should provide biologically meaningful information about  
45 the lineage's historical cellular and external environment.

46 Most studies implicitly assume that the CUB of a genome is shaped by a single cellular environment.  
47 As genes are horizontally transferred, introgress, or combined to form novel hybrid species, one would  
48 expect to see the influence of multiple cellular environments on a genomes codon usage pattern (Médigue  
49 *et al.*, 1991; Lawrence and Ochman, 1997). Given that transferred genes are likely to be less adapted than  
50 endogenous genes to their new cellular environment, we expect a greater selection against mismatched  
51 codon usage in transferred genes if donor and recipient environment differ greatly in their selection bias,  
52 making such transfers less likely. More practically, if differences in codon usage of transferred genes are  
53 unaccounted for, they may distort the interpretation of codon usage patterns. Such distortion could  
54 lead to the wrong inference of codon preference for an amino acid (Shah and Gilchrist, 2011; Gilchrist  
55 *et al.*, 2015), underestimate the variation in protein synthesis rate, or influence mutation estimates when  
56 analyzing a genome.

57 To illustrate these ideas, we analyze the CUB of the genome of *Lachancea kluyveri*, which is sister  
58 to all other Lachancea species. The Lachancea clade diverged from the *Saccharomyces* clade, prior to  
59 its whole genome duplication  $\sim 100$  Mya ago (Marcet-Houben and Gabaldón, 2015; Beimforde *et al.*,

60 Since that time, *L. kluyveri* has experienced a large introgression of exogenous genes which is  
61 found in all of its populations (Friedrich *et al.*, 2015), but in no other known Lachancea species (Vakirlis  
62 *et al.*, 2016). The introgression replaced the left arm of the C chromosome and displays a 13% higher GC  
63 content than the endogenous *L. kluyveri* genome (Payen *et al.*, 2009; Friedrich *et al.*, 2015). Previous  
64 studies suggest that the source of the introgression is likely a currently unknown or potentially extinct  
65 Lachancea lineage based on gene concatenation or synteny relationships (Payen *et al.*, 2009; Friedrich  
66 *et al.*, 2015; Vakirlis *et al.*, 2016; Brion *et al.*, 2017). These characteristics make *L. kluyveri* an ideal  
67 model to study the effects of an introgressed cellular environment and the resulting mismatch in codon  
68 usage.

69 Using ROC SEMPPR, a Bayesian population genetics model based on a mechanistic description of  
70 ribosome movement along an mRNA, allows us to quantify the cellular environment in which genes have  
71 evolved by separately estimating the effects of mutation bias and selection bias on codon usage. ROC  
72 SEMPPR's resulting predictions of protein synthesis rates have been shown to be on par with laboratory  
73 measurements (Shah and Gilchrist, 2011; Gilchrist *et al.*, 2015). In contrast to often used heuristic ap-  
74 proaches to study codon usage (Sharp and Li, 1987; Wright, 1990; dos Reis *et al.*, 2004), ROC SEMPPR  
75 explicitly incorporates and distinguishes between mutation and selection effects on codon usage and prop-  
76 erly weights by amino acid usage (Cope *et al.*, 2018). We use ROC SEMPPR to independently describe  
77 two cellular environments reflected in the *L. kluyveri* genome; the signature of the current environment  
78 in the endogenous genes and the decaying signature of the exogenous environment in the introgressed  
79 genes. Our results indicate that the difference in GC content between endogenous and exogenous genes  
80 is mostly due to the differences in mutation bias of their ancestral environments. Accounting for these  
81 different signatures of mutation bias and selection bias of the endogenous and exogenous sets of genes  
82 substantially improves our ability to predict present day protein synthesis rates. These endogenous and  
83 exogenous gene set specific estimates of mutation bias and selection bias, in turn, allow us to address  
84 more refined questions of biological importance. For example, they allow us to provide an alternative  
85 hypothesis about the origin of the introgression and identify *E. gossypii* as the nearest sampled relative  
86 of the source of the introgressed genes out of the 332 budding yeast lineages with sequenced genomes  
87 (Shen *et al.*, 2018). While this hypothesis is in contrast previous work (Payen *et al.*, 2009; Friedrich *et al.*,  
88 2015; Vakirlis *et al.*, 2016; Brion *et al.*, 2017), we find support for it in gene trees and synteny. We also  
89 estimate the age of the introgression to be on the order of 0.2 - 1.7 Mya, estimate the selection against  
90 these genes, both at the time of introgression and now, and predict a detectable signature of CUB to  
91 persist in the introgressed genes for another 0.3 - 2.8 Mya, highlighting the sensitivity of our approach.

Table 1: Model selection of the two competing hypothesis. Combined: mutation bias and selection bias for synonymous codons is shared between endogenous and exogenous genes. Separated: mutation bias and selection bias for synonymous codons is allowed to vary between endogenous and exogenous genes. Reported are the log-likelihood,  $\log(\mathcal{L})$ , the number of parameters estimated  $n$ , the log-marginal likelihood  $\log(\mathcal{L}_M)$ , and Bayes Factor  $K$ .

Hypothesis	$\log(\mathcal{L})$	$n$	$\log(\mathcal{L}_M)$	$\log(K)$
Combined	-2,650,047	5,483	-2,657,582	—
Separated	-2,612,397	5,402	-2,615,288	42,294

## Results

### The Signatures of two Cellular Environments within *L. kluyveri*'s Genome

We used our software package AnaCoDa (Landerer *et al.*, 2018) to compare model fits of ROC SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two sets of 4,864 endogenous and 497 exogenous genes. ROC SEMPPR is a statistical model that relates the effects of mutation bias  $\Delta M$  and selection bias  $\Delta\eta$  between synonymous codons, and protein synthesis rate  $\phi$  to explain the observed codon usage patterns. Bayes factor strongly support the hypothesis that the *L. kluyveri* genome consists of genes with two different and distinct patterns of codon usage bias rather than a single ( $K = \exp(42,294)$ ; Table 1). We find additional support for this hypothesis when we compare our predictions of protein synthesis rate to empirically observed mRNA expression values as proxy for protein synthesis. Specifically, the explanatory power between our predictions and observed values improved by  $\sim 42\%$ , from  $R^2 = 0.33$  to 0.46 (Figure 1).

### Comparing Differences in the Endogenous and Exogenous Codon Usage

To better understand the differences in the endogenous and exogenous cellular environments, we compared our parameter estimates of mutation bias  $\Delta M$  and selection  $\Delta\eta$  for the two sets of genes. Our estimates of  $\Delta M$  for the endogenous and exogenous genes were negatively correlated ( $\rho = -0.49$ ), indicating weak similarity with only  $\sim 5\%$  of the codons share the same sign between the two mutation environments (Figure 2a). Overall, the endogenous genes only show a selection preference for C and G ending codons in  $\sim 58\%$  of the codon families. In contrast, the exogenous genes display a strong preference for A and T ending codons in  $\sim 89\%$  of the codon families.

For example, the endogenous genes show a mutational bias for A and T ending codons in  $\sim 95\%$  of the codon families (the exception being Phe, F). The exogenous genes display an equally consistent mutational bias towards C and G ending codons (Table S1). In contrast to  $\Delta M$ , our estimates of  $\Delta\eta$  for the endogenous and exogenous genes were positively correlated ( $\rho = 0.69$ ) and showing the same sign in  $\sim 53\%$  of codons between the two selection environments (Figure 2). ROC SEMPPR constraints

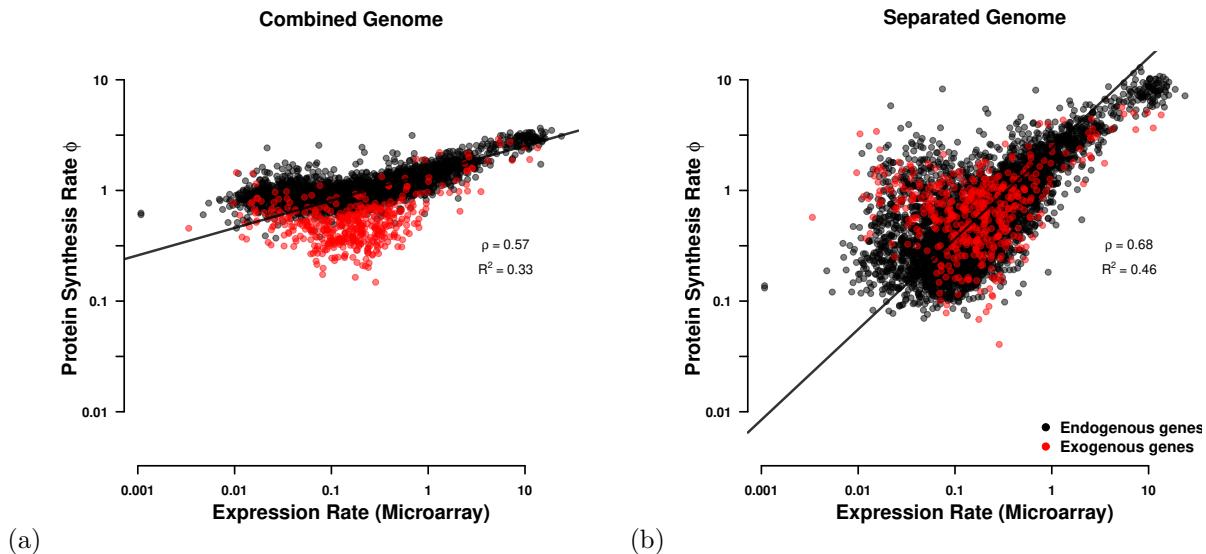


Figure 1: Comparison of predicted protein synthesis rate  $\phi$  to microarray data from Tsankov *et al.* (2010) for (a) the combined genome and (b) the separated endogenous and exogenous genes. Endogenous genes are displayed in black and exogenous genes in red. Black line indicates type II regression line assuming noise in the dependent and independent variable (Sokal and Rohlf, 1981).

