

## RESEARCH

1  
2  
3  
4  
5  
6  
7  
8

# Unlocking a signal of introgression from codons in *Lachancea kluyveri* using a mutation-selection model

9 Cedric Landerer<sup>1,2,3\*</sup>, Brian C O'Meara<sup>1,2</sup>, Russell Zaretzki<sup>2,4</sup> and Michael A Gilchrist<sup>1,2</sup>

|    |  |
|----|--|
| 10 |  |
| 11 |  |
| 12 |  |
| 13 |  |
| 14 |  |
| 15 |  |
| 16 |  |
| 17 |  |
| 18 |  |
| 19 |  |
| 20 |  |
| 21 |  |
| 22 |  |
| 23 |  |
| 24 |  |
| 25 |  |
| 26 |  |
| 27 |  |
| 28 |  |
| 29 |  |
| 30 |  |
| 31 |  |
| 32 |  |
| 33 |  |

Correspondence:  
edric.landerer@gmail.com  
Max-Planck Institute of  
Molecular Cell Biology and  
Genetics, Pfotenhauerstr. 108,  
1307, Dresden, Germany  
Full list of author information is  
available at the end of the article  
Correspondence

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

## Abstract

**Background:** For decades, codon usage has been used as a measure of adaptation for translational efficiency and translation accuracy 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 about protein synthesis and grounded in population genetics.

**Results:** We find substantial differences in codon usage between the endogenous and exogenous genes, and show that these differences can be largely attributed to differences in mutation bias favoring A/T ending codons in the endogenous genes while favoring 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 the decaying signal of endogenous codon mutation and 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 that the introgression occurred  $\sim 6 \times 10^8$  generation ago, and estimate its historic and current selection against mismatched codon usage.

**Conclusions:** Our work illustrates how mechanistic, population genetic models like ROC SEMPPR can separate the effects of mutation and selection on codon usage and provide quantitative estimates from sequence data.

**Keywords:** codon usage; population genetics; introgression; mutation; selection

## 1      **Background**

2      Synonymous codon usage patterns varies within a genome and between taxa, re-  
3      flecting differences in mutation bias, selection, and genetic drift. The signature of  
4      mutation bias is largely determined by the organism's internal or cellular environ-  
5      ment, such as their DNA repair genes or UV exposure. While this mutation bias  
6      is an omnipresent evolutionary force, its impact can be obscured or amplified by  
7      selection. The signature of selection on codon usage is largely determined by an or-  
8      ganism's cellular environment alone, such as, but not limited to, its tRNA species,  
9      their copy number, and their post-transcriptional modifications. In general, the  
10     strength of selection on codon usage is assumed to increase with its expression level  
11     [1–3], specifically its protein synthesis rate [4]. Thus as protein synthesis increases,  
12     codon usage shifts from a process dominated by mutation to a process dominated  
13     by selection. The overall efficacy of mutation and selection on codon usage is a  
14     function of the organism's effective population size  $N_e$ . ROC SEMPPR allows us  
15     to disentangle the evolutionary forces responsible for the patterns of codon usage  
16     bias [5–7] (CUB) encoded in an species' genome, by explicitly modeling the com-  
17     bined evolutionary forces of mutation, selection, and drift [4, 8–10]. In turn, these  
18     evolutionary parameters should provide biologically meaningful information about  
19     the lineage's historical cellular and external environment.

20     Most studies implicitly assume that the CUB of a genome is shaped by a single  
21     cellular and external environment. However, this assumption is clearly violated to  
22     increasing degrees via horizontally gene transfer, large scale introgressions, and hy-  
23     brid specie formation. In these scenarios, one would expect to see the signature of  
24     multiple cellular environments in a genome's CUB [11, 12]. Indeed, differences in  
25     CUB between linages have been proposed to have a major effect on their rates of  
26     gene transfer with rates declining with differences in their CUB. On a more practical  
27     level, if differences in codon usage of transferred genes are not taken into account  
28     for, they may distort the interpretation of codon usage patterns. Such distortion  
29     could lead to the wrong inference of codon preference for an amino acid [8, 10], un-  
30     derestimate the variation in protein synthesis rate, or distort estimates of mutation  
31     bias when analyzing a genome.

32     To illustrate these ideas, we analyze the CUB of the genome of the yeast *Lachancea*  
33     *kluyveri* using ROC SEMPPR, a population genetics based model of synonymous

1 codon usage evolution that accounts for and, in turn, can estimate the contribution  
2 of mutation bias  $\Delta M$ , selection bias. The mathematics of ROC SEMPPR are de-  
3 rived on a mechanistic description of ribosome movement along an mRNA, although  
4 the approximation of other biological mechanisms could also be consistent with the  
5 model. Broadly speaking, ROC SEMPPR allows us to quantify the cellular environ-  
6 ment in which genes have evolved by separately estimating the effects of mutation  
7 bias and selection bias on codon usageDE between synonymous codons and pro-  
8 tein synthesis rate  $\phi$  to the patterns of codon usage observed within a set of genes.  
9 Briefly, the set of  $\Delta M$  for an amino acid quantifies the relative differences in muta-  
10 tional stability or bias between the synonymous codons of the amino acid  $S$ . In the  
11 absence of selection bias (or equivalently when gene expression  $\phi = 0$ ), the equilib-  
12 rium frequency of synonymous codon  $i$  is simply  $\exp[-\Delta M_i] / \left( \sum_{j \in S} \exp[-\Delta M_j] \right)$ .  
13 Because the time units of protein production rate have no intrinsic time scale, we  
14 define the average protein production rate for a set of genes to be one, i.e.  $\bar{\phi} = 1$   
15 by definition [10]. In order to facilitate comparisons between gene sets, we express  
16 both,  $\Delta M$  and  $\Delta \eta$ , as deviation from the mean of each synonymous codon family  
17 (see Materials and Methods for details). Nevertheless, the difference  $\Delta \eta$  describes  
18 the difference in fitness between two synonymous codons relative to drift for a gene  
19 whose protein production rate  $\phi$  is equal to the the average rate of protein produc-  
20 tion  $\bar{\phi}$  across the set of genes. In other words, for a gene whose protein is expressed  
21 at the average rate, for any two given synonymous codons  $i$  and  $j$ ,  $\Delta \eta_i - \Delta \eta_j = N_e s$ .

22  
23 The *Lachancea* clade diverged from the *Saccharomyces* clade, prior to its whole  
24 genome duplication  $\sim 100$  Mya ago [13, 14]. Since that time, *L. kluyveri*, which is  
25 sister species to all other *Lachancea* spp., has experienced a large introgression of  
26 exogenous genes (1 Mb, 457 genes) which is found in all of its populations [15, 16],  
27 but in no other known *Lachancea* species [17]. The introgression replaced the left  
28 arm of the C chromosome and displays a 13% higher GC content than the en-  
29 doogenous *L. kluyveri* genome [15, 16]. Previous studies suggest that the source of  
30 the introgression is probably a currently unknown or potentially extinct *Lachancea*  
31 lineage based on gene concatenation or synteny relationships [15–18]. These char-  
32 acteristics make *L. kluyveri* an ideal model to study the effects of an introgressed  
33 cellular environment and the resulting mismatch in codon usage.

1 While previous studies [8, 9] have used information on gene expression to sepa-  
2 rate the effects of mutation and selection on codon usage, ROC SEMPPR does not  
3 need such information but can provide it. ROC SEMPPR's resulting predictions  
4 of protein synthesis rates have been shown to be on par with laboratory measure-  
5 ments [8, 10]. In contrast to often used heuristic approaches to study codon usage  
6 [5, 6, 19], ROC SEMPPR explicitly incorporates and distinguishes between mu-  
7 tation and selection effects on codon usage and properly weights its estimates by  
8 amino acid usage [20]. We use ROC SEMPPR to separately describe the two cellular  
9 environments reflected in the *L. kluyveri* genome; the signature of the endogenous  
10 environment reflected in the larger set of non-introgressed genes and the decaying  
11 signature of the ancestral, exogenous environment in the smaller set of introgressed  
12 genes. Our results indicate that the current difference in GC content between en-  
13 dogenous and exogenous genes is mostly due to the differences in mutation bias  
14  $\Delta M$  of their respective cellular environments. Taking the different signatures of  
15  $\Delta M$  and selection bias  $\Delta \eta$  of the endogenous and exogenous sets of genes substan-  
16 tially improves our ability to predict present day protein synthesis rates  $\phi$ . These  
17 endogenous and exogenous gene set specific estimates of  $\Delta M$  and  $\Delta \eta$ , in turn, allow  
18 us to address more refined biological questions. For example, we find support for  
19 an alternative origin of the exogenous genes and identify *E. gossypii* as the nearest  
20 sampled relative of the source of the introgressed genes out of the 332 budding yeast  
21 lineages with sequenced genomes [21]. While this inference is in contrast to previous  
22 work [15–18], we find additional phylogenetic support for via gene tree reconstruc-  
23 tion and gene synteny. We also estimate the age of the introgression to be on the  
24 order of 0.2 - 1.7 Mya, estimate the selection against these genes, both at the time  
25 of introgression and now, and predict a detectable signature of CUB to persist in  
26 the introgressed genes for another 0.3 - 2.8 Mya, highlighting the sensitivity of our  
27 approach.  
28

## 29 Results

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

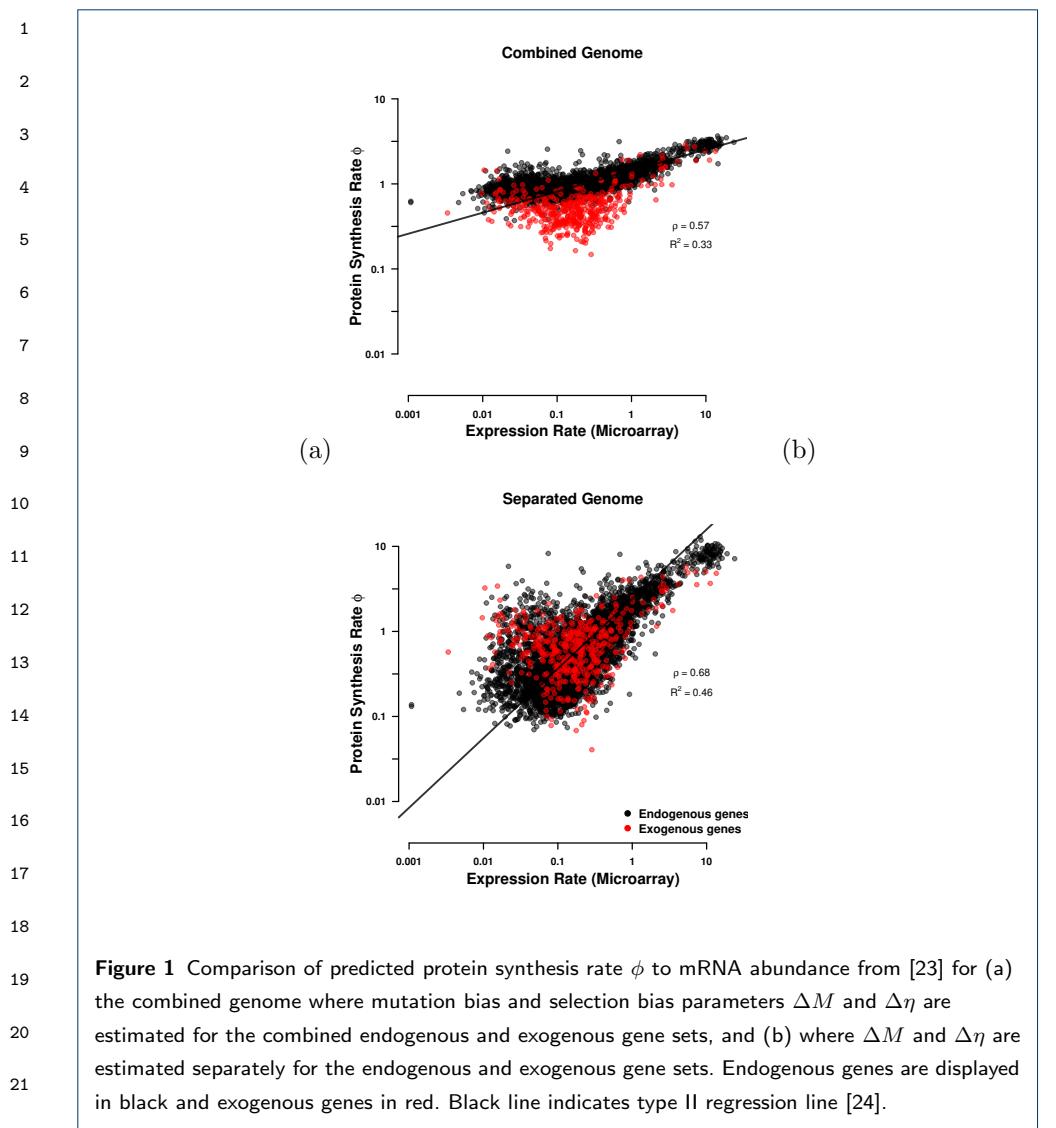
31 We used our software package AnaCoDa [22] to compare model fits of ROC  
32 SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two  
33 sets of 4,864 endogenous and 497 exogenous genes. These two set where initially

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

| Hypothesis | $\log(\mathcal{L})$ | $n$   | $\log(\mathcal{L}_M)$ | $\log(K)$ | $p$ |
|------------|---------------------|-------|-----------------------|-----------|-----|
| Combined   | -2,650,047          | 5,483 | -2,657,582            | —         | —   |
| Separated  | -2,612,397          | 5,402 | -2,615,288            | 42,294    | 0   |

7  
 8 identified based on their striking difference in GC content [15], with very little over-  
 9 lap in GC content between the two sets (Figure S1a). ROC SEMPPR is a statistical  
 10 model that relates the effects of mutation bias  $\Delta M$ , selection bias  $\Delta\eta$  between syn-  
 11 onymous codons and protein synthesis rate  $\phi$ , to explain the observed codon usage  
 12 patterns. Thus, the probability of observing a synonymous codon is proportional  
 13 to  $p \propto \exp(-\Delta M - \Delta\eta\phi)$  [10]. Briefly,  $\Delta M$  describes the mutation bias between  
 14 two synonymous codons at stationarity under a time reversible mutation model.  
 15 Because ROC SEMPPR only considers the stationary probabilities, only variation  
 16 in mutation bias, not absolute mutation rates can be detected.  $\Delta\eta$  describes the  
 17 fitness difference between two synonymous codons relative to drift [10]. Since  $\Delta\eta$  is  
 18 scaled by protein synthesis rate  $\phi$ , this term is dominant in highly expressed genes  
 19 and tends towards 0 in low expression genes, allowing us to separate the effect of  
 20 mutation bias and selection bias on codon usage. We express both,  $\Delta M$  and  $\Delta\eta$ ,  
 21 as deviation from the mean of each synonymous codon family which prevents that  
 22 the choice of the reference codon affects our results (see Materials and Methods for  
 23 details).

24  
 25 Bayes factor strongly support the hypothesis that the *L. kluyveri* genome consists  
 26 of genes with two different and distinct patterns of codon usage bias rather than a  
 27 single ( $K = \exp(42,294)$ ; Table 1). We find additional support for this hypothesis  
 28 when we compare our predictions of protein synthesis rate to empirically observed  
 29 mRNA expression values as a proxy for protein synthesis. Specifically, we improve  
 30 the variance explained by our predicted protein synthesis rates by  $\sim 42\%$ , from  $R^2 =$   
 31  $0.33 (p < 10^{10})$  to  $0.46 (p < 10^{10})$  (Figure 1). While the implicit consideration of GC  
 32 content in this analysis certainly plays a roll, it does not explain the improvement  
 33 in  $R^2$  (Figure S1b).



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

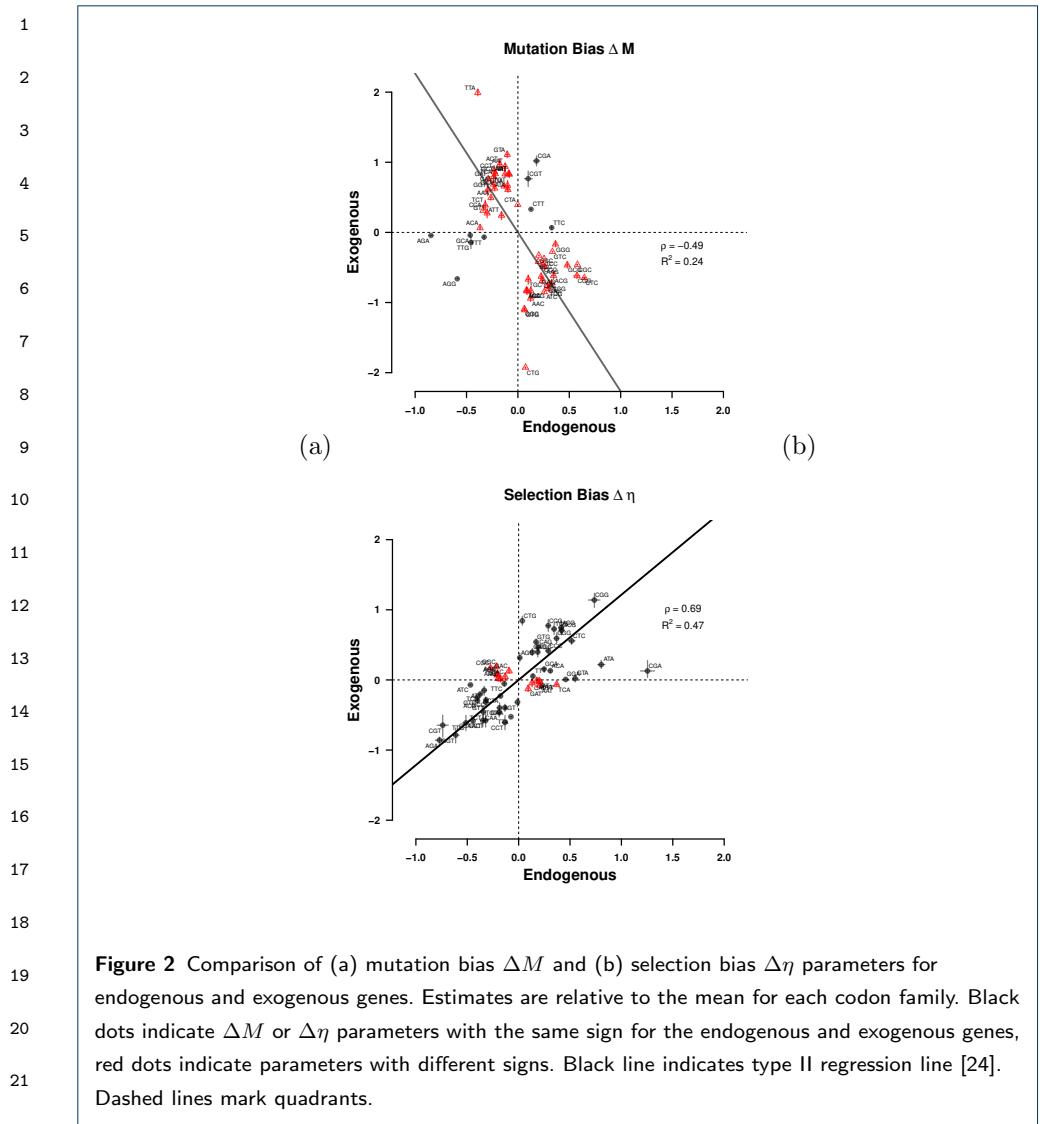
Because ROC SEMPPR defines  $\bar{\phi} = 1$ , it makes the interpretation of  $\Delta \eta$  as selection on codon usage of the average gene with  $\phi = 1$  straightforward and gives us the ability to compare the efficacy of selection  $sN_e$  across genomes. While it may be expected for the endogenous and exogenous genes to differ in their codon usage pattern due to the large difference in GC content it is not clear how much of this difference is due to differences in the mutation bias  $\Delta M$  or selection bias  $\Delta \eta$  between the gene sets. To better understand the differences in the endogenous and exogenous cellular environments, we compared our parameter estimates of  $\Delta M$  and  $\Delta \eta$  for the two sets of genes. Our estimates of  $\Delta M$  for the endogenous and

1 exogenous genes were negatively correlated ( $\rho = -0.49, p = 3.56 \times 10^{-5}$ ), indicating  
2 weak similarity with only  $\sim 5\%$  of the codons share the same sign between the two  
3 mutation environments (Figure 2a). Overall, the endogenous genes only show a  
4 selection preference for C and G ending codons in  $\sim 58\%$  of the codon families.  
5 In contrast, the exogenous genes display a strong preference for A and T ending  
6 codons in  $\sim 89\%$  of the codon families.

7 For example, the endogenous genes show a mutational bias for A and T ending  
8 codons in  $\sim 95\%$  of the codon families. The exception is Leucine (Leu, L), where  
9 mutation appears to favor the codon TTG over TTA (Figure 3, Table S1). The  
10 exogenous genes display an equally consistent mutational bias towards C and G  
11 ending codons (the exception being Phe, F). In contrast to  $\Delta M$ , our estimates of  $\Delta\eta$   
12 for the endogenous and exogenous genes were positively correlated ( $\rho = 0.69$ ) and  
13 showing the same sign in  $\sim 53\%$  of codons between the two selection environments  
14 (Figure 2).

15 We find that the signature of selection bias  $\Delta\eta$  also differs substantially between  
16 the endogenous and exogenous gene sets. The difference in codon usage between  
17 endogenous and exogenous genes is striking as the sign for  $\Delta\eta$  changes, indicating a  
18 change in codon preference for some amino acids. As a result, our estimates of the  
19 optimal codon differ in nine cases between endogenous and exogenous genes (Figure  
20 3, Table S2). For example, the usage of the Asparagine (Asn, N) codon AAC is  
21 increased in highly expressed endogenous genes but the same codon is depleted in  
22 highly expressed exogenous genes. For Aspartic acid (Asp, D), the combined genome  
23 shows the same codon preference in highly expressed genes as the exogenous gene  
24 set. Generally, fits to the complete *L. kluyveri* genome reveal that the relatively  
25 small exogenous gene set ( $\sim 10\%$  of genes) has a disproportionate effect on the  
26 model fit (Figure S2, S3).

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



**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 [24]. Dashed lines mark quadrants.

to  $\Delta M$ , our estimates of selection bias  $\Delta \eta$  for the endogenous and exogenous genes are positively correlated ( $\rho = 0.69$ ,  $p = 9.76 \times 10^{-10}$ ) and show the same sign in  $\sim 53\%$  of the cases (Figure 2).

This striking difference in codon usage was noted previously. For example, using RSCU [5], GAA (coding for Glu, E) was identified as the optimal synonymous codon in the whole genome and GAG as the optimal codon in the exogenous genes [15]. Our results, however, indicate that GAA is the optimal codon in both, endogenous and exogenous genes, and that the high RSCU in the exogenous genes of GAG is driven by mutation bias (Table S1 and S2). Similar effects are observed for other amino acids.

1 The effect of the small exogenous gene set on the fit to the complete *L. kluyveri*  
2 genome is smaller for our estimates of selection bias  $\Delta\eta$  than  $\Delta M$ , but still large.  
3 We find that the complete *L. kluyveri* genome is estimated to share the selectively  
4 preferred codon with the exogenous genes in  $\sim 60\%$  of codon families that show dis-  
5 similarity between endogenous and exogenous genes. We also find that the complete  
6 *L. kluyveri* genome fit shares mutationally preferred codons with the exogenous  
7 genes in  $\sim 78\%$  of the 19 codon families showing a difference in mutational codon  
8 preference between the endogenous and exogenous genes. In two cases, Isoleucine  
9 (Ile, I) and Arginine (Arg, R), the strong dissimilarity in mutation preference results  
10 in an estimated codon preference in the complete *L. kluyveri* genome that differs  
11 from both the endogenous, and the exogenous genes. These results clearly show that  
12 it is important to recognize the difference in endogenous and exogenous genes and  
13 treat these genes as separate sets to avoid the inference of incorrect synonymous  
14 codon preferences and better predict protein synthesis.

15 In order to test the robustness of our results we randomized the di-nucleotides  
16 of the endogenous and exogenous genes while maintaining their frequency and the  
17 GC-content. We find that we can still separate the effect of mutation and selection  
18 on codon usage. The randomized gene sets show a similar pattern of mutation bias  
19 and selection bias as the canonical endogenous (Figure S10 a,b) and the exogenous  
20 (Figure S11 a,b) gene sets. Similar, the predicted protein synthesiis rate  $\phi$  of the  
21 endogenous ( $\rho = 0.73, p < 10^{-10}$ ) and exogenous ( $\rho = 0.88, p = p < 10^{-10}$ )  
22 randomized gene set is highly correlated with our canonical predictions of  $\phi$  (Figure  
23 S13 a,b).

## 25 Can Codon Usage Help Determine the Source of the Exogenous Genes

26 Since the origin of the exogenous genes is currently unknown, we explored if the  
27 information on codon usage extracted from the exogenous genes can be used to  
28 identify a potential source lineage. We combined our estimates of mutation bias  
29  $\Delta M$  and selection bias  $\Delta\eta$  with synteny information and searched for potential  
30 source lineages of the introgressed exogenous region. We used  $\Delta M$  to identify can-  
31 didate lineages as the endogenous and exogenous genes show greater dissimilarity  
32 in mutation bias than in selection bias. We examined 332 budding yeasts [21] and,  
33 identified the ten lineages with the highest correlation to the exogenous  $\Delta M$  pa-

**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 exogenous  $\Delta M$  and  $\Delta \eta$  with the indicated species', 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

parameters as potential source lineages (Figure 4, Table 2). 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 from our comparison of codon families; however, there was no need to exclude Serine as CTG is not a one step neighbor of the remaining Serine codons. A mutation between CTG and the remaining Serine codons would require two mutations with one of them being non-synonymous, which would violate the weak mutation assumption of ROC SEMPPr.

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 S4). Only ~ 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 (~ 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 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. Previous studies which studied the exogenous genes and

1 chromosome recombination in the Lachancea clade concluded that the exogenous  
2 region originated from within the Lachancea clade, from an unknown or potentially  
3 extinct lineage [15–17]. While it is not possible for us to dispute this hypothesis,  
4 our results provide a novel hypothesis about the origin of the exogenous genes.  
5

6 To further test the plausibility of *E. gossypii* as potential source lineage, we iden-  
7 tified 127 genes in our dataset [21] with homologous genes in *E. gossypii* and other  
8 Lachancea and used IQTree [25] to infer the phylogenetic relationship of the exoge-  
9 nous genes. Our results show that at least  $\sim 45\%$  of exogenous genes (57/127) are  
10 more closely related to *E. gossypii* than to other Lachancea S5. Interestingly, our re-  
11 sults also indicate that codon usage does not necessarily correlate with phylogenetic  
12 distance (Table 2).  
13

### 14 Estimating Introgression Age

### 15

16 If we assume that the exogenous genes originated from the *E. gossypii* lineage, we  
17 can estimate the age of the introgression based on our estimates of mutation bias  
18  $\Delta M$ . We modeled the change in codon frequency over time as exponential decay,  
19 and estimated the age of the introgression assuming that *E. gossypii* still represents  
20 the mutation bias of its ancestral source lineage at the time of the introgression and  
21 a constant mutation rate. We infer the age of the introgression to be on the order  
22 of  $6.2 \pm 1.2 \times 10^8$  generations. Assuming *L. kluyveri* experiences between one and  
23 eight generations per day, we estimate the introgression to have occurred between  
24 212,000 to 1,700,000 years ago. Our estimate places the time of the introgression  
earlier than the previous estimate of 19,000 - 150,000 years by [16].  
25

26 Using our model of exponential decay model, we also estimated the persistence of  
27 the signal of the exogenous cellular environment. We predict that the  $\Delta M$  signal of  
28 the source cellular environment will have decayed to be within one percent of the  
29 *L. kluyveri* environment in  $\sim 5.4 \pm 0.2 \times 10^9$  generations, or between 1,800,000 and  
15,000,000 years. Together, these results indicate that the mutation signature of  
the exogenous genes will persist for a very long time.  
30

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

### 32

33 We define the selection against inefficient codon usage as the difference between the  
fitness on the log 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  
34

1 details). As the introgression occurred before the diversification of *L. kluyveri* and  
2 has fixed throughout all populations [16], we can not observe the original endogenous  
3 sequences that have been replaced by the introgression. Overall, we predict that a  
4 small number of low expression genes ( $\phi < 1$ ) were weakly exapted at the time of the  
5 introgression (Figure 5a). Thus, they appear to provide a small fitness advantage  
6 due to the accordance of exogenous mutation bias with endogenous selection bias  
7 (compare Figure S2 and S3). High expression genes ( $\phi > 1$ ) are predicted to have  
8 faced the largest selection against their mismatched codon usage in the novel cellular  
9 environment. In order to account for differences in the efficacy of selection on codon  
10 usage either due to the cost of pausing, differences in the effective population size,  
11 or the decline in fitness with every ATP wasted between the donor lineage and *L.*  
12 *kluyveri* we added a linear scaling factor  $\kappa$  to scale our estimates of  $\Delta\eta$  between the  
13 donor lineage and *L. kluyveri* and searched for the value that minimized the cost of  
14 the introgression, thus giving us the best case scenario (See Methods for details).

15 Using our estimates of  $\Delta M$  and  $\Delta\eta$  from the endogenous genes and assuming the  
16 current exogenous amino acid composition of genes is representative of the replaced  
17 endogenous genes, we estimate the strength of selection against the exogenous genes  
18 at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of  
19 selection bias for the exogenous genes show that, while well correlated with the  
20 endogenous genes, only nine amino acids share the same selectively preferred codon.  
21 Exogenous genes are, therefore, expected to represent a significant reduction in  
22 fitness for *L. kluyveri* due to mismatch in codon usage. Since  $\Delta\eta$  is proportional  
23 to the difference in fitness between the wild type and a mutant, we can use our  
24 estimates of  $\Delta\eta$  to approximate the selection against the exogenous genes  $\Delta s$  [10,  
25 26]. We estimate that the selection against all exogenous genes due to mismatched  
26 codon usage to have been  $\Delta s \approx -0.0008$  at the time of the introgression and  
27  $\approx -0.0003$  today. This reduction in  $\Delta s$  is primarily due to adaptive changes to the  
28 codon usage of the most highly expressed, introgressed genes (Figures 5a & S8).  
29 Based on the selection against the codon mismatch at the time of the introgression  
30 and assuming an effective population size  $N_e$  on the order of  $10^7$  [27], we estimate  
31 a fixation probability of  $(1 - \exp[-\Delta s])/(1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$  [26] for the  
32 exogenous genes. Clearly, the possibility of fixation under this simple scenario is  
33 effectively zero. In order for the exogenous genes to have reached fixation one or

1 more exogenous loci must have provided a selective advantage not considered in  
2 this study (See Discussion).

3

4

## 5 Discussion

6 In order to study the evolutionary effects of the large scale introgression of the left  
7 arm of chromosome C, we used ROC SEMPPR, a mechanistic model of ribosome  
8 movement along an mRNA. The usage of a mechanistic model rooted in popula-  
9 tion genetics allows us generate more nuanced quantitative parameter estimates  
10 and separate the effects of mutation and selection on the evolution of codon usage.  
11 This allowed us to calculate the selection against the introgression, and provides *E.*  
12 *gossypii* as a potential source lineage of the introgression which was previously not  
13 considered. Our parameter estimates indicate that the *L. kluyveri* genome contains  
14 distinct signatures of mutation and selection bias from both an endogenous and ex-  
15 ogenous cellular environment. By fitting ROC SEMPPR separately to *L. kluyveri*'s  
16 endogenous and exogenous sets of genes we generate a quantitative description of  
17 their signatures of mutation bias and natural selection for efficient protein transla-  
18 tion.

19 In contrast to other methods such as RSCU, CAI, or tAI, ROC SEMPPR does  
20 not rely on external information such as gene expression or tRNA gene copy number  
21 [5, 19]. Instead, ROC SEMPPR allows for the estimation of protein synthesis rate  $\phi$   
22 and separates the effects of mutation and selection on codon usage. In addition, [20]  
23 showed that approaches like CAI are sensitive to amino acid composition, another  
24 property that distinguishes the endogenous and exogenous genes [15].

25 Previous work by [15] showed an increased bias towards GC rich codons in the  
26 exogenous genes but our results provide more nuanced insights by separating the  
27 effects of mutation bias and selection. We are able to show that the difference in GC  
28 content between endogenous and exogenous genes is mostly due to differences in  
29 mutation bias as 95% of exogenous codon families show a strong mutation bias to-  
30 wards GC ending codons (Table S1). However, the exogenous genes show a selective  
31 preference for AT ending codons for 90% of codon families (Table S2). Acknowl-  
32 edging the increased mutation bias towards GC ending codons and the difference in  
33 strength of selection between endogenous and exogenous genes by separating them

1 also improves our estimates of protein synthesis rate  $\phi$  by 42% relative to the full  
2 genome estimate ( $R^2 = 0.46, p = 0$  vs.  $0.32, p = 0$ , respectively).

3 Previous studies showed that nucleotide composition can be strongly affected by  
4 biased gene conversion, which, in turn would affect codon usage. Biased gene conver-  
5 sion is thought to act similar to directional selection, typically favoring the fixation  
6 of G/C alleles [28, 29]. Further, [30, Harrison & Charlesworth] suggested that bi-  
7 ased gene conversion affects codon usage in *S. cerevisiae*. ROC SEMPPR, however,  
8 does not explicitly account for biased gene conversion. If biased gene conversion is  
9 independent of gene expression, as in the case of DNA repair, it will be absorbed  
10 in our estimates of  $\Delta M$ . If instead biased gene conversion forms hotspots, and  
11 thus becomes gene specific, it will affect our estimates of protein synthesis  $\phi$ . This  
12 might be the case at recombination hotspots. Recombination, however, is very low  
13 in the introgressed region (discussed below) [15, 18]. The low recombination rate  
14 also indicates that the GC content had to be high before the introgression occurred.

15 The estimates of mutation and selection bias parameters,  $\Delta M$  and  $\Delta \eta$ , are ob-  
16 tained under an equilibrium assumption. Given that the introgression is still adapt-  
17 ing to its new environment, this assumption is clearly violated. However, the adap-  
18 tation of the exogenous genes progresses very slowly as a quasi-static process as  
19 shown in this work as well as [16]. Therefore, the genome can be assumed to main-  
20 tain an internal equilibrium at any given time. We see empirical evidence for this  
21 behavior in our ability to predict gene expression and to correctly identify the low  
22 expression genes (Figure 1b).

23 Despite the violation of the equilibrium assumption, the mutation and selection  
24 bias parameters  $\Delta M$  and  $\Delta \eta$  of the introgressed exogenous genes contain informa-  
25 tion, albeit decaying, about its previous cellular environment. We selected the top  
26 ten lineages with the highest similarity in  $\Delta M$  to see if our parameters estimates  
27 would allow us to identify a potential source lineage. The synteny relationship of  
28 these lineages with the exogenous genes was calculated as a point of comparison as  
29 it provides orthogonal information to our parameter estimates. Synteny with the  
30 exogenous genes is limited to the Saccharomycetaceae clade, excluding all of the  
31 potential source lineages identified using codon usage but *E. gossypii* (Table 2). In-  
32 terestingly, this also showed that similarity in codon usage does not correlate with  
33 phylogenetic distance.

1 Previous work indicated that the donor lineage of the exogenous genes has to be  
2 a, potentially unknown, Lachancea lineage [15–18]. These previous results, however,  
3 are based on species rather than gene trees, ignoring the differential adaptation rate  
4 to their novel cellular environment between genes or do not consider lineages outside  
5 of the Lachancea clade. Considering the similarity in selection bias (Figure 2b) and  
6 our calculation of selection on the exogenous genes (Figure 5b), both of which  
7 are free of any assumption about the origin of the exogenous genes, a species tree  
8 estimated from the exogenous genes will be biased towards the Lachancea clade.  
9 Estimating individual gene trees rather than relying on a species tree provided  
10 further evidence that the exogenous genes could originate from a lineage that does  
11 not belong to the Lachancea clade. As we highlighted in this study, relatively small  
12 sets of genes with a signal of a foreign cellular environment can significantly bias  
13 the outcome of a study. The same holds true for phylogenetic inferences [31], and as  
14 we showed the signal of the original endogenous cellular environment that shaped  
15 CUB is at different stages of decay in high and low expression genes (Figure S8).  
16 In summary, our work does not dispute an unknown Lachancea as possible origin,  
17 but provides an alternative hypothesis based on the codon usage of the exogenous  
18 genes, phylogenetic analysis, and synteny.

19 In terms of understanding the spread of the introgression, we calculated the ex-  
20 pected selective cost of codon mismatch between the *L. kluyveri* and *E. gossypii*  
21 lineages. Under our working hypothesis, the majority of the introgressed would have  
22 imposed a selective cost due to codon mismatch. Nevertheless, ~30% of low expres-  
23 sion exogenous genes ( $\phi < 1$ ) appeared to be exapted at the time of the introgres-  
24 sion. This exaptation is due to the mutation bias in the endogenous genes matching  
25 the selection bias in the exogenous genes for GC ending codons. Our estimate of  
26 the selective cost of codon mismatch on the order of  $-0.0008$ . While this selective  
27 cost may not seem very large, assuming *L. kluyveri* had a large  $N_e$ , the fixation  
28 probability of the introgression is the astronomically small value of  $\approx 10^{-6952} \approx 0$ .  
29 While this estimate heavily depends on the working hypothesis that the exogenous  
30 genes originated from the *E. gossypii* lineage, we can also calculate the hypothetical  
31 fixation probability if the current exogenous genes would introgress into *L. kluyveri*.  
32 Our estimate of the current selective cost of the mismatch of codon usage is on the  
33 order of  $-0.0003$ . The fixation probability of the current exogenous genes would

1 still be astronomically small  $\approx 10^{-2609} \approx 0$ . These results are in accordance with  
2 previous work, highlighting the necessity of codon usage compatibility between endo-  
3 genous and transferred exogenous genes [32, 33]. Thus, the basic scenario of an  
4 introgression between two yeast species with large  $N_e$  and where the introgression  
5 solely imposes a selective cost due to codon mismatch is clearly too simplistic.

6 One or more loci with a combined selective advantage on the order of 0.0008  
7 or greater would have made the introgression change from disadvantageous to ef-  
8 fectively neutral or advantageous. While this scenario seems plausible, it raises  
9 the question as to why recombination events did not limit the introgression to  
10 only the adaptive loci. A potential answer is the low recombination rate between  
11 the endogenous and exogenous regions [15, 18]. Estimates of the recombination  
12 rate as measured by crossovers (COs) for *L. kluyveri* are almost four times lower  
13 than for *S. cerevisiae* and about half that of *Schizosaccharomyces pombe* ( $\approx 1.6$   
14 COs/Mb/meiosis,  $\approx 6$  COs/Mb/meiosis,  $\approx 3$  COs/Mb/meiosis) with no observed  
15 crossovers in the introgressed region [18], and no observed transposable elements  
16 [15]. This is presumably due to the dissimilarity in GC content and/or a lower than  
17 average sequence homology between the exogenous region and the one it replaced.  
18 A population bottleneck reducing the  $N_e$  of the *L. kluyveri* lineage around the time  
19 of the introgression could also help explain the spread of the introgression. Compati-  
20 ble with these explanation is the possibility of several advantageous loci distributed  
21 across the exogenous region drove a rapid selective sweep and/or the population  
22 through a bottleneck speciation process.

23 Assuming *E. gossypii* as potential source lineage of the exogenous region, we  
24 illustrated how information on codon usage can be used to infer the time since  
25 the introgression occurred using our estimates of mutation bias  $\Delta M$ . The  $\Delta M$   
26 estimates are well suited for this task as they are free of the influence of selection  
27 and unbiased by  $N_e$  and other scaling terms, which is in contrast to our estimates of  
28  $\Delta\eta$  [10]. Our estimated age of the introgression of  $6.2 \pm 1.2 \times 10^8$  generations is  $\sim 10$   
29 times longer than a previous minimum estimate by [16] of  $5.6 \times 10^7$  generations,  
30 which was based on the effective population recombination rate and the population  
31 mutation parameter [34]. Furthermore, these estimates assume that the current *E.*  
32 *gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the  
33 time of the introgression. Thus, if the ancestral mutation environments were more

1 similar (dissimilar) at the time of the introgression then our result is an overestimate  
2 (underestimate).

3 Further, the presented work provides a template to explore the evolution of codon  
4 usage. This applies not only to species who experienced an introgression but is more  
5 generally applicable to any species.

## 6 Conclusion

7 Overall, our results show the usefulness of the separation of mutation bias and  
8 selection bias and the importance of recognizing the presence of multiple cellular  
9 environments in the study of codon usage. We also illustrate how a mechanistic  
10 model like ROC SEMPPR and the quantitative estimates it provides can be used for  
11 more sophisticated hypothesis testing in the future. In contrast to other approaches  
12 used to study codon usage like CAI [5] or tAI [19], ROC SEMPPR incorporates the  
13 effects of mutation bias and amino acid composition explicitly [20]. We highlight  
14 potential issues when estimating codon preferences, as estimates can be biased by  
15 the signature of a second, historical cellular environment. In addition, we show  
16 how quantitative estimates of mutation bias and selection relative to drift can be  
17 obtained from codon data and used to infer the fitness cost of an introgression as  
18 well as its history and potential future.

## 20 Materials and Methods

### 21 Separating Endogenous and Exogenous Genes

22 A GC-rich region was identified by [15] in the *L. kluyveri* genome extending from  
23 position 1 to 989,693 of chromosome C. This region was later identified as an  
24 introgression by [16]. We obtained the *L. kluyveri* genome from SGD Project  
25 <http://www.yeastgenome.org/download-data/> (on 09-27-2014) and the annotation  
26 for *L. kluyveri* NRRL Y-12651 (assembly ASM14922v1) from NCBI (on 12-09-  
27 2014). We assigned 457 genes located on chromosome C with a location within the  
28 ~ 1 Mb window to the exogenous gene set. All other 4864 genes of the *L. kluyveri*  
29 genome were assigned to the exogenous genes.

### 31 Model Fitting with ROC SEMPPR

32 ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) [22] and R (3.4.1)  
33 [35]. ROC SEMPPR was run from 10 different starting values for at least 250,000

1 iterations and thinned to keep every 50th iteration. After manual inspection to  
2 verify that the MCMC had converged, parameter posterior means, log posterior  
3 probability and log likelihood were estimated from the last 500 samples (last 10%  
4 of samples).

5

## 6 Model selection

7 The marginal likelihood of the combined and separated model fits was calculated  
8 using a generalized harmonic mean estimator [36]. A variance scaling of 1.1 was  
9 used to scale the important density of the estimator. Using the estimated marginal  
10 likelihoods, we calculated the Bayes factor to assess model performance. Increases  
11 in the variance scaling increase the estimated Bayes factor, therefore we report a  
12 conservative Bayes factor bases on a small variance scaling S9.

13

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

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

18

$$19 \Delta M_i = \Delta M_{i,1} - \overline{\Delta M_i} \quad (1) \quad 19$$

20

$$21 \Delta \eta_i = \Delta \eta_{i,1} - \overline{\Delta \eta_i} \quad (2) \quad 21$$

22 Comparison of codon specific parameters ( $\Delta M$  and  $\Delta \eta = 2N_e q(\eta_i - \eta_j)$ ) was per-  
23 formed using the function lmodel2 in the R package lmodel2 (1.7.3) [37] and R  
24 version 3.4.1 [35]. The parameter  $\Delta \eta$  can be interpreted as the difference in fitness  
25 between codon  $i$  and  $j$  for the average gene with  $\phi = 1$  scaled by the effective pop-  
26 ulation size  $N_e$ , and the selective cost of an ATP  $q$  [4, 10]. Type II regression was  
27 performed with re-centered parameter estimates, accounting for noise in dependent  
28 and independent variable [24].

29

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

1 Phylogenetic Analysis

2 Using the dataset from [21], we first identified 129 alignments for exogenous genes  
 3 that further contained homologous genes for *E. gossypii*, and at least one other  
 4 Lachancea species. We excluded all species from the alignments that do not belong  
 5 to the Saccharomycetaceae clade. IQTree [25] was used to identify the best fit-  
 6 ting model for each gene and to estimate the individual gene trees. Each gene tree  
 7 was rooted using either *Saccharomyces cerevisiae*, *Saccharomyces uvarum*, *Saccha-*  
 8 *romyces eubayanus* as outgroup. We calculated the most recent common ancestor  
 9 (MRCA) of *L. kluyveri* and *E. gossypii* as well as the MRCA of *L. kluyveri* and the  
 10 remaining Lachancea. The distance between the MRCA and the root was used to  
 11 asses which pairs (*L. kluyveri* and *E. gossypii*, or *L. kluyveri* and other Lachancea)  
 12 have a more recent common ancestor.

13

14 Synteny Comparison

15 We obtained complete genome sequences for all 10 candidate lineages (Table 2)  
 16 from NCBI (on: 02-05-2017). Genomes were aligned and checked for synteny using  
 17 SyMAP (4.2) with default settings [38, 39]. We assess synteny as percentage coverage  
 18 of the exogenous gene region.

19

20 Estimating Age of Introgression

21 We modeled the change in codon frequency over time using an exponential model  
 22 for all two codon amino acids. While our approach is equivalent to [40], we want  
 23 to explicitly state the relationship between the change in codon frequency  $c_1$  as a  
 24 function of mutation bias  $\Delta M$  as

$$25 \quad \frac{dc_1}{dt} = -\mu_{1,2}c_1 - \mu_{2,1}(1 - c_1) \quad (3)$$

26 where  $\mu_{i,j}$  is the rate at which codon  $i$  mutates to codon  $j$  and  $c_1$  is the fre-  
 27 quency of the reference codon. Initial codon frequencies  $c_1(0)$  for each codon  
 28 family were taken from our mutation parameter estimates for *E. gossypii* where  
 29  $c_1(0) = \exp[\Delta M_{gos}] / (1 + \exp[\Delta M_{gos}])$ . Our estimates of  $\Delta M_{endo}$  can be used to  
 30 calculate the steady state of equation 3 were  $\frac{dc_1}{dt} = 0$  to obtain the equality

$$31 \quad \frac{\mu_{2,1}}{\mu_{1,2} + \mu_{2,1}} = \frac{1}{1 + \exp[\Delta M_{endo}]} \quad (4)$$

32

33

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

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

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

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

6 Equation 5 was solved with a mutation rate  $\mu_{2,1}$  of  $3.8 \times 10^{-10}$  per nucleotide per  
 7 generation [41]. Current codon frequencies for each codon family were taken from  
 8 our estimates of  $\Delta M$  from the exogenous genes. Mathematica (11.3) [42] was used  
 9 to calculate the time  $t_{\text{intro}}$  it takes for the initial codon frequencies  $c_1(0)$  for each  
 10 codon family to equal the current exogenous codon frequencies. The same equation  
 11 was used to determine the time  $t_{\text{decay}}$  at which the signal of the exogenous cellular  
 12 environment has decayed to within 1% of the endogenous environment.

#### 14 Estimating Selection against Codon Mismatch

15 In order to estimate the selection against codon mismatch, we had to make three  
 16 key assumptions. First, we assumed that the current exogenous amino acid sequence  
 17 of a gene is representative of its ancestral state and the replaced endogenous gene  
 18 it replaced. Second, we assume that the currently observed cellular environment of  
 19 *E. gossypii* reflects the cellular environment that the exogenous genes experienced  
 20 before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy  
 21 of selection between the cellular environments due to differences in either effective  
 22 population size  $N_e$  or the selective cost of an ATP  $q$  of the source lineage and *L.*  
 23 *kluyveri* can be expressed as a scaling constant and that protein synthesis rate  $\phi$   
 24 has not changed between the replaced endogenous and the introgressed exogenous  
 25 genes. Using estimates for  $N_e = 1.36 \times 10^7$  [27] for *Saccharomyces paradoxus* we  
 26 scale our estimates of  $\Delta\eta$  which explicitly contains the effective population size  $N_e$   
 27 [10] and define  $\Delta\eta' = \frac{\Delta\eta}{N_e}$ .

28 All of our genome parameter estimations are scaled by lineage specific effects such  
 29 as  $N_e$ , the average, absolute gene expression level, and/or the proportionate fitness  
 30 value of an ATP. In order to account for these genome specific differences in scaling,  
 31 we scale the difference in the efficacy of selection on codon usage between the donor  
 32 lineage and *L. kluyveri* using a linear scaling factor  $\kappa$ . As  $\Delta\eta$  is defined as  $\Delta\eta =$   
 33  $2N_e q(\eta_i - \eta_j)$ , we cannot distinguish if  $\kappa$  is a scaling on protein synthesis rate  $\phi$ ,

1 effective population size  $N_e$ , or the selective cost of an ATP  $q$  [4, 10]. We calculated  
 2 the selection against each genes codon mismatch