

<sup>1</sup> Unlocking a signal of introgression from codons in *Lachancea*  
<sup>2</sup> *kluveryi* 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> May 23, 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 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 genetic load. 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 this  
31 mutation bias is an omnipresent evolutionary force, its impact can be obscured or amplified by selection.  
32 The signature of selection on codon usage is also largely determined by an organism's cellular environ-  
33 ment, such as its tRNA species, their copy number, and their post-transcriptional modifications. The  
34 strength of selection on the codon usage of an individual gene is largely determined by its expression and  
35 synthesis rate which, in turn, is largely determined by the organism's external environment. In general,  
36 the strength of selection on codon usage increases with its expression level (Gouy and Gautier, 1982;  
37 Ikemura, 1985; Bulmer, 1990), specifically its protein synthesis rate (Gilchrist, 2007). Thus as protein  
38 synthesis increases, codon usage shifts from a process dominated by mutation to a process dominated  
39 by selection. The overall efficacy of selection on codon usage is a function of the organism's effective  
40 population size  $N_e$  which, in turn, is largely determined by its external environment. ROC SEMPPR  
41 allows us disentangle the evolutionary forces responsible for the patterns of codon usage bias (CUB)  
42 encoded in an species' genome, by explicitly modeling the combined forces of mutation, selection, and  
43 drift (Gilchrist, 2007; Shah and Gilchrist, 2011; Wallace *et al.*, 2013; Gilchrist *et al.*, 2015). In turn,  
44 these forces should provide biologically meaningful information about the lineage's historical cellular and  
45 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  
50 than endogenous genes to their new cellular environment, we expect a greater genetic load of transferred  
51 genes if donor and recipient environment differ greatly in their selection bias, making such transfers less  
52 likely. More practically, if differences in codon usage of transferred genes are unaccounted for, they may  
53 distort the interpretation of codon usage patterns. Such distortion could lead to the wrong inference of  
54 codon preference for an amino acid (Shah and Gilchrist, 2011; Gilchrist *et al.*, 2015)., underestimate the  
55 variation in protein synthesis rate, or influence mutation estimates when analyzing a genome.

56 To illustrate these ideas, we analyze the CUB of the genome of *Lachancea kluyveri*, which is sister to all  
57 other Lachancea. The Lachancea clade diverged from the *Saccharomyces* clade, prior to its whole genome  
58 duplication ~ 100 Mya ago (Marcet-Houben and Gabaldón, 2015; Beimforde *et al.*, 2014). Since that  
59 time, *L. kluyveri* has experienced a large introgression of exogenous genes found in all of its populations

(Friedrich *et al.*, 2015), but in no other known *Lachancea* species (Vakirlis *et al.*, 2016). The introgression replaced the left arm of the C chromosome and displays a 13% higher GC content than the endogenous *L. kluyveri* genome (Payen *et al.*, 2009; Friedrich *et al.*, 2015). The origin the introgression is currently unknown, but previous studies suggest that the source is likely a currently unknown or potentially extinct *Lachancea* lineage (Payen *et al.*, 2009; Friedrich *et al.*, 2015; Vakirlis *et al.*, 2016; Brion *et al.*, 2017). These characteristics make *L. kluyveri* an ideal model to study the effects of an introgressed cellular environment and the resulting mismatch in codon usage.

Using ROC SEMPPR, a Bayesian population genetics model based on a mechanistic description of ribosome movement along an mRNA, allows us to quantify the cellular environment in which genes have evolved by separately estimating the effects of mutation bias and selection bias on codon usage. ROC SEMPPR's resulting predictions of protein synthesis rates have been shown to be on par with laboratory measurements (Shah and Gilchrist, 2011; Gilchrist *et al.*, 2015). In contrast to often used heuristic approaches to study codon usage (Sharp and Li, 1987; Wright, 1990; dos Reis *et al.*, 2004), ROC SEMPPR explicitly incorporates and distinguishes between mutation and selection effects on codon usage and properly weights by amino acid usage (Cope *et al.*, 2018). We use ROC SEMPPR to independently describe two cellular environments reflected in the *L. kluyveri* genome; the signature of the current environment in the endogenous genes and the decaying signature of the exogenous environment in the introgressed genes. Our results indicate that the difference in GC content between endogenous and exogenous genes is mostly due to the differences in mutation bias of their ancestral environments. Accounting for these different signatures of mutation bias and selection bias of the endogenous and exogenous sets of genes substantially improves our ability to predict present day protein synthesis rates. These endogenous and exogenous gene set specific estimates of mutation bias and selection bias, in turn, allow us to address more refined questions of biological importance. For example, it allows us to provide an alternative hypothesis about the origin of the introgression and identify *E. gossypii* as the nearest sampled relative of the source of the introgressed genes out of the 332 budding yeast lineages with sequenced genomes (Shen *et al.*, 2018). While this hypothesis is in contrast previous work (Payen *et al.*, 2009; Friedrich *et al.*, 2015; Vakirlis *et al.*, 2016; Brion *et al.*, 2017), we find support for it in gene trees and synteny. We also estimate the age of the introgression to be on the order of 0.2 - 1.7 Mya, estimate the genetic load of these genes, both at the time of introgression and now, and predict a detectible signature of CUB to 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  $\Delta M$  selection bias  $\Delta\eta$  between 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). For example, the endogenous genes show a mutational bias for A and T ending codons in  $\sim 95\%$  of the codon families. In contrast, the exogenous genes display an equally consistent mutational bias towards C and G ending codons (Table S1). As a result, only the two codon amino acid Phenylalanine (Phe, F) shares the same rank order across the endogenous and exogenous  $\Delta M$  estimates.

In contrast, our estimates of  $\Delta\eta$  for the endogenous and exogenous genes were positively correlated ( $\rho = 0.69$ ) and showing similarity of  $\sim 53\%$  between the two selection environments (Figure 2). ROC SEMPPR constraints  $E[\phi] = 1$ , allowing us to interpret  $\Delta\eta$  as selection on codon usage of the average gene with  $\phi = 1$  and gives us the ability to compare the efficacy of selection  $sN_e$  across genomes. We

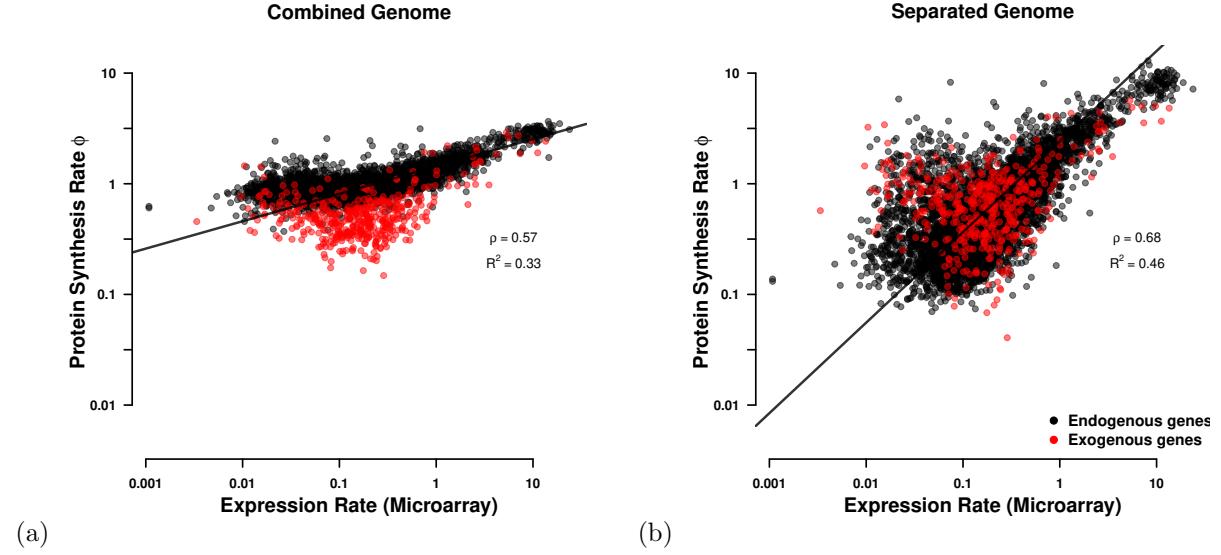


Figure 1: Comparison of predicted protein synthesis rate  $\phi$  to microarray data from Tsankov *et al.* (2010) for (a) the combined genome and (b) the separated endogenous and exogenous genes. Endogenous genes are displayed in black and exogenous genes in red. Black line indicates type II regression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

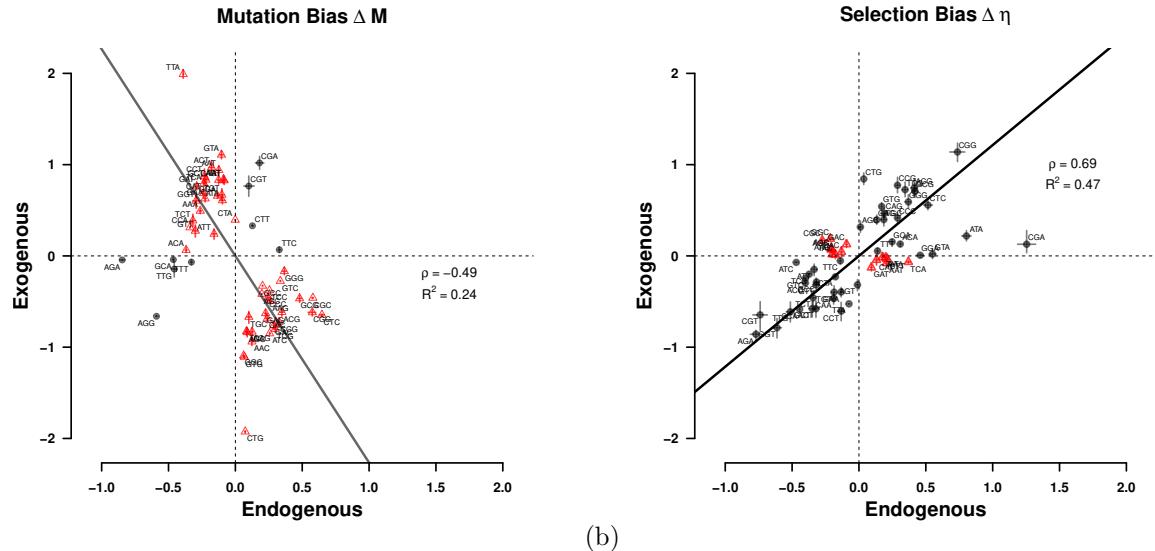


Figure 2: Comparison of (a) mutation bias  $\Delta M$  and (b) selection bias  $\Delta \eta$  parameters for endogenous and exogenous genes. Estimates are relative to the mean for each codon family. Black dots indicate  $\Delta M$  or  $\Delta \eta$  parameters with the same sign for the endogenous and 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.

115 find that the strength of selection within each codon family differs between sets of genes. Overall, the  
116 endogenous genes only show a selection preference for C and G ending codons in  $\sim$  58% of the codon  
117 families. In contrast, the exogenous genes display a strong preference for A and T ending codons in  
118  $\sim$  89% of the codon families.

119 The difference in codon usage between endogenous and exogenous genes is striking as some amino  
120 acids have opposite codon preferences. As a result, our estimates of the optimal codon differ in nine  
121 cases between endogenous and exogenous genes (Figure 3, Table S2). For example, Asparagine (Asn, N)  
122 shows strong preference for AAC in highly expressed endogenous genes but the same codon is depleted  
123 in highly expressed exogenous genes. In addition, fits to the complete *L. kluyveri* genome reveal that  
124 the relatively small exogenous gene set ( $\sim$  10% of genes) has a disproportional effect on the model fit  
125 (Figure S1, S2). We find that the combined genome shows the same codon preference in highly expressed  
126 genes as the exogenous gene set for Aspartic acid (Asp, D), despite the gene set only representing  $\sim$  10%  
127 of the genes. In the nine cases the endogenous and exogenous genes show difference in the selectively  
128 favored codon. In five of these cases (Asp, D; His, H; Lys, K; Asn, N; and Pro, P) the endogenous genes  
129 favor the codon with the most abundant tRNA. For the remaining four cases (Ile, I; Ser, S; Thr, T; and  
130 Val, V), there are no tRNA genes for the wobble free cognate codon encoded in the *L. kluyveri* genome.  
131 We find that the complete *L. kluyveri* genome fit shares mutational preference with the exogenous genes  
132 in  $\sim$  78% of the 19 codon families that are discordant between the endogenous and exogenous genes.  
133 In two cases, Isoleucine (Ile, I) and Arginine (Arg, R), the strong dissimilarity in mutation preference  
134 results in an estimated codon preference in the complete *L. kluyveri* genome that differs from both the  
135 endogenous, and the exogenous genes.

136 The effect of the small exogenous gene set on the fit to the complete *L. kluyveri* genome is smaller  
137 in our estimates of selection bias  $\Delta\eta$  than  $\Delta M$ , but still large. We find that the complete *L. kluyveri*  
138 genome is estimated to share the selection preference with the exogenous genes in  $\sim$  60% of codon  
139 families that show dissimilarity between endogenous and exogenous genes. These results clearly show  
140 that it is important to recognize the difference in endogenous and exogenous genes and treat these genes  
141 as separate sets to avoid the inference of incorrect synonymous codon preferences and better predict  
142 protein synthesis.

## 143 Determining Source of Exogenous Genes

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

### Endogenous and Exogenous Codon Usage

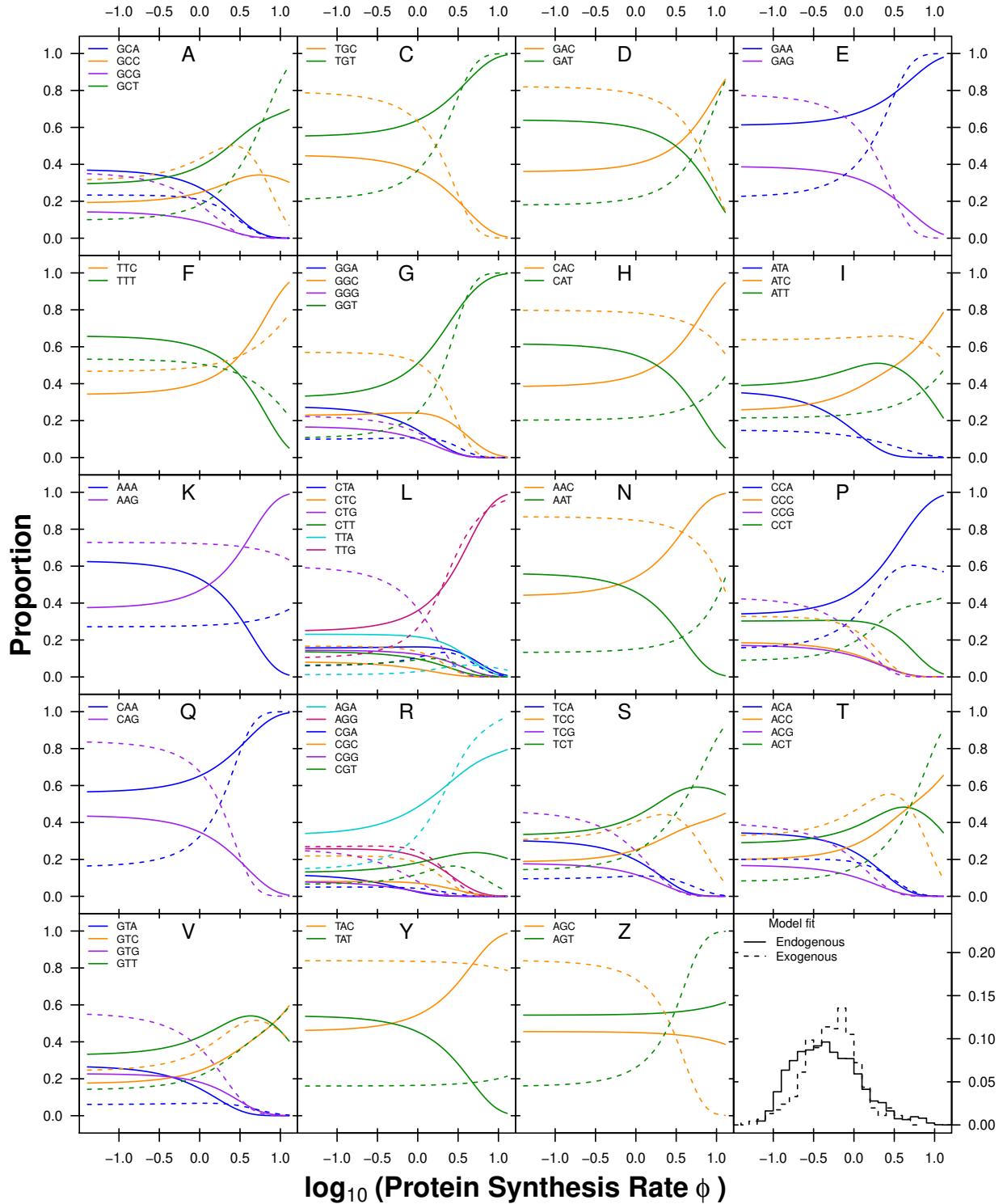


Figure 3: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dashed 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 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 non-standard codon table 12

as the endogenous and exogenous genes show greater dissimilarity in mutation bias than in selection bias.

Two of the ten candidate lineages utilise the alternative yeast nuclear code (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 neighbour 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 ~ 17% of all examined yeast show a positive correlation in both,  $\Delta M$  and  $\Delta \eta$  with the exogenous genes, whereas most lineages (~ 83%) show a negative correlation for  $\Delta M$  but only 21 % a negative correlation for  $\Delta \eta$ . This indicates that information on mutation bias provides more information about a potential origin of the exogenous genes.

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. However, this result is in contrast to previous studies suggesting that the introgressioned region originated from within the Lachancea clade (Payen *et al.*, 2009; Vakirlis *et al.*, 2016). Therefore, 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 ~ 60% of exogenous genes (73/121) are more closely related to *E. gossypii* than to other Lachancea.

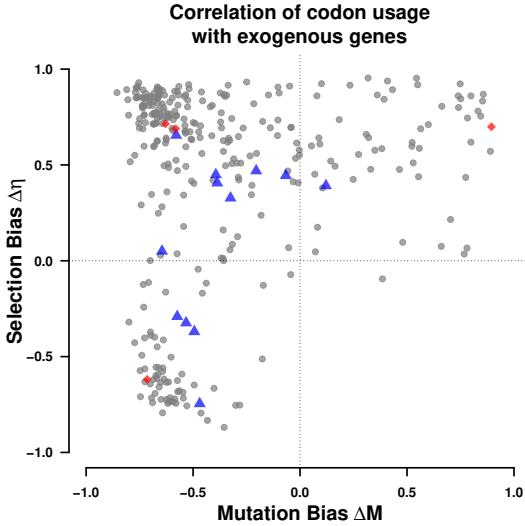


Figure 4: Correlation coefficients of  $\Delta M$  and  $\Delta \eta$  of the exogenous genes with 332 examined budding yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta \eta$  of the lineages with the exogenous parameter estimates. Blue triangles indicate the *Lachancea* and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

## 169 Estimating Introgression Age

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

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

183      **Estimating Genetic Load of Codon Mismatch of the Exogenous Genes**

184      We define genetic load as the difference between the fitness of an expected, replaced endogenous gene  
185      and the exogenous gene,  $sN_e \propto \phi\Delta\eta$  due to the mismatch in codon usage parameters (See Methods for  
186      details). As the introgression occurred before the diversification of *L. kluyveri* and has fixed throughout  
187      all populations (Friedrich *et al.*, 2015), we can not observe the original endogenous sequences that have  
188      been replaced by the introgression. Using our estimates of  $\Delta M$  and  $\Delta\eta$  from the endogenous genes  
189      and assuming the current exogenous amino acid composition of genes is representative of the replaced  
190      endogenous genes, we estimate the genetic load of the exogenous genes at the time of introgression (Figure  
191      5a) and currently (Figure 5b). Estimates of selection bias for the exogenous genes show that, while well  
192      correlated with the endogenous genes, only nine amino acids share the same optimal codon. Exogenous  
193      genes are, therefore, expected to represent a significant reduction in fitness, or genetic load for *L. kluyveri*  
194      due to this mismatch in codon usage. We find that the expected genetic load  $E[s_g]$  due to mismatched  
195      codon usage was -0.0008 at the time of the introgression and still represents a genetic load of -0.0003  
196      today.

197      In order to account for differences in the efficacy of selection on codon usage either due to the cost of  
198      pausing or differences in the effective population size or the value of an ATP between the donor lineage  
199      and *L. kluyveri* we added a linear scaling factor  $\kappa$  to scale our estimates of  $\Delta\eta$  between the donor lineage  
200      and *L. kluyveri* (See Methods for details). We predict that a small number of low expression genes  
201      ( $\phi < 1$ ) were weakly exapted at the time of the introgression (Figure 5a). High expression genes ( $\phi > 1$ )  
202      are predicted to have carried the largest genetic load in the novel cellular environment. These highly  
203      expressed genes are inferred to have the greatest degree of adaptation since the time of the introgression  
204      to the *L. kluyveri* cellular environment (Figures 5a & S6).

205      **Discussion**

206      In order to study the evolutionary effects of the large scale introgression of C-left, we used ROC SEMPPR,  
207      a mechanistic model of ribosome movement along an mRNA. Our parameter estimates indicate that the  
208      *L. kluyveri* genome contains distinct signatures of mutation and selection bias from both an endogenous  
209      and exogenous cellular environment. By fitting ROC SEMPPR separately to *L. kluyveri*'s endogenous  
210      and exogenous sets of genes we generate a quantitative description of their signatures of mutation bias  
211      and natural selection for efficient protein translation. While previous work by Payen *et al.* (2009) showed  
212      a preference for GC rich codons in the exogenous genes our results provide more nuanced insights. Our  
213      results indicate that the difference in GC content between endogenous and exogenous genes is mostly due  
214      to differences in mutation bias as 95% of exogenous codon families show a strong mutation bias towards

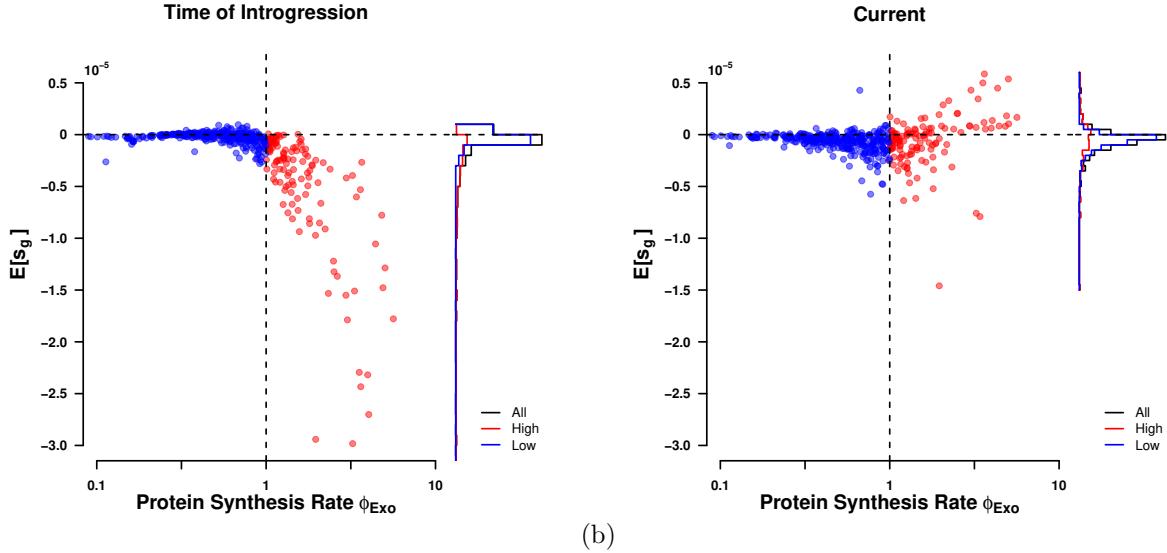


Figure 5: Genetic load  $s = \Delta\eta\phi$  (a) at the time of introgression ( $\kappa = 5$ ), and (b) currently ( $\kappa = 1$ ). Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates a genetic load of 0.

GC ending codons. We also show that the strength and rank order of selection within a codon family differ between endogenous and exogenous cellular environments. Even though the exogenous genes make up only  $\sim 10\%$  of the *L. kluyveri* genome, when we fail to recognize these differences our estimates of  $\Delta M$  and  $\Delta\eta$  deviate substantial from their actual values (Figure S4). This sensitivity of our parameters to a second cellular environment is surprising and highlights the importance of recognizing different cellular environments reflected within a genome. Furthermore, our results indicate that we can attribute the increased GC content in the exogenous genes mostly to differences in mutation bias favoring GC ending codons rather than selection.

The separation of the endogenous and exogenous genes improves our estimates of protein synthesis rate  $\phi$  by 42% relative to the full genome estimate ( $R^2 = 0.46$  vs. 0.32, respectively). Furthermore, failing to separately analyze the endogenous and exogenous genes results in a reduced intergenic variation in ROC SEMPPR estimates of  $\phi$  (compare Figure 1a & b). This behavior is due, in part, to constraining  $E[\phi] = 1$ , which is necessary in order to compare the efficacy of selection across genomes. Thus, extremely small variances in the  $\phi$  values estimated by ROC SEMPPR could indicate that a genome contains the signature of multiple cellular environments.

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. However, due to the greater dissimilarity of the  $\Delta M$  estimates between the endogenous and exogenous genes, and the slower decay

rate of mutation bias, identified potential source lineages using  $\Delta M$  we decided to focus on our estimates of mutation bias to identify potential source lineages. The top ten lineages with the highest similarity in  $\Delta M$  to the exogenous genes were selected as potential candidates (Figure 2). In terms of gene order, we found that synteny with the exogenous genes is limited to the Saccharomycetaceae clade, which excludes all of the potential source lineages identified using codon usage but *E. gossypii* (Table 2). Previous work indicated that the donor lineage of the exogenous genes has to be a, potentially unknown, Lachancea (Payen *et al.*, 2009; Vakirlis *et al.*, 2016). These previous results, however, are based on species rather than genes trees. This becomes important since, as we highlight here, relatively small sets of genes with a differing signal can significantly bias the outcome of a study. The same holds true for phylogenetic inferences (Salichos and Rokas, 2013), and as we showed the signal of the original endogenous cellular environment that shaped CUB is at different stages of decay in high and low expression genes (Figure S6).

Considering the similarity in selection bias (Figure 2b) and our calculation of the genetic load of the exogenous genes (Figure 5b), both of which are free of any assumption about the origin of the exogenous genes, a species tree estimated from the exogenous genes may be biased towards the Lachancea clade. While we were only able to consider a subset of exogenous genes for our phylogenetic analysis, these genes provide further evidence that the exogenous genes could originate from a lineage that does not belong to the Lachancea clade. In addition, the synteny coverage extends along most of the exogenous regions with the exception of the 3' and 5' ends of the exogenous region, and of the 13 candidate lineages, *E. gossypii* is the only lineage with a GC content > 50%, making it most similar to the exogenous genes. Thus, only the *E. gossypii* genome displays strong correlations in  $\Delta M$  and  $\Delta \eta$ , synteny, and similar GC content with the exogenous genes. In summary, our work does not dispute an unknown Lachancea as possible origin, but provides an alternative hypothesis based on the codon usage of the exogenous genes and synteny.

Interestingly, our results also indicate that codon usage does not necessarily correlate with phylogenetic distance. As codon usage is influenced by an organism's environment and it is common to find lineages evolving in very different environments, this is not surprising. Nevertheless, it shows that the approach presented here could be used to study the evolution of codon usage in more detail.

Our results showed that the endogenous and exogenous genes show similarity in  $\Delta \eta$  but not in  $\Delta M$ . This can be caused by either similarity in  $\Delta \eta$  between the source lineage of the exogenous genes and *L. kluyveri* or by a fast decay of the original signature of  $\Delta \eta$ . In our estimation of a source lineage, we explicitly assumed the first case. While our estimate is still supported by GC content and synteny, we cannot rule out that this biased our inference of the source lineage.

Assuming *E. gossypii* as potential source lineage of the introgressed region, we illustrated how info-

mation on codon usage can be used to infer the time since the introgression occurred using our estimates of mutation bias  $\Delta M$ . Our  $\Delta M$  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 time than a previous minimum estimate by Friedrich *et al.* (2015) of  $5.6 \times 10^7$  generations. Our estimate assumes that the current *E. gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the time of the introgression. If the ancestral mutation environments were more similar (dissimilar) at the time of the introgression than now our result is an overestimate (underestimate).

In order to estimate the introgression's genetic load due to codon mismatch, we had to make three key assumptions: 1) at the time of introgression the amino acid sequences of the endogenous genes and exogenous genes were highly similar, 2) the current *L. kluyveri* cellular environment is reflective of the cellular environment at the time of the introgression, and 3) the *E. gossypii* cellular environment reflects its ancestral environment at the time of the introgression. In general due to their very nature, low expression genes contribute little to the genetic load. Indeed,  $\sim 30\%$  of low expression exogenous genes ( $\phi < 1$ ) appeared to be exapted at the time of the introgression. This exaptation is due to the mutation bias in the endogenous genes matching the selection bias in the exogenous genes for GC ending codons. In contrast, highly expressed genes are predicted to have imposed a large genetic load. Many of these genes appear to still represent a significant genetic load. Overall, our estimates of codon mismatch genetic load, therefore, suggest strong selection against the introgression.

It is hard to contextualize the probability of this introgression being fixed as we are not aware of any estimates of the frequency at which such large scale introgressions of genes occur. A related example of a large scale merger of genomic material can be found in *S. pastorianus*, which is currently believed to be a hybrid of *S. cerevisiae* and *S. eubayanus* lineages, (Baker *et al.*, 2015). Unlike with *L. kluyveri* and *E. gossypii*, the progenitor lineages of *S. pastorianus* have similar codon usage parameters. The correlation between  $\Delta M$  and  $\Delta\eta$  for these two lineages are  $\rho = 0.83$  and  $0.98$  (data not shown). These similarities in  $\Delta M$  and  $\Delta\eta$  parameters suggest that the genetic load for *S. pastorianus* due to codon usage mismatch is small relative to the exogenous genes considered here. The large genetic load of the exogenous genes due to codon mismatch at the time of the introgression would seem to indicate that the fixation of the introgression was either a fluke event or the codon mismatch genetic load was countered by one or more highly advantageous loci within the introgression.

Under the first scenario, our best estimate of the selection coefficient against the introgression based on expected codon mismatch at that time is  $s = -0.0008$  and an effective population size  $N_e$  on the order of  $10^7$  (Wagner, 2005) yields an approximate fixation probability of  $(1-\exp[-s])/(1-\exp[-2sN_e]) \approx 10^{-6952}$  (Sella and Hirsh, 2005). The astronomically small fixation probability indicates that there was virtually

no chance of fixation and was likely not a fluke. Even though *L. kluyveri* diverged from the rest of the Lachancea clade around 85 Mya (Kensche *et al.*, 2008; Marcet-Houben and Gabaldón, 2015), if we assume 1 to 8 generations/day, which implies  $10^{10}$  to  $10^{11}$  generations since the time of divergence, one round of meiosis for every 1000 rounds of mitosis of which only one in 100 meiosis events lead to outcrossing based on *S. paradoxus* (Tsai *et al.*, 2008), and  $N_e \approx 10^8$  there were only  $10^{13}$  to  $10^{14}$  opportunities for such an introgression to have occurred and fixed. Clearly, unless there was a severe bottleneck with  $N_e < 1/|s| \approx 1,250$  around the time of introgression, which conceivably could have been triggered by a speciation event, this scenario seems very unlikely.

In the second scenario, where we assume the introgression contained advantageous loci, 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.

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.

## Materials and Methods

### Separating Endogenous and Exogenous Genes

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

332 (on 12-09-2014). We assigned 457 genes located on chromosome C (position 1 to 989,693 of chromosome  
333 C) with a location within the ~ 1 Mb window to the exogenous gene set. All other 4864 genes of the *L.*  
334 *kluyveri* genome were assigned to the exogenous genes. All genes could be assigned to one or the other  
335 gene set unambiguously.

## 336 Model Fitting with ROC SEMPPR

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

## 342 Model selection

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

## 348 Comparing Codon Specific Parameter Estimates

Choice of reference codon does reorganize codon families coding for an amino acid relative to each other,  
therefore all parameter estimates are relative to the mean for each codon family.

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

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

349 Comparison of codon specific parameters ( $\Delta M$  and  $\Delta \eta = 2N_e q(\eta_i - \eta_j)$ ) was performed using the function  
350 lmodel2 in the R package lmodel2 (1.7.3) (Legendre, 2018) and R version 3.4.1 (R Core Team, 2013).  
351 The parameter  $\Delta \eta$  can be interpreted as the difference in fitness between codon  $i$  and  $j$  for the average  
352 gene with  $\phi = 1$  scaled by the effective population size  $N_e$ , and the selective cost of an ATP  $q$  (Gilchrist,  
353 2007; Gilchrist *et al.*, 2015). Type II regression was performed with re-centered parameter estimates,  
354 accounting for noise in dependent and independent variable (Sokal and Rohlf, 1981).

355 **Phylogenetic Analysis**

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

362 **Synteny Comparison**

363 We obtained complete genome sequences from NCBI (on: 02-05-2017). Genomes were aligned and  
364 checked for synteny using SyMAP (4.2) with default settings (Soderlund *et al.*, 2006, 2011). We assess  
365 synteny as percentage coverage of the exogenous gene region.

366 **Estimating Age of Introgression**

We modelled 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 where taken from our mutation parameter estimates for *E. gossypii* where  $c_1(0) = \exp[\Delta M_{\text{gos}}]/(1 + \exp[\Delta M_{\text{gos}}])$ . Our estimates of  $\Delta M_{\text{endo}}$  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

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

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

368 Equation 5 was solved with a mutation rate  $\mu_{2,1}$  of  $3.8 \times 10^{-10}$  per nucleotide per generation (Lang  
369 and Murray, 2008). Current codon frequencies for each codon family where taken from our estimates of  
370  $\Delta M$  from the exogenous genes. Mathematica (11.3) (Wolfram Research Inc., 2017) was used to calculate

371 the time  $t_{\text{intro}}$  it takes for the initial codon frequencies  $c_1(0)$  for each codon family to equal the current  
 372 exogenous codon frequencies. The same equation was used to determine the time  $t_{\text{decay}}$  at which the  
 373 signal of the exogenous cellular environment has decayed to within 1% of the endogenous environment.

### 374 Estimating Genetic Load

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

385 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  
 386 if  $\kappa$  is a scaling on protein synthesis rate  $\phi$ , effective population size  $N_e$ , or the selective cost of an ATP  
 387  $q$  (Gilchrist, 2007; Gilchrist *et al.*, 2015). We calculated the genetic load each gene represents due to its  
 388 mismatched codon usage assuming additive fitness effects as

$$389 s_g = \sum_{i=1}^{L_g} -\kappa \phi_g \Delta\eta'_i \quad (6)$$

390 where  $s_g$  is the overall strength of selection for translational efficiency on gene,  $g$  in the exogenous gene  
 391 set,  $\kappa$  is a constant, scaling the efficacy of selection between the endogenous and exogenous cellular  
 392 environments,  $L_g$  is length of the protein in codons,  $\phi_g$  is the estimated protein synthesis rate of the gene  
 393 in the endogenous environment, and  $\Delta\eta'_i$  is the  $\Delta\eta'$  for the codon at position  $i$ . As stated previously, our  
 394  $\Delta\eta$  are relative to the mean of the codon family. We find that the genetic load of 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  
 395 endogenous and the current exogenous genes, and  $\kappa = 5$  for  $s_g$  for the genetic load at the time of  
 396 introgression.

397 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 genetic load due to mismatched codon usage of the introgression as  $E[s_g] = s_{\text{intro},g} - s_{\text{endo},g}$  where  $s_{\text{intro},g}$  is the genetic load of 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.

## References

- Baker, E. C., Wang, B., Bellora, N., *et al.* 2015. The genome sequence of *saccharomyces eubayanus* and the domestication of lager-brewing yeasts. *Molecular Biology and Evolution*, 32(11): 2818–2831.
- Beimforde, C., Feldberg, K., Nylander, S., *et al.* 2014. Estimating the phanerozoic history of the ascomycota lineages: combining fossil and molecular data. *Mol. Phylogenet. Evol.*, 78: 386–398.
- Brion, C., Legrand, S., Peter, J., *et al.* 2017. Variation of the meiotic recombination landscape and properties over a broad evolutionary distance in yeasts. *PLoS Genetics*, 13(8): e1006917.

- 412 Bulmer, M. 1990. The selection-mutation-drift theory of synonymous codon usage. *Genetics*, 129: 897–  
413 907.
- 414 Cope, A. L., Hettich, R. L., and Gilchrist, M. A. 2018. Quantifying codon usage in signal peptides:  
415 Gene expression and amino acid usage explain apparent selection for inefficient codons. *Biochimica et  
416 Biophysica Acta (BBA) - Biomembranes*, 1860(12): 2479–2485.
- 417 dos Reis, M., Savva, R., and Wernisch, L. 2004. Solving the riddle of codon usage preferences: a test for  
418 translational selection. *Nucleic Acids Research*, 32(17): 5036–5044.
- 419 Friedrich, A., Reiser, C., Fischer, G., and Schacherer, J. 2015. Population genomics reveals chromosome-  
420 scale heterogeneous evolution in a protoploid yeast. *Molecular Biology and Evolution*, 32(1): 184 –  
421 192.
- 422 Gilchrist, M. A. 2007. Combining models of protein translation and population genetics to predict protein  
423 production rates from codon usage patterns. *Molecular Biology and Evolution*, 24(11): 2362–2372.
- 424 Gilchrist, M. A., Chen, W. C., Shah, P., Landerer, C. L., and Zaretzki, R. 2015. Estimating gene  
425 expression and codon-specific translational efficiencies, mutation biases, and selection coefficients from  
426 genomic data alone. *Genome Biology and Evolution*, 7: 1559–1579.
- 427 Gouy, M. and Gautier, C. 1982. Codon usage in bacteria: correlation with gene expressivity. *Nucleic  
428 Acids Research*, 10: 7055–7074.
- 429 Gronau, Q. F., Sarafoglou, A., Matzke, D., *et al.* 2017. Ta tutorial on bridge sampling. *Journal of  
430 Mathematical Psychology*, 81: 80–97.
- 431 Ikemura, T. 1985. Codon usage and tRNA content in unicellular and multicellular organisms. *Molecular  
432 Biology and Evolution*, 2: 13–34.
- 433 Kensche, P. R., Oti, M., Dutilh, B. E., and Huynen, M. A. 2008. Conservation of divergent transcription  
434 in fungi. *Trends Genet.*, 5(24): 207–211.
- 435 Landerer, C., Cope, A., Zaretzki, R., and Gilchrist, M. A. 2018. Anacoda: analyzing codon data with  
436 bayesian mixture models. *Bioinformatics*, 34(14): 2496–2498.
- 437 Lang, G. I. and Murray, A. W. 2008. Estimating the per-base-pair mutation rate in the yeast saccha-  
438 romyces cerevisiae. *Genetics*, 178(1): 67 – 82.
- 439 Lawrence, J. G. and Ochman, H. 1997. Amelioration of bacterial genomes: Rates of change and exchange.  
440 *Journal of Molecular Miology*, 44: 383–397.

- 441 Legendre, P. 2018. *lmodel2: Model II Regression*. R package version 1.7-3.
- 442 Marcet-Houben, M. and Gabaldón, T. 2015. Beyond the whole-genome duplication: Phylogenetic ev-  
443 idence for an ancient interspecies hybridization in the baker's yeast lineage. *PLoS Biology*, 13(8):  
444 e1002220.
- 445 Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., and Danchin, A. 1991. Evidence for horizontal gene  
446 transfer in escherichia coli speciation. *Journal of Molecular Miology*, 222(4): 851–856.
- 447 Nguyen, L. T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. 2015. Iq-tree: A fast and effective  
448 stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution*,  
449 32(1): 268–274.
- 450 Payen, C., Fischer, G., Marck, C., *et al.* 2009. Unusual composition of a yeast chromosome arm is  
451 associated with its delayed replication. *Genome Research*, 19(10): 1710–1721.
- 452 R Core Team 2013. *R: A Language and Environment for Statistical Computing*. R Foundation for  
453 Statistical Computing, Vienna, Austria.
- 454 Salichos, L. and Rokas, A. 2013. Inferring ancient divergences requires genes with strong phylogenetic  
455 signals. *Nature*, 497: 327–331.
- 456 Sella, G. and Hirsh, A. E. 2005. The application of statistical physics to evolutionary biology. *Proceedings  
457 of the National Academy of Sciences of the United States of America*, 102: 9541–9546.
- 458 Shah, P. and Gilchrist, M. A. 2011. Explaining complex codon usage patterns with selection for trans-  
459 lational efficiency, mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences  
460 U.S.A*, 108(25): 10231–10236.
- 461 Sharp, P. M. and Li, W. H. 1987. The codon adaptation index - a measure of directional synonymous  
462 codon usage bias, and its potential applications. *Nucleic Acids Research*, 15: 1281–1295.
- 463 Shen, X. X., Opulente, D. A., Kominek, J., *et al.* 2018. Tempo and mode of genome evolution in the  
464 budding yeast subphylum. *Cell*, 175(6): 1533–1545.e20.
- 465 Soderlund, C., Nelson, W., Shoemaker, A., and Paterson, A. 2006. Symap A system for discovering and  
466 viewing syntenic regions of fpc maps. *Genome Research*, 16: 1159 – 1168.
- 467 Soderlund, C., Bomhoff, M., and Nelson, W. 2011. Symap v3.4: a turnkey synteny system with applica-  
468 tion to plant genomes. *Nucleic Acids Research*, 39(10): e68.

- 469 Sokal, R. R. and Rohlf, F. J. 1981. *Biometry - The principles and practice of statistics in biological,*  
470 pages 547–555. W. H. Freeman.
- 471 Tsai, I. J., Bensasson, D., Burt, A., and Koufopanou, V. 2008. Population genomics of the wild yeast  
472 *saccharomyces paradoxus*: quantifying the life cycle. *Proc Natl Acad Sci U.S.A.*, 105: 4957–4962.
- 473 Tsankov, A. M., Thompson, D. A., Socha, A., Regev, A., and Rando, O. J. 2010. The role of nucleosome  
474 positioning in the evolution of gene regulation. *PLoS Biol*, 8(7): e1000414.
- 475 Vakirlis, N., Sarilar, V., Drillon, G., *et al.* 2016. Reconstruction of ancestral chromosome architecture  
476 and gene repertoire reveals principles of genome evolution in a model yeast genus. *Genome research*,  
477 26(7): 918–32.
- 478 Wagner, A. 2005. Energy constraints on the evolution of gene expression. *Molecular Biology and Evolu-*  
479 *tion*, 22: 1365–1374.
- 480 Wallace, E. W., Airoldi, E. M., and Drummond, D. A. 2013. Estimating selection on synonymous codon  
481 usage from noisy experimental data. *Molecular Biology and Evolution*, 30: 1438–1453.
- 482 Wolfram Research Inc. 2017. *Mathematica 11*.
- 483 Wright, F. 1990. The ‘effective number of codons’ used in a gene. *Genet*, 87: 23–29.

## Supplementary Material

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

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

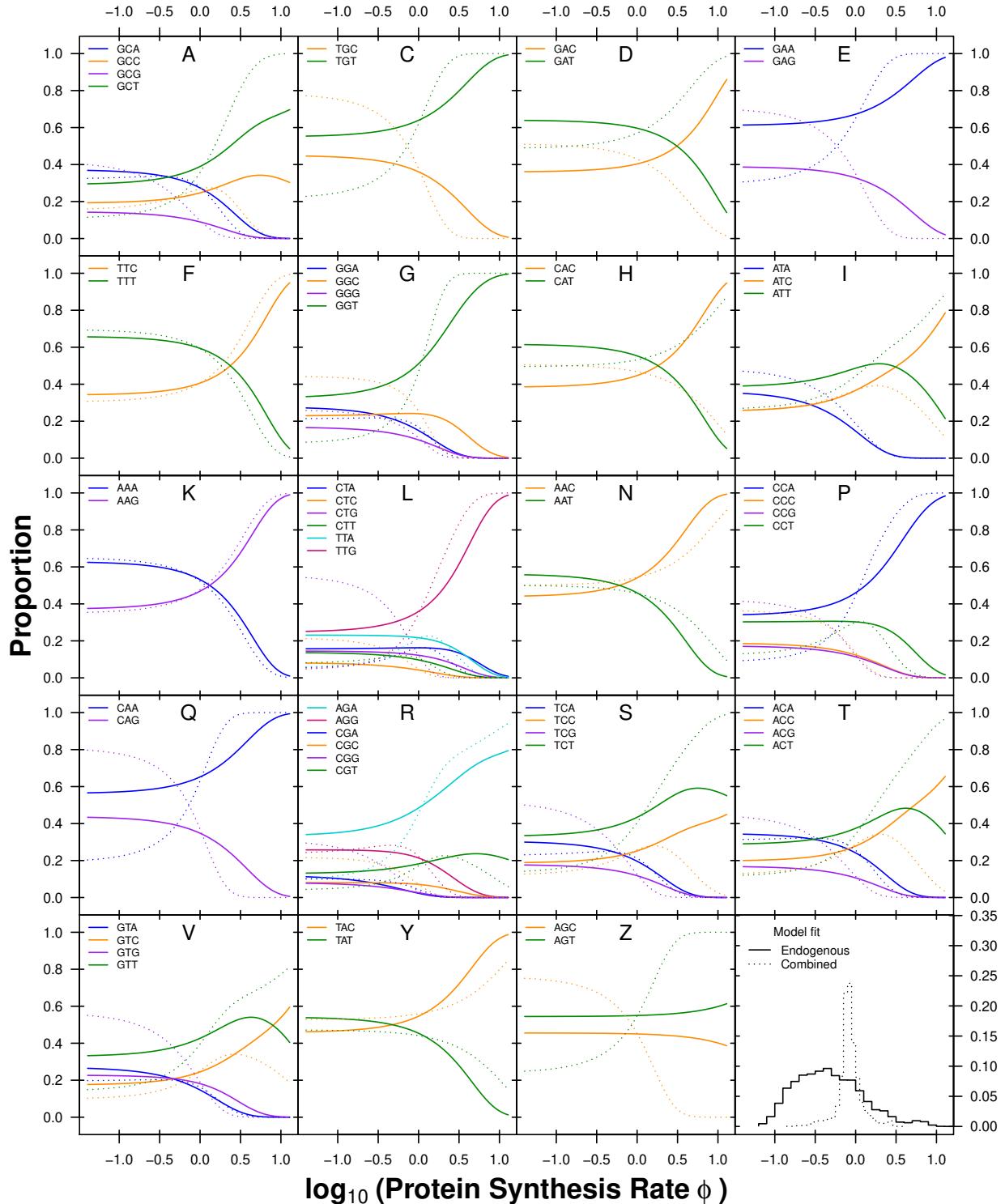


Figure S1: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dotted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

## Exogenous and Combined Codon Usage

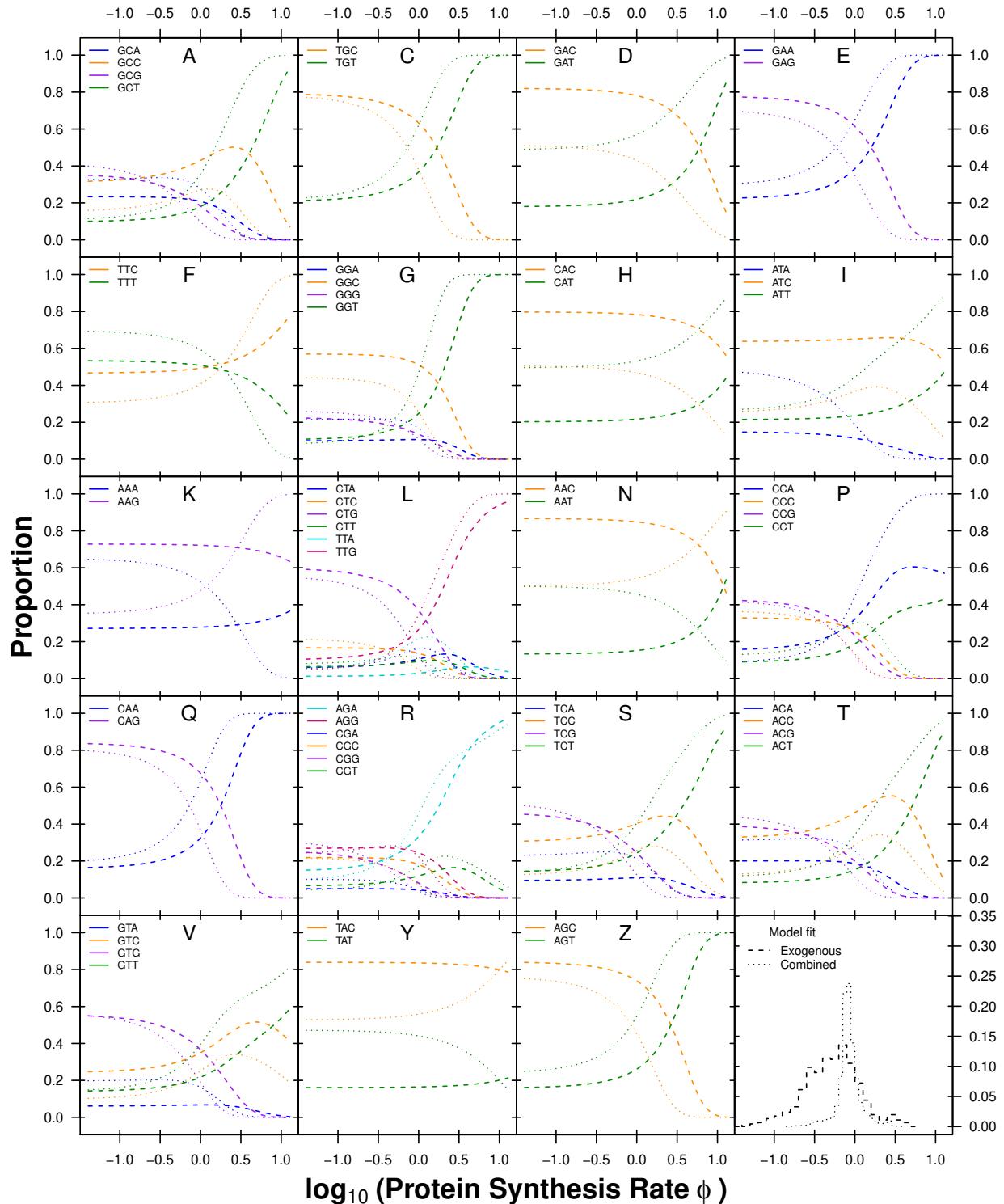


Figure S2: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. dashed line indicates the exogenous codon usage, dotted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

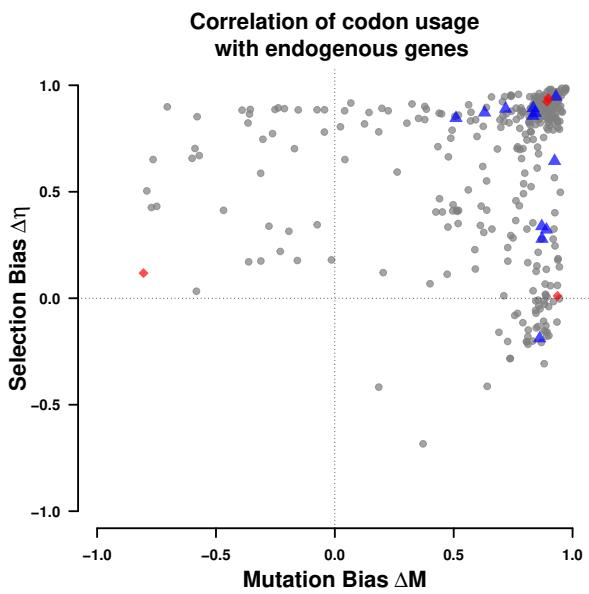


Figure S3: Correlation coefficients of  $\Delta M$  and  $\Delta\eta$  of the endogenous genes with 332 examined budding yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta\eta$  of the lineages with the exogenous parameter estimates. Blue triangles indicate the Lachancea and red 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).

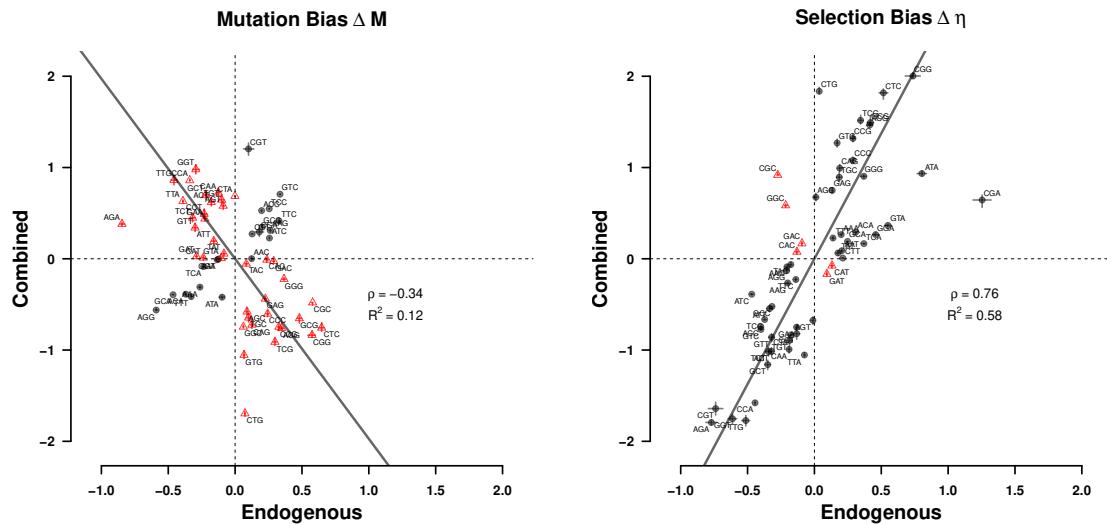


Figure S4: Comparison of (a) mutation bias  $\Delta M$  and (b) selection bias  $\Delta \eta$  parameters for endogenous genes and combined gene sets. Estimates are relative to the mean for each codon family. Black dots indicate  $\Delta M$  or  $\Delta \eta$  parameters with the same sign for the endogenous 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.

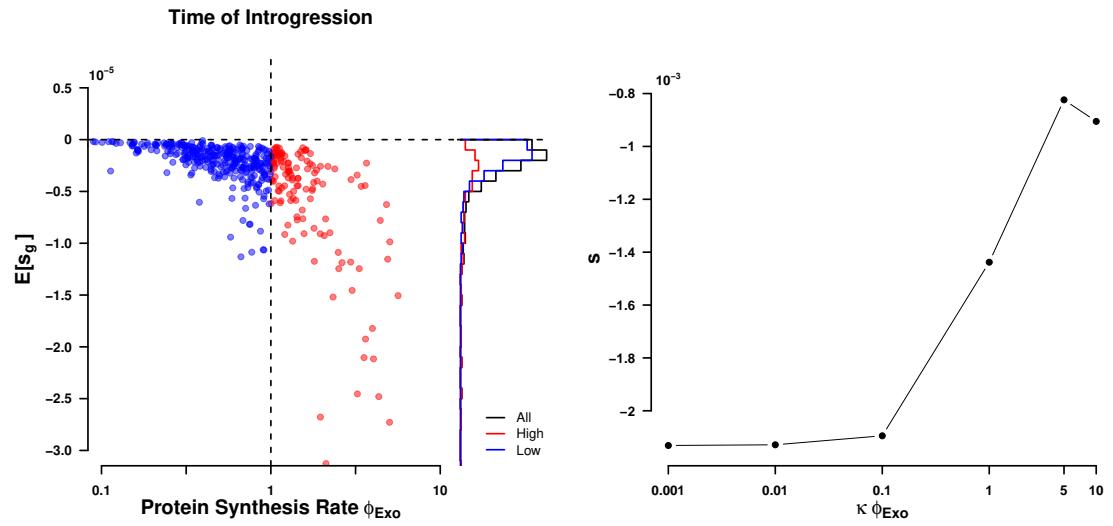


Figure S5: Genetic load (left) without scaling of  $\phi$  per gene. Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates a genetic load of 0. (Right) Change of total genetic load with scaling term  $\kappa$  between *E. gossypii* and *L. kluyveri*

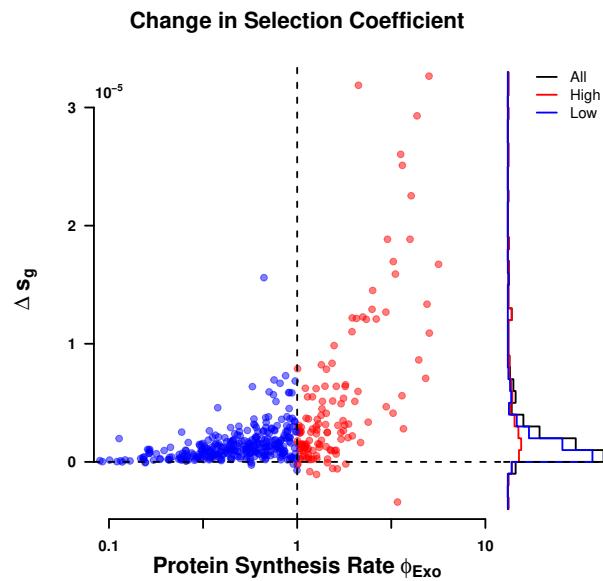


Figure S6: Total amount of adaptation estimated to have occurred between time of introgression and currently observed per gene. Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates a genetic load of 0.

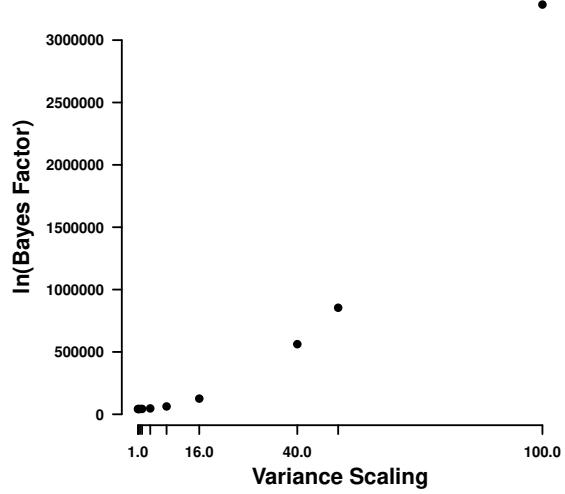


Figure S7: Influence of the variance scaling of the importance distribution on the estimated Bayes factor.