117      $E[\phi] = 1$ , allowing us to interpret  $\Delta\eta$  as selection on codon usage of the average gene with  $\phi = 1$  and  
118     gives us the ability to compare the efficacy of selection  $sN_e$  across genomes.

119     We find that the efficacy of selection within each codon family differs between sets of genes. The  
120     difference in codon usage between endogenous and exogenous genes is striking as some amino acids have  
121     opposite codon preferences. As a result, our estimates of the optimal codon differ in nine cases between  
122     endogenous and exogenous genes (Figure 3, Table S2). For example, the usage of the Asparagine (Asn,  
123     N) codon AAC is increased in highly expressed endogenous genes but the same codon is depleted in highly  
124     expressed exogenous genes. For Aspartic acid (Asp, D), the combined genome shows the same codon  
125     preference in highly expressed genes as the exogenous gene set. Generally, fits to the complete *L. kluyveri*  
126     genome reveal that the relatively small exogenous gene set ( $\sim 10\%$  of genes) has a disproportional effect  
127     on the model fit (Figure S1, S2).

128     Of the nine cases in which the endogenous and exogenous genes show differences in the selectively  
129     most favored codon five cases (Asp, D; His, H; Lys, K; Asn, N; and Pro, P) the endogenous genes favor  
130     the codon with the most abundant tRNA. For the remaining four cases (Ile, I; Ser, S; Thr, T; and Val, V),  
131     there are no tRNA genes for the wobble free cognate codon encoded in the *L. kluyveri* genome. However,  
132     the codon preference of these four amino acids in the exogenous genes matches the most abundant tRNA  
133     encoded in the *L. kluyveri* genome.

134     The effect of the small exogenous gene set on the fit to the complete *L. kluyveri* genome is smaller

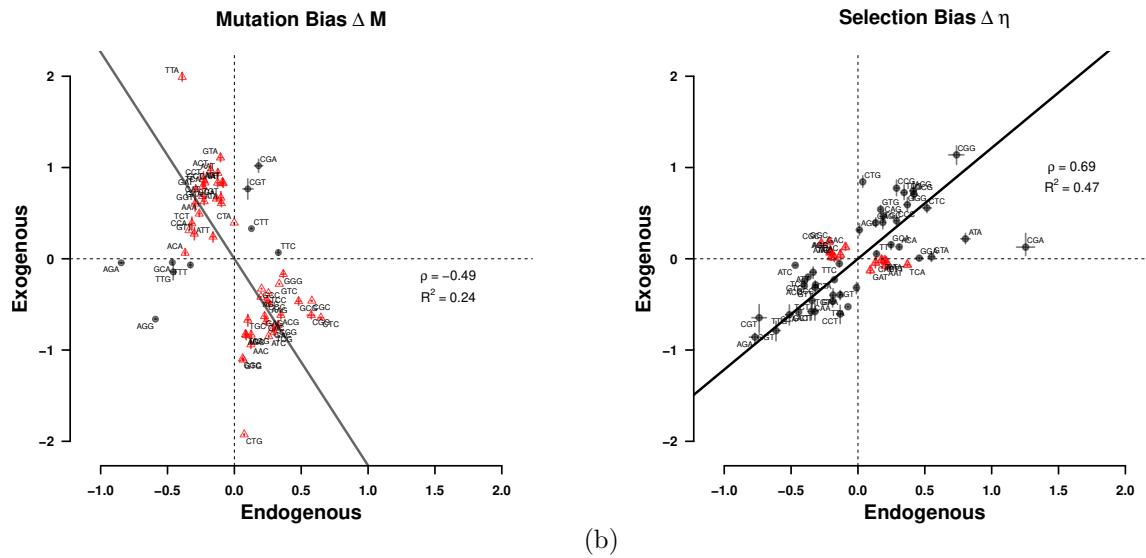


Figure 2: Comparison of (a) mutation bias  $\Delta M$  and (b) selection bias  $\Delta\eta$  parameters for endogenous and exogenous genes. Estimates are relative to the mean for each codon family. Black dots indicate  $\Delta M$  or  $\Delta\eta$  parameters with the same sign for the endogenous and exogenous genes, red dots indicate parameters with different signs. Black line indicates type II regression line assuming noise in the dependent and independent variable (Sokal and Rohlf, 1981). Dashed lines mark quadrants.

in our estimates of selection bias  $\Delta\eta$  than  $\Delta M$ , but still large. We find that the complete *L. kluyveri* genome is estimated to share the selection preference with the exogenous genes in  $\sim 60\%$  of codon families that show dissimilarity between endogenous and exogenous genes. We find that the complete *L. kluyveri* genome fit shares mutational preference with the exogenous genes in  $\sim 78\%$  of the 19 codon families showing a difference in mutational codon preference between the endogenous and exogenous genes. In two cases, Isoleucine (Ile, I) and Arginine (Arg, R), the strong dissimilarity in mutation preference results in an estimated codon preference in the complete *L. kluyveri* genome that differs from both the endogenous, and the exogenous genes. These results clearly show that it is important to recognize the difference in endogenous and exogenous genes and treat these genes as separate sets to avoid the inference of incorrect synonymous codon preferences and better predict protein synthesis.

## Determining Source of Exogenous Genes

We combined our estimates of mutation bias  $\Delta M$  and selection bias  $\Delta \eta$  with synteny information and searched for potential source lineages of the introgressed exogenous region. We examined 332 budding yeasts (Shen *et al.*, 2018) and, identified the ten lineages with the highest correlation for the  $\Delta M$  parameters as potential source lineages (Figure 4, Table 2). We used  $\Delta M$  to identify candidate lineages as the endogenous and exogenous genes show greater dissimilarity in mutation bias than in selection

## Endogenous and Exogenous Codon Usage

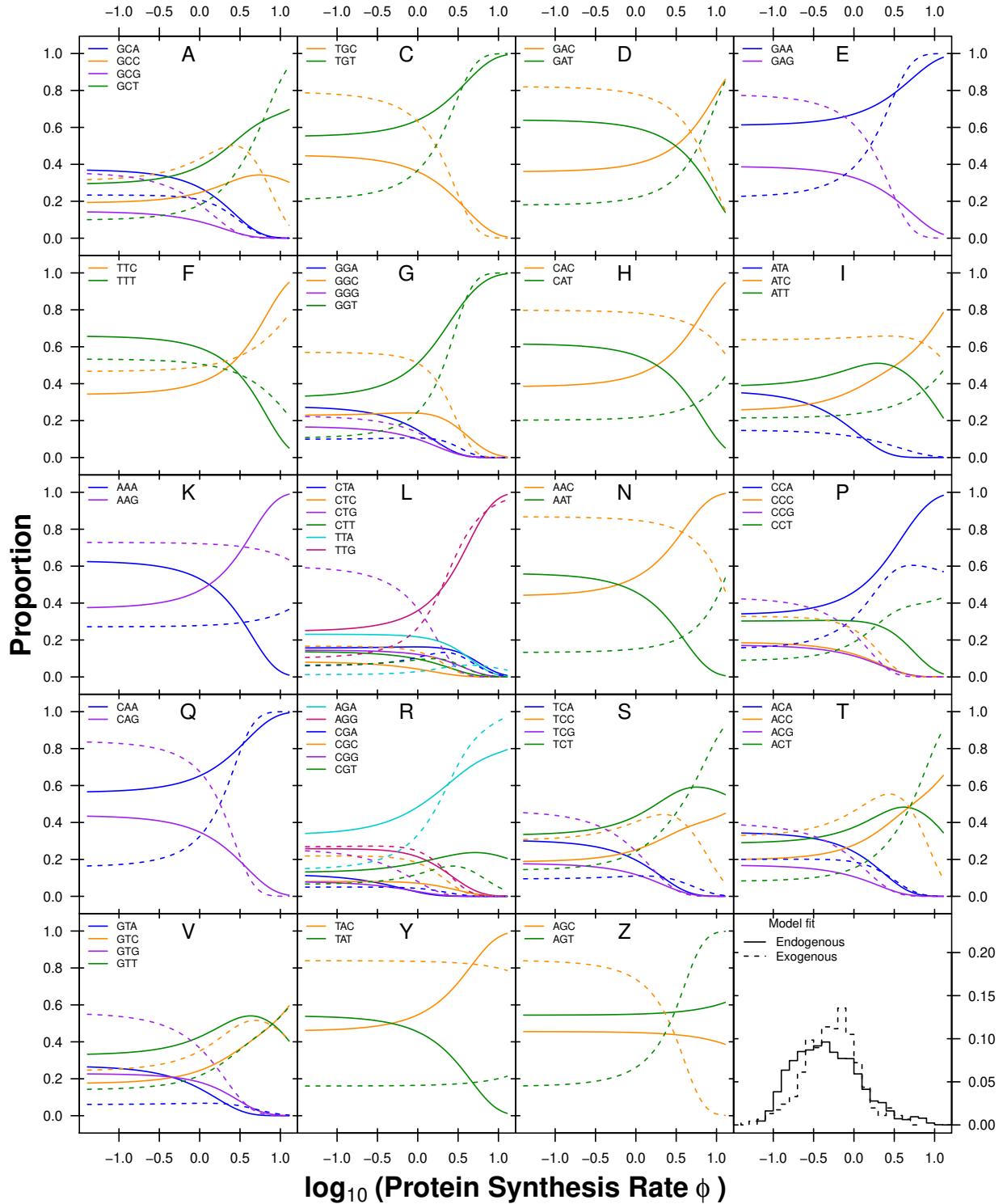


Figure 3: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dashed line indicates the exogenous codon usage.

Table 2: Budding yeast lineages showing similarity in codon usage with the exogenous genes.  $\rho_{\Delta M}$  and  $\rho_{\Delta \eta}$  represent the Pearson correlation coefficient for  $\Delta M$  and  $\Delta \eta$ , respectively. GC content is the average GC content of the whole genome. Synteny is the percentage of the exogenous genes found in the listed lineage. Only one lineage (*E. gossypii*) shows a similar GC content > 50%.

Species	$\rho_{\Delta M}$	$\rho_{\Delta \eta}$	GC content	Synteny %	Distance [Mya]
<i>Eremothecium gossypii</i>	0.89	0.70	51.7	75	211.0847
<i>Danielozyma ontarioensis</i>	0.75	0.92	46.6	3	470.1043
<i>Metschnikowia shivogae</i>	0.86	0.87	49.8	0	470.1043
<i>Babjeviella inositovora</i>	0.83	0.78	48.1	0	470.1044
<i>Ogataea zsoltii</i>	0.75	0.85	47.7	0	470.1042
<i>Metschnikowia hawaiiensis</i>	0.80	0.86	44.4	0	470.1042
<i>Candida succiphila</i>	0.85	0.83	40.9	0	470.1042
<i>Middlehovenomyces tepae</i>	0.80	0.62	40.8	0	651.9618
<i>Candida albicans</i> *	0.84	0.75	33.7	0	470.1043
<i>Candida dubliniensis</i> *	0.78	0.75	33.1	0	470.1043

\* Lineages use the alternative yeast nuclear code

bias. Two of the ten candidate lineages utilize the alternative yeast nuclear code (NCBI codon table 12). In this case, the codon CTG codes for Serine instead of Leucine. We therefore excluded the Leucine codon family in our comparison of codon families, however, there was no need to exclude Serine as well as CTG is not a one step neighbor of the remaining Serine codons. The endogenous *L. kluyveri* genome exhibits codon usage very similar to most (77 %) yeast lineages examined, indicating that most of the examined yeasts share a similar codon usage (Figure S3). Only  $\sim 17\%$  of all examined yeast show a positive correlation in both,  $\Delta M$  and  $\Delta \eta$  with the exogenous genes, whereas the vast majority of lineages ( $\sim 83\%$ ) show a negative correlation for  $\Delta M$ , only 21 % show a negative correlation for  $\Delta \eta$ .

Comparing synteny between the exogenous genes, which are restricted to the left arm of chromosome C, and the determined candidate yeast species we find that *E. gossypii* is the only species that displays high synteny (Table 2). Furthermore, the synteny relationship between the exogenous region and other yeasts appears to be limited to Saccharomycetaceae clade. Given these results, we conclude that of the 332 examined yeast lineages the *E. gossypii* lineage is the most likely source of the introgressed exogenous genes. This result is in contrast to previous studies which studied the exogenous genes and chromosome recombination in the Lachancea clade and concluded that the introgressed region originated from within the Lachancea clade (Payen *et al.*, 2009; Friedrich *et al.*, 2015; Vakirlis *et al.*, 2016). To validate our results, we identified 121 genes in our dataset (Shen *et al.*, 2018) with homologous gene in *E. gossypii* and *L. thermotolerance* and used IQTree (Nguyen *et al.*, 2015) to infer the phylogenetic relationship of the exogenous genes. Our results show that  $\sim 60\%$  of exogenous genes (73/121) are more closely related to *E. gossypii* than to other Lachancea. Interestingly, our results also indicate that codon usage does not necessarily correlate with phylogenetic distance (Table 2).

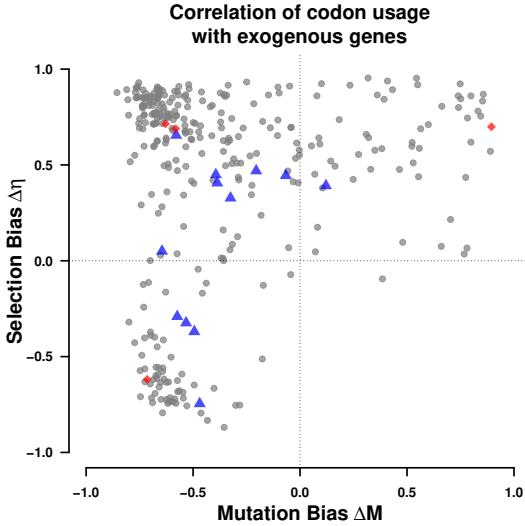


Figure 4: Correlation coefficients of  $\Delta M$  and  $\Delta \eta$  of the exogenous genes with 332 examined budding yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta \eta$  of the lineages with the exogenous parameter estimates. Blue triangles indicate the *Lachancea* and red diamonds indicate *Eremothecium* species. All regressions were performed using a type II regression assuming noise in the dependent and independent variable (Sokal and Rohlf, 1981).

## 172 Estimating Introgression Age

173 We modeled the change in codon frequency over time as exponential decay, and estimated the age of the  
 174 introgression assuming that *E. gossypii* still represents the mutation bias of its ancestral source lineage  
 175 at the time of the introgression and a constant mutation rate. We infer the age of the introgression to  
 176 be on the order of  $6.2 \pm 1.2 \times 10^8$  generations. Assuming *L. kluyveri* experiences between one and eight  
 177 generations per day, we estimate the introgression to have occurred between 212,000 to 1,700,000 years  
 178 ago. Our estimate places the time of the introgression earlier than the previous estimate of 19,000 -  
 179 150,000 years by Friedrich *et al.* (2015).

180 Using our model of exponential decay model, we also estimated the persistence of the signal of the  
 181 exogenous cellular environment. We predict that the  $\Delta M$  signal of the source cellular environment will  
 182 have decayed to be within one percent of the *L. kluyveri* environment in  $\sim 5.4 \pm 0.2 \times 10^9$  generations, or  
 183 between 1,800,000 and 15,000,000 years. Together, these results indicate that the mutation signature  
 184 of the exogenous genes will persist for a very long time.

## 185 Estimating Selection against Codon Mismatch of the Exogenous Genes

186 We define the selection against inefficient codon usage as the difference between the fitness on the log  
 187 scale of an expected, replaced endogenous gene and the exogenous gene,  $s \propto \phi \Delta \eta$  due to the mismatch in

codon usage parameters (See Methods for details). As the introgression occurred before the diversification of *L. kluyveri* and has fixed throughout all populations (Friedrich *et al.*, 2015), we can not observe the original endogenous sequences that have been replaced by the introgression. Using our estimates of  $\Delta M$  and  $\Delta\eta$  from the endogenous genes and assuming the current exogenous amino acid composition of genes is representative of the replaced endogenous genes, we estimate the selection against the exogenous genes at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of selection bias for the exogenous genes show that, while well correlated with the endogenous genes, only nine amino acids share the same selectively preferred codon. Exogenous genes are, therefore, expected to represent a significant reduction in fitness for *L. kluyveri* due to mismatch in codon usage. We find that the selection against the exogenous genes  $\Delta s$  due to mismatched codon usage was -0.0008 at the time of the introgression and still is -0.0003 today. Based on the selection against the codon mismatch at the time of the introgression  $\Delta s = -0.0008$  and an effective population size  $N_e$  on the order of  $10^7$  (Wagner, 2005) we approximate a fixation probability of  $(1 - \exp[-\Delta s])/(1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$  (Sella and Hirsh, 2005) for the exogenous genes.

In order to account for differences in the efficacy of selection on codon usage either due to the cost of pausing, differences in the effective population size or the decline in fitness with every ATP wasted between the donor lineage and *L. kluyveri* we added a linear scaling factor  $\kappa$  to scale our estimates of  $\Delta\eta$  between the donor lineage and *L. kluyveri* (See Methods for details). We predict that a small number of low expression genes ( $\phi < 1$ ) were weakly exapted at the time of the introgression (Figure 5a). High expression genes ( $\phi > 1$ ) are predicted to have faced the largest selection against their mismatched codon usage in the novel cellular environment. These highly expressed genes are inferred to have the greatest degree of adaptation since the time of the introgression to the *L. kluyveri* cellular environment (Figures 5a & S6).

## Discussion

In order to study the evolutionary effects of the large scale introgression of the left arm of chromosome C, we used ROC SEMPPR, a mechanistic model of ribosome movement along an mRNA. The usage of a mechanistic model rooted in population genetics allows us generate more nuanced quantitative parameter estimates and separate the effects of mutation and selection on the evolution of codon usage. This allowed us to calculate the selection against the introgression, and provides *E. gossypii* as a potential source lineage of the introgression which was previously not considered. Our parameter estimates indicate that the *L. kluyveri* genome contains distinct signatures of mutation and selection bias from both an endogenous and exogenous cellular environment. By fitting ROC SEMPPR separately to *L. kluyveri*'s

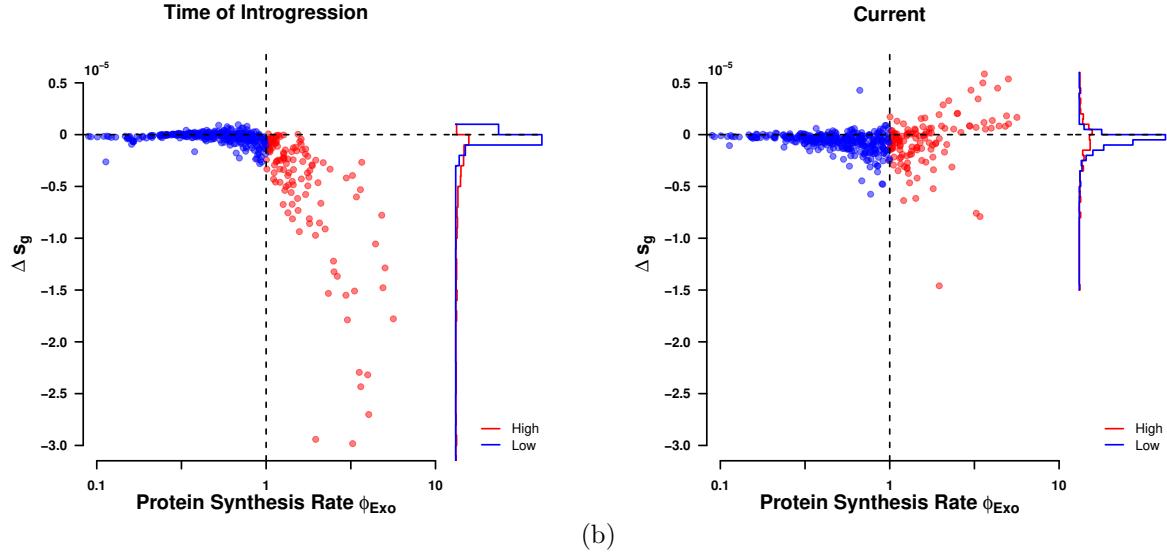


Figure 5: Selection against mismatched codon usage  $s = \Delta\eta\phi$  (a) at the time of introgression ( $\kappa = 5$ ), and (b) currently ( $\kappa = 1$ ). Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates neutrality.

endogenous and exogenous sets of genes we generate a quantitative description of their signatures of mutation bias and natural selection for efficient protein translation.

Previous work by Payen *et al.* (2009) showed an increased preference for GC rich codons in the exogenous genes but our results provide more nuanced insights by separating the effects of mutation bias and selection. We are able to show that the difference in GC content between endogenous and exogenous genes is mostly due to differences in mutation bias as 95% of exogenous codon families show a strong mutation bias towards GC ending codons (Table S1). However, the exogenous genes show a selective preference for AT ending codons for 90% of codon families (Table S2). Acknowledging the increased mutation bias towards GC ending codons and the difference in strength of selection between endogenous and exogenous genes by separating them also improves our estimates of protein synthesis rate  $\phi$  by 42% relative to the full genome estimate ( $R^2 = 0.46$  vs. 0.32, respectively).

The mutation and selection bias parameters  $\Delta M$  and  $\Delta\eta$  of the introgressed exogenous genes contain information, albeit decaying, about its previous cellular environment. We selected the top ten lineages with the highest similarity in  $\Delta M$  to see if our parameter estimates would allow us to identify a potential source lineage. The synteny relationship of these lineages with the exogenous genes was calculated as a point of comparison as it provides orthogonal information to our parameter estimates. Synteny with the exogenous genes is limited to the Saccharomycetaceae clade, excluding all of the potential source lineages identified using codon usage but *E. gossypii* (Table 2). Interestingly, this also showed that similarity in

238 codon usage does not correlate with phylogenetic distance.

239 Previous work indicated that the donor lineage of the exogenous genes has to be a, potentially un-  
240 known, Lachancea lineage (Payen *et al.*, 2009; Friedrich *et al.*, 2015; Vakirlis *et al.*, 2016; Brion *et al.*,  
241 2017). These previous results, however, are based on species rather than genes trees ignoring the dif-  
242 ferential adaptation rate to their novel cellular environment between genes or due not consider lineages  
243 outside of the Lachancea clade. Considering the similarity in selection bias (Figure 2b) and our calcula-  
244 tion of selection on the exogenous genes (Figure 5b), both of which are free of any assumption about the  
245 origin of the exogenous genes, a species tree estimated from the exogenous genes may be biased towards  
246 the Lachancea clade. Estimating individual gene trees rather than relying on a species tree provided  
247 further evidence that the exogenous genes could originate from a lineage that does not belong to the  
248 Lachancea clade. As we highlighted in this study, relatively small sets of genes with a signal of a foreign  
249 cellular environment can significantly bias the outcome of a study. The same holds true for phylogenetic  
250 inferences (Salichos and Rokas, 2013), and as we showed the signal of the original endogenous cellular  
251 environment that shaped CUB is at different stages of decay in high and low expression genes (Figure  
252 S6). In summary, our work does not dispute an unknown Lachancea as possible origin, but provides  
253 an alternative hypothesis based on the codon usage of the exogenous genes, phylogenetic analysis, and  
254 synteny.

255 We find large selection against the exogenous genes due to mismatched codon usage, resulting in  
256 an astronomically small fixation probability, and making fixation extremely unlikely. With such a small  
257 fixation probability,  $((1 - \exp[-\Delta s]) / (1 - \exp[-2\Delta s N_e]) \approx 10^{-6952})$  one or more advantageous loci had to  
258 drive the introgressed genes to fixation. Even though *L. kluyveri* diverged from the rest of the Lachancea  
259 clade around 85 Mya (Kensche *et al.*, 2008; Marcket-Houben and Gabaldón, 2015), if we assume 1 to 8  
260 generations/day, which implies  $10^{10}$  to  $10^{11}$  generations since the time of divergence, one round of meiosis  
261 for every 1000 rounds of mitosis of which only one in 100 meiosis events lead to outcrossing based on  
262 *S. paradoxus* (Tsai *et al.*, 2008), and  $N_e \approx 10^8$  there were only  $10^{13}$  to  $10^{14}$  opportunities for such an  
263 introgression to have occurred and fixed. In this context, the astronomically small fixation probability  
264 indicates that there was virtually no chance of fixation and was likely not a fluke, and unless there was  
265 a severe bottleneck with  $N_e < 1/|s| \approx 1,250$  around the time of introgression, which conceivably could  
266 have been triggered by a speciation event, this scenario seems, thus, very unlikely.

267 If, on the other hand, the introgressed region contained advantageous loci, the effective selection  
268 against the exogenous genes would be reduced and the introgression could have even been advantageous.  
269 Indeed,  $\sim 30\%$  of low expression exogenous genes ( $\phi < 1$ ) appeared to be exapted at the time of the  
270 introgression. This exaptation is due to the mutation bias in the endogenous genes matching the selection  
271 bias in the exogenous genes for GC ending codons. These exapted loci, however, are already accounted

for in our calculation of  $\Delta s$ . Thus, loci that, at the time of the introgression, provided a benefit to *L. kluyveri* must have been contained in the introgressed region Alternatively, an environmental shift may have caused the replaced endogenous region to suddenly contain deleterious loci making the reduction in translation efficiency the lesser of two evil.

However, one may wonder why recombination events did not limit the introgression to only the adaptive loci. A potential answer is the low recombination rate between the endogenous and exogenous regions Payen *et al.* (2009); Brion *et al.* (2017) . This is presumably due to the dissimilarity in GC content and/or a lower than average sequence homology between the exogenous region and the one it replaced. Compatible with this explanation is the possibility of several highly advantageous loci distributed across the region which then drove a rapid selective sweep and/or the population through a bottleneck speciation process.

Assuming *E. gossypii* as potential source lineage of the introgressed region, we illustrated how information on codon usage can be used to infer the time since the introgression occurred using our estimates of mutation bias  $\Delta M$ . The  $\Delta M$  estimates are well suited for this task as they are free of the influence of selection and unbiased by  $N_e$  and other scaling terms, which is in contrast to our estimates of  $\Delta \eta$  (Gilchrist *et al.*, 2015). Our estimated age of the introgression of  $6.2 \pm 1.2 \times 10^8$  generations is  $\sim 10$  times longer than a previous minimum estimate by Friedrich *et al.* (2015) of  $5.6 \times 10^7$  generations, which was based on the effective population recombination rate and the population mutation parameter (Ruderfer *et al.*, 2006). Furthermore, these estimates assume that the current *E. gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the time of the introgression. Thus, If the ancestral mutation environments were more similar (dissimilar) at the time of the introgression then our result is an overestimate (underestimate).

Overall, our results show the usefulness of the separation of mutation bias and selection bias and the importance of recognizing the presence of multiple cellular environments in the study of codon usage. We also illustrate how a mechanistic model like ROC SEMPPR and the quantitative estimates it provides can be used for more sophisticated hypothesis testing in the future. In contrast to other approaches used to study codon usage like CAI (Sharp and Li, 1987) or tAI (dos Reis *et al.*, 2004), ROC SEMPPR incorporates the effects of mutation bias and amino acid composition explicitly (Cope *et al.*, 2018). We highlight potential issues when estimating codon preferences, as estimates can be biased by the signature of a second, historical cellular environment. In addition, we show how quantitative estimates of mutation bias and selection relative to drift can be obtained from codon data and used to infer the fitness cost of an introgression as well as its history and potential future.

304

## Materials and Methods

305

### Separating Endogenous and Exogenous Genes

306

A GC-rich region was identified by Payen *et al.* (2009) in the *L. kluyveri* genome extending from position 1 to 989,693 of chromosome C. This region was later identified as an introgression by Friedrich *et al.* (2015). We obtained the *L. kluyveri* genome from SGD Project <http://www.yeastgenome.org/download-data/> (on 09-27-2014) and the annotation for *L. kluyveri* NRRL Y-12651 (assembly ASM14922v1) from NCBI (on 12-09-2014). We assigned 457 genes located on chromosome C with a location within the ~ 1 Mb window to the exogenous gene set. All other 4864 genes of the *L. kluyveri* genome were assigned to the exogenous genes.

313

### Model Fitting with ROC SEMPPR

314

ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) (Landerer *et al.*, 2018) and R (3.4.1) (R Core Team, 2013). ROC SEMPPR was run from 10 different starting values for at least 250,000 iterations and thinned to every 50th iteration. After manual inspection to verify that the MCMC had converged, parameter posterior means, log posterior probability and log likelihood were estimated from the last 500 samples (last 10% of samples).

319

### Model selection

320

The marginal likelihood of the combined and separated model fits was calculated using a generalized harmonic mean estimator (Gronau *et al.*, 2017). A variance scaling of 1.1 was used to scale the important density of the estimator. Using the estimated marginal likelihoods, we calculated the Bayes factor to assess model performance. Increases in the variance scaling increase the estimated Bayes factor, therefore we report a conservative Bayes factor bases on a small variance scaling S7.

325

### Comparing Codon Specific Parameter Estimates and Selecting Candidate lineages

326

As the choice of reference codon can reorganize codon families coding for an amino acid relative to each other, all parameter estimates were interpreted relative to the mean for each codon family.

$$\Delta M_i = \Delta M_{i,1} - \overline{\Delta M_i} \quad (1)$$

$$\Delta \eta_i = \Delta \eta_{i,1} - \overline{\Delta \eta_i} \quad (2)$$

327 Comparison of codon specific parameters ( $\Delta M$  and  $\Delta\eta = 2N_e q(\eta_i - \eta_j)$ ) was performed using the function  
328 lmodel2 in the R package lmodel2 (1.7.3) (Legendre, 2018) and R version 3.4.1 (R Core Team, 2013).  
329 The parameter  $\Delta\eta$  can be interpreted as the difference in fitness between codon  $i$  and  $j$  for the average  
330 gene with  $\phi = 1$  scaled by the effective population size  $N_e$ , and the selective cost of an ATP  $q$  (Gilchrist,  
331 2007; Gilchrist *et al.*, 2015). Type II regression was performed with re-centered parameter estimates,  
332 accounting for noise in dependent and independent variable (Sokal and Rohlf, 1981).

333 Due to the greater dissimilarity of the  $\Delta M$  estimates between the endogenous and exogenous genes,  
334 and the slower decay rate of mutation bias, we decided to focus on our estimates of mutation bias  
335 to identify potential source lineages. The top ten lineages with the highest similarity in  $\Delta M$  to the  
336 exogenous genes were selected as potential candidates (Figure 2).

### 337 Phylogenetic Analysis

338 Using the dataset from Shen *et al.* (2018), we first identified 121 alignments for exogenous genes and  
339 further contained homologous genes for *E. gossypii*, and *L. thermotolerance*. We excluded all species  
340 from the alignments that do not belong to the Saccharomyctaceae clade. IQTree (Nguyen *et al.*, 2015)  
341 was used to identify the best fitting model for each gene and to estimate the individual gene trees. The  
342 distance between *L. kluyveri*, *E. gossypii*, and *L. thermotolerance* was calculated for each tree to identify  
343 genes for which exogenous genes are more closely related to *E. gossypii* or *L. thermotolerance*.

### 344 Synteny Comparison

345 We obtained complete genome sequences for all 10 candidate lineages (Table 2) from NCBI (on: 02-  
346 05-2017). Genomes were aligned and checked for synteny using SyMAP (4.2) with default settings  
347 (Soderlund *et al.*, 2006, 2011). We assess synteny as percentage coverage of the exogenous gene region.

### 348 Estimating Age of Introgression

We modeled the change in codon frequency over time using an exponential model for all two codon amino acids, and describing the change in codon  $c_1$  as

$$\frac{dc_1}{dt} = -\mu_{1,2}c_1 - \mu_{2,1}(1 - c_1) \quad (3)$$

where  $\mu_{i,j}$  is the rate at which codon  $i$  mutates to codon  $j$  and  $c_1$  is the frequency of the reference codon. Initial codon frequencies  $c_1(0)$  for each codon family were taken from our mutation parameter estimates for *E. gossypii* where  $c_1(0) = \exp[\Delta M_{gos}] / (1 + \exp[\Delta M_{gos}])$ . Our estimates of  $\Delta M_{endo}$  can be used to

calculate the steady state of equation 3 were  $\frac{dc_1}{dt} = 0$  to obtain the equality

$$\frac{\mu_{2,1}}{\mu_{1,2} + \mu_{2,1}} = \frac{1}{1 + \exp[\Delta M_{\text{endo}}]} \quad (4)$$

Solving for  $\mu_{1,2}$  gives us  $\mu_{1,2} = \Delta M_{\text{endo}} \exp[\mu_{2,1}]$  which allows us to rewrite and solve equation 3 as

$$c_1(t) = \frac{1 + \exp[-X](K - 1)}{1 + \Delta M_{\text{endo}}} \quad (5)$$

where  $X = (1 + \Delta M_{\text{endo}})\mu_{2,1}t$  and  $K = c_1(0)(1 + \Delta M_{\text{endo}})$ .

Equation 5 was solved with a mutation rate  $\mu_{2,1}$  of  $3.8 \times 10^{-10}$  per nucleotide per generation (Lang and Murray, 2008). Current codon frequencies for each codon family where taken from our estimates of  $\Delta M$  from the exogenous genes. Mathematica (11.3) (Wolfram Research Inc., 2017) was used to calculate the time  $t_{\text{intro}}$  it takes for the initial codon frequencies  $c_1(0)$  for each codon family to equal the current exogenous codon frequencies. The same equation was used to determine the time  $t_{\text{decay}}$  at which the signal of the exogenous cellular environment has decayed to within 1% of the endogenous environment.

## Estimating Selection against Codon Mismatch

In order to estimate the selection against codon mismatch, we had to make three key assumptions. First, we assumed that the current exogenous amino acid sequence of a gene is representative of its ancestral state and the replaced endogenous gene it replaced. Second, we assume that the currently observed cellular environment of *E. gossypii* reflects the cellular environment that the exogenous genes experienced before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy of selection between the cellular environments due to differences in either effective population size  $N_e$  or the selective cost of an ATP  $q$  of the source lineage and *L. kluyveri* can be expressed as a scaling constant and that protein synthesis rate  $\phi$  has not changed between the replaced endogenous and the introgressed exogenous genes. Using estimates for  $N_e = 1.36 \times 10^7$  (Wagner, 2005) for *Saccharomyces paradoxus* we scale our estimates of  $\Delta\eta$  which explicitly contains the effective population size  $N_e$  (Gilchrist *et al.*, 2015) and define  $\Delta\eta' = \frac{\Delta\eta}{N_e}$ .

We scale the difference in the efficacy of selection on codon usage between the donor lineage and *L. kluyveri* using a linear scaling factor  $\kappa$ . As  $\Delta\eta$  is defined as  $\Delta\eta = 2N_e q(\eta_i - \eta_j)$ , we cannot distinguish if  $\kappa$  is a scaling on protein synthesis rate  $\phi$ , effective population size  $N_e$ , or the selective cost of an ATP  $q$  (Gilchrist, 2007; Gilchrist *et al.*, 2015). We calculated the selection against each genes codon mismatch

assuming additive fitness effects as

$$s_g = \sum_{i=1}^{L_g} -\kappa \phi_g \Delta\eta'_i \quad (6)$$

where  $s_g$  is the overall strength of selection for translational efficiency on gene,  $g$  in the exogenous gene set,  $\kappa$  is a constant, scaling the efficacy of selection between the endogenous and exogenous cellular environments,  $L_g$  is length of the protein in codons,  $\phi_g$  is the estimated protein synthesis rate of the gene in the endogenous environment, and  $\Delta\eta'_i$  is the  $\Delta\eta'$  for the codon at position  $i$ . As stated previously, our  $\Delta\eta$  are relative to the mean of the codon family. We find that the selection against the introgressed genes is minimized at  $\kappa \sim 5$  (Figure S5b). Thus, we expect a five fold difference in the efficacy of selection between *L. kluyveri* and *E. gossypii*, due to differences in either protein synthesis rate  $\phi$ , effective population size  $N_e$ , and/or the selective cost of an ATP  $q$ . Therefore, we set  $\kappa = 1$  if we calculate the  $s_g$  for the endogenous and the current exogenous genes, and  $\kappa = 5$  for  $s_g$  for selection calculations at the time of introgression.

However, since we are unable to observe codon sequences of the replaced endogenous genes and for the exogenous genes at the time of introgression, instead of summing over the sequence, we calculate the expected codon count  $E[n_{g,i}]$  for codon  $i$  in gene  $g$  simply as the probability of observing codon  $i$  multiplied by the number of times the corresponding amino acids is observed in gene  $g$ , yielding:

$$E[n_{g,i}] = P(c_i | \Delta M, \Delta\eta, \phi) \times m_{a_i} \quad (7)$$

$$E[n_{g,i}] = \frac{\exp[-\Delta M_i - \Delta\eta_i \phi_g]}{\sum_j^C \exp[-\Delta M_j - \Delta\eta_j \phi_g]} \times m_{a_i} \quad (8)$$

where  $m_{a_i}$  is the number of occurrences of amino acid  $a$  that codon  $i$  codes for. Thus replacing the summation over the sequence length  $L_g$  in equ. (6) by a summation over the codon set  $C$  and calculating  $s_g$  as

$$s_g = \sum_{i=1}^C -\kappa \phi_g \Delta\eta'_i E[n_{g,i}] \quad (9)$$

We report the selection due to mismatched codon usage of the introgression as  $\Delta s_g = s_{\text{intro},g} - s_{\text{endo},g}$  where  $s_{\text{intro},g}$  is the selection against an introgressed gene  $g$  either at the time of the introgression or presently.

## Acknowledgments

This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ) and DEB-1355033 (BCO, MAG, and RZ) with additional support from The University of Tennessee Knoxville. CL received support as a Graduate Student Fellow at the National Institute for Mathematical and Biological Synthesis, an Institute sponsored by the National Science Foundation through NSF Award DBI-1300426, with additional support from UTK. The authors would like to thank Alexander Cope for helpful criticisms and suggestions for this work.

388 References

- 389 Baker, E. C., Wang, B., Bellora, N., *et al.* 2015. The genome sequence of *Saccharomyces eubayanus* and  
390 the domestication of lager-brewing yeasts. *Molecular Biology and Evolution*, 32(11): 2818–2831.

391 Beimforde, C., Feldberg, K., Nylander, S., *et al.* 2014. Estimating the phanerozoic history of the ascomy-  
392 cota lineages: combining fossil and molecular data. *Mol. Phylogenet. Evol.*, 78: 386–398.

393 Brion, C., Legrand, S., Peter, J., *et al.* 2017. Variation of the meiotic recombination landscape and  
394 properties over a broad evolutionary distance in yeasts. *PLoS Genetics*, 13(8): e1006917.

395 Bulmer, M. 1990. The selection-mutation-drift theory of synonymous codon usage. *Genetics*, 129: 897–  
396 907.

397 Cope, A. L., Hettich, R. L., and Gilchrist, M. A. 2018. Quantifying codon usage in signal peptides:  
398 Gene expression and amino acid usage explain apparent selection for inefficient codons. *Biochimica et  
399 Biophysica Acta (BBA) - Biomembranes*, 1860(12): 2479–2485.

400 dos Reis, M., Savva, R., and Wernisch, L. 2004. Solving the riddle of codon usage preferences: a test for  
401 translational selection. *Nucleic Acids Research*, 32(17): 5036–5044.

402 Friedrich, A., Reiser, C., Fischer, G., and Schacherer, J. 2015. Population genomics reveals chromosome-  
403 scale heterogeneous evolution in a protoploid yeast. *Molecular Biology and Evolution*, 32(1): 184 –  
404 192.

405 Gilchrist, M. A. 2007. Combining models of protein translation and population genetics to predict protein  
406 production rates from codon usage patterns. *Molecular Biology and Evolution*, 24(11): 2362–2372.

407 Gilchrist, M. A., Chen, W. C., Shah, P., Landerer, C. L., and Zaretzki, R. 2015. Estimating gene  
408 expression and codon-specific translational efficiencies, mutation biases, and selection coefficients from  
409 genomic data alone. *Genome Biology and Evolution*, 7: 1559–1579.

- 410 Gouy, M. and Gautier, C. 1982. Codon usage in bacteria: correlation with gene expressivity. *Nucleic*  
411 *Acids Research*, 10: 7055–7074.
- 412 Gronau, Q. F., Sarafoglou, A., Matzke, D., *et al.* 2017. A tutorial on bridge sampling. *Journal of*  
413 *Mathematical Psychology*, 81: 80–97.
- 414 Ikemura, T. 1985. Codon usage and tRNA content in unicellular and multicellular organisms. *Molecular*  
415 *Biology and Evolution*, 2: 13–34.
- 416 Kensche, P. R., Oti, M., Dutilh, B. E., and Huynen, M. A. 2008. Conservation of divergent transcription  
417 in fungi. *Trends Genet.*, 5(24): 207–211.
- 418 Landerer, C., Cope, A., Zaretzki, R., and Gilchrist, M. A. 2018. AnaCoDa: analyzing codon data with  
419 bayesian mixture models. *Bioinformatics*, 34(14): 2496–2498.
- 420 Lang, G. I. and Murray, A. W. 2008. Estimating the per-base-pair mutation rate in the yeast *Saccha-*  
421 *romyces cerevisiae*. *Genetics*, 178(1): 67 – 82.
- 422 Lawrence, J. G. and Ochman, H. 1997. Amelioration of bacterial genomes: Rates of change and exchange.  
423 *Journal of Molecular Biology*, 44: 383–397.
- 424 Legendre, P. 2018. *lmodel2: Model II Regression*. R package version 1.7-3.
- 425 Marcet-Houben, M. and Gabaldón, T. 2015. Beyond the whole-genome duplication: Phylogenetic ev-  
426 idence for an ancient interspecies hybridization in the baker’s yeast lineage. *PLoS Biology*, 13(8):  
427 e1002220.
- 428 Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., and Danchin, A. 1991. Evidence for horizontal gene  
429 transfer in *Escherichia coli* speciation. *Journal of Molecular Biology*, 222(4): 851–856.
- 430 Nguyen, L. T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. 2015. Iq-tree: A fast and effective  
431 stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution*,  
432 32(1): 268–274.
- 433 Payen, C., Fischer, G., Marck, C., *et al.* 2009. Unusual composition of a yeast chromosome arm is  
434 associated with its delayed replication. *Genome Research*, 19(10): 1710–1721.
- 435 R Core Team 2013. *R: A Language and Environment for Statistical Computing*. R Foundation for  
436 Statistical Computing, Vienna, Austria.
- 437 Ruderfer, D. M., Pratt, S. C., Seidl, H. S., and Kruglyak, L. 2006. Population genomic analysis of  
438 outcrossing and recombination in yeast. *Nature Genetics*, 38(9): 1077–1081.

- 439 Salichos, L. and Rokas, A. 2013. Inferring ancient divergences requires genes with strong phylogenetic  
440 signals. *Nature*, 497: 327–331.
- 441 Sella, G. and Hirsh, A. E. 2005. The application of statistical physics to evolutionary biology. *Proceedings*  
442 *of the National Academy of Sciences of the United States of America*, 102: 9541–9546.
- 443 Shah, P. and Gilchrist, M. A. 2011. Explaining complex codon usage patterns with selection for trans-  
444 lational efficiency, mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences*  
445 *U.S.A.*, 108(25): 10231–10236.
- 446 Sharp, P. M. and Li, W. H. 1987. The codon adaptation index - a measure of directional synonymous  
447 codon usage bias, and its potential applications. *Nucleic Acids Research*, 15: 1281–1295.
- 448 Shen, X. X., Opulente, D. A., Kominek, J., *et al.* 2018. Tempo and mode of genome evolution in the  
449 budding yeast subphylum. *Cell*, 175(6): 1533–1545.e20.
- 450 Soderlund, C., Nelson, W., Shoemaker, A., and Paterson, A. 2006. Symap A system for discovering and  
451 viewing syntenic regions of fpc maps. *Genome Research*, 16: 1159 – 1168.
- 452 Soderlund, C., Bomhoff, M., and Nelson, W. 2011. Symap v3.4: a turnkey synteny system with applica-  
453 tion to plant genomes. *Nucleic Acids Research*, 39(10): e68.
- 454 Sokal, R. R. and Rohlf, F. J. 1981. *Biometry - The principles and practice of statistics in biological*,  
455 pages 547–555. W. H. Freeman.
- 456 Tsai, I. J., Bensasson, D., Burt, A., and Koufopanou, V. 2008. Population genomics of the wild yeast  
457 *Saccharomyces paradoxus*: quantifying the life cycle. *Proc Natl Acad Sci U.S.A.*, 105: 4957–4962.
- 458 Tsankov, A. M., Thompson, D. A., Socha, A., Regev, A., and Rando, O. J. 2010. The role of nucleosome  
459 positioning in the evolution of gene regulation. *PLoS Biol*, 8(7): e1000414.
- 460 Vakirlis, N., Sarilar, V., Drillon, G., *et al.* 2016. Reconstruction of ancestral chromosome architecture  
461 and gene repertoire reveals principles of genome evolution in a model yeast genus. *Genome research*,  
462 26(7): 918–32.
- 463 Wagner, A. 2005. Energy constraints on the evolution of gene expression. *Molecular Biology and Evolu-*  
464 *tion*, 22: 1365–1374.
- 465 Wallace, E. W., Aioldi, E. M., and Drummond, D. A. 2013. Estimating selection on synonymous codon  
466 usage from noisy experimental data. *Molecular Biology and Evolution*, 30: 1438–1453.
- 467 Wolfram Research Inc. 2017. *Mathematica 11*.

Wright, F. 1990. The ‘effective number of codons’ used in a gene. *Genet.*, 87: 23–29.

469

## Supplementary Material

470

Supporting Materials for *Decomposing Mutation and Selection to Identify Mismatched Codon Usage* by  
Landerer *et al.*

471

Table S1: Synonymous mutation codon preference based on our estimates of  $\Delta M$ . Shown are the most likely codon in low expression genes for each amino acid in: *E. gossypii*, in the endogenous and exogenous genes of *L. kluyveri*, and in the combined *L. kluyveri* genome without accounting for the two cellular environments.

Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	Combined
Ala A	GCG	GCA	GCG	GCG
Cys C	TGC	TGT	TGC	TGC
Asp D	GAC	GAT	GAC	GAC
Glu E	GAG	GAA	GAG	GAG
Phe F	TTC	TTT	TTT	TTT
Gly G	GGC	GGT	GGC	GGC
His H	CAC	CAT	CAC	CAC
Ile I	ATC	ATT	ATC	ATA
Lys K	AAG	AAA	AAG	AAA
Leu L	CTG	TTG	CTG	CTG
Asn N	AAC	AAT	AAC	AAT
Pro P	CCG	CCA	CCG	CCG
Gln Q	CAG	CAA	CAG	CAG
Arg R	CGC	AGA	AGG	CGG
Ser <sub>4</sub> S	TCG	TCT	TCG	TCG
Thr T	ACG	ACA	ACG	ACG
Val V	GTG	GTT	GTG	GTG
Tyr Y	TAC	TAT	TAC	TAC
Ser <sub>2</sub> Z	AGC	AGT	AGC	AGC

Table S2: Synonymous selection codon preference based on our estimates of  $\Delta\eta$ . Shown are the most likely codon in high expression genes for each amino acid in: *E. gossypii*, in the endogenous and exogenous genes of *L. kluyveri*, and in the combined *L. kluyveri* genome without accounting for the two cellular environments.

Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	Combined
Ala A	GCT	GCT	GCT	GCT
Cys C	TGT	TGT	TGT	TGT
Asp D	GAT	GAC	GAT	GAT
Glu E	GAA	GAA	GAA	GAA
Phe F	TTT	TTC	TTC	TTC
Gly G	GGA	GGT	GGT	GGT
His H	CAT	CAC	CAT	CAT
Ile I	ATA	ATC	ATT	ATT
Lys K	AAA	AAG	AAA	AAG
Leu L	TTA	TTG	TTG	TTG
Asn N	AAT	AAC	AAT	AAC
Pro P	CCA	CCA	CCT	CCA
Gln Q	CAA	CAA	CAA	CAA
Arg R	AGA	AGA	AGA	AGA
Ser <sub>4</sub> S	TCA	TCC	TCT	TCT
Thr T	ACT	ACC	ACT	ACT
Val V	GTT	GTC	GTT	GTT
Tyr Y	TAT	TAC	TAT	TAC
Ser <sub>2</sub> Z	AGT	AGT	AGT	AGT



## Endogenous and Combined Codon Usage

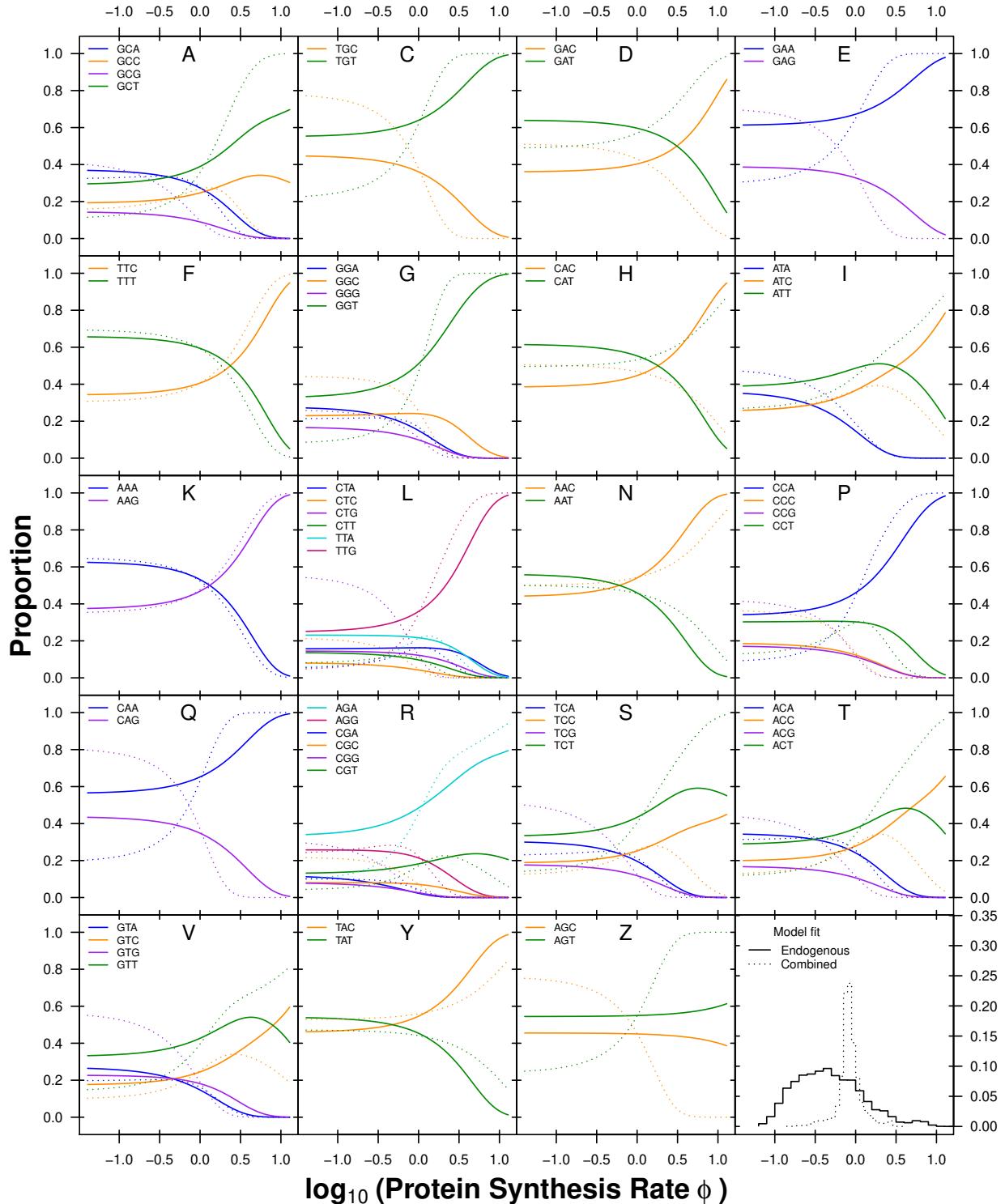


Figure S1: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dotted line indicates the combined codon usage.

## Exogenous and Combined Codon Usage

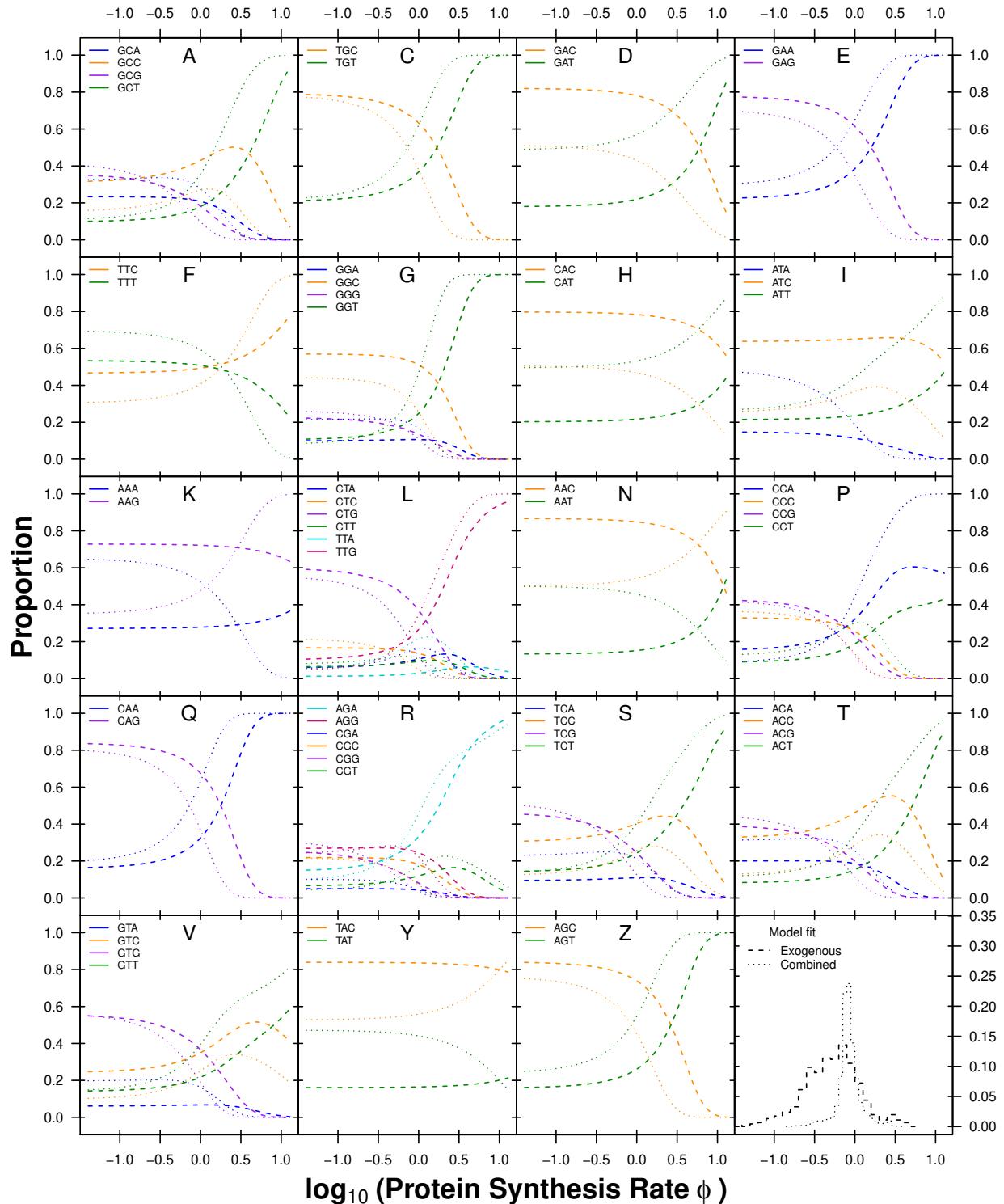


Figure S2: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. dashed line indicates the exogenous codon usage, dotted line indicates the combined codon usage.

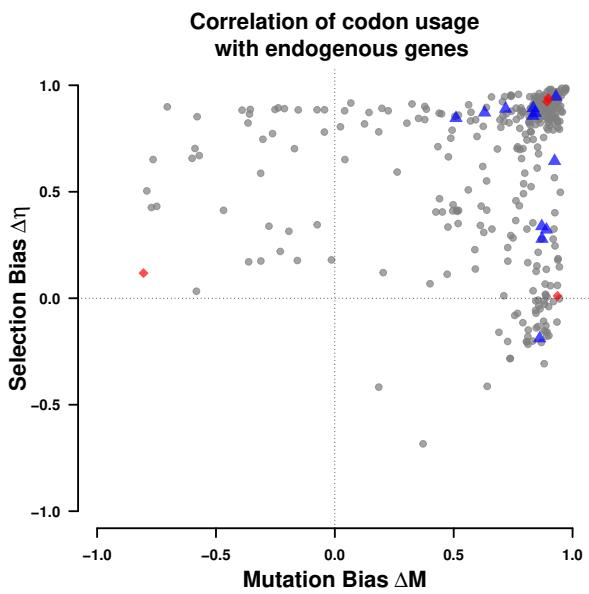


Figure S3: Correlation coefficients of  $\Delta M$  and  $\Delta\eta$  of the endogenous genes with 332 examined budding yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta\eta$  of the lineages with the exogenous parameter estimates. Blue triangles indicate the Lachancea and red diamonds indicate Eremothecium lineages. All regressions were performed using a type II regression assuming noise in the dependent and independent variable (Sokal and Rohlf, 1981).

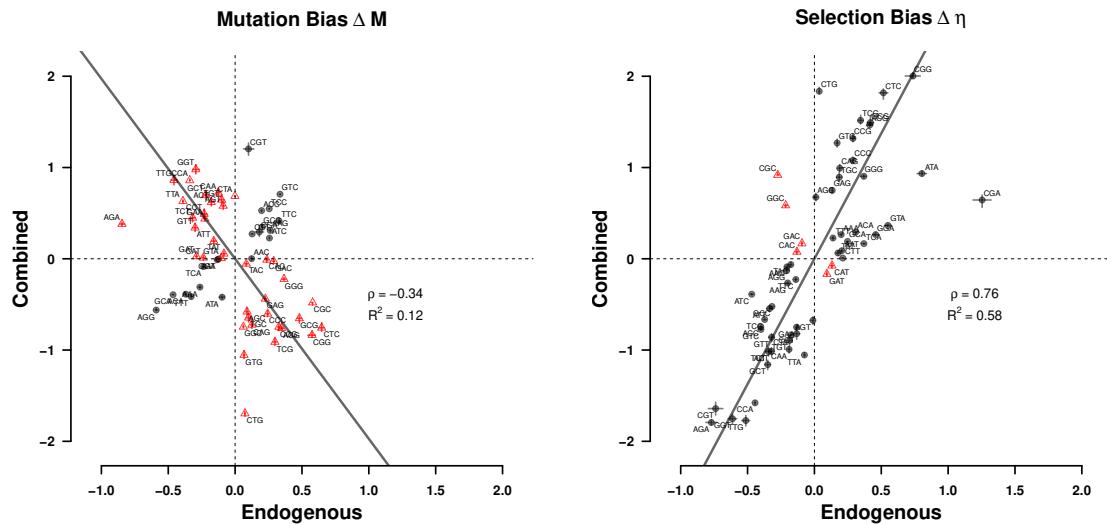


Figure S4: Comparison of (a) mutation bias  $\Delta M$  and (b) selection bias  $\Delta \eta$  parameters for endogenous genes and combined gene sets. Estimates are relative to the mean for each codon family. Black dots indicate  $\Delta M$  or  $\Delta \eta$  parameters with the same sign for the endogenous and exogenous genes, red dots indicate parameters with different signs. Black line indicates type II regression line assuming noise in the dependent and independent variable (Sokal and Rohlf, 1981). Dashed lines mark quadrants.

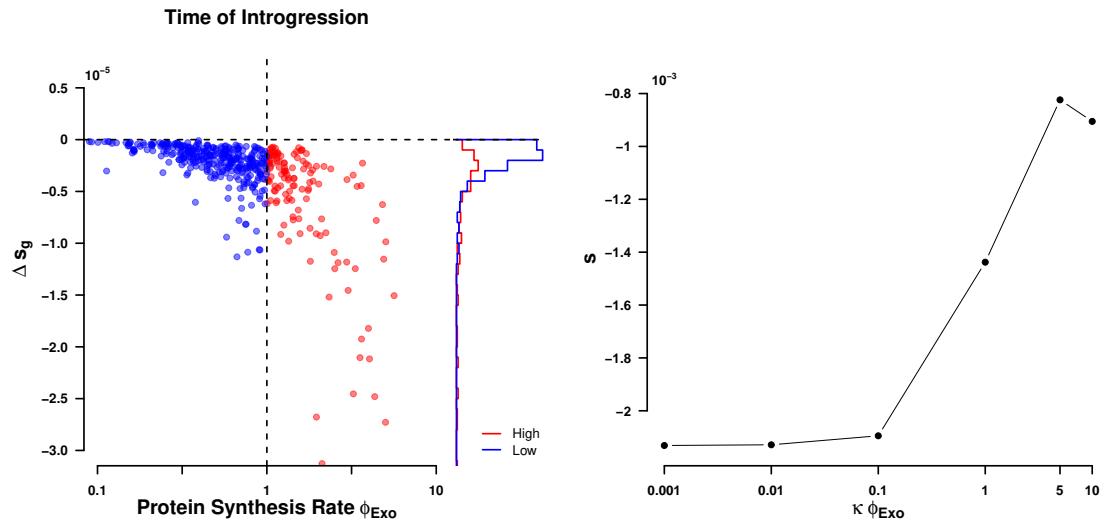


Figure S5: Selection against mismatched codon usage (left) without scaling of  $\phi$  per gene. Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates neutrality. (Right) Change of total selection against mismatched codon usage with scaling term  $\kappa$  between *E. gossypii* and *L. kluyveri*

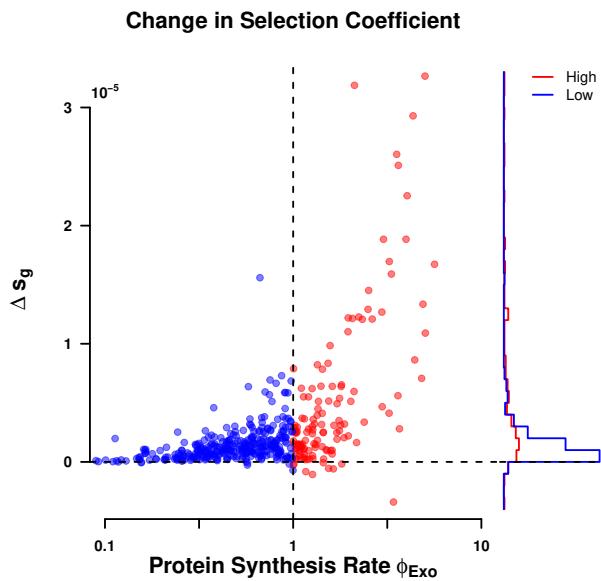


Figure S6: Total amount of adaptation estimated to have occurred between time of introgression and currently observed per gene. Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates no change in selection against mismatched codon usage.

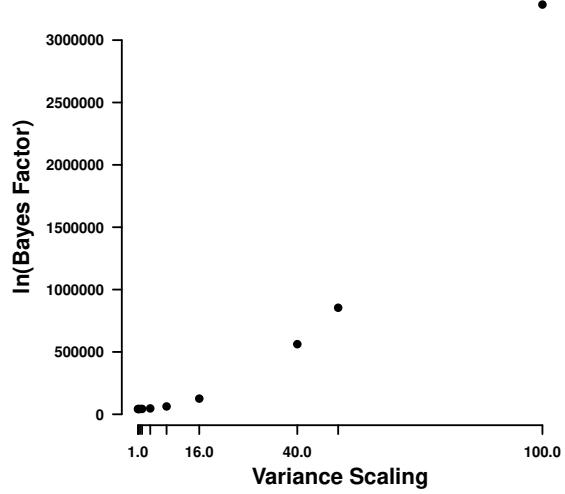


Figure S7: Influence of the variance scaling of the importance distribution on the estimated Bayes factor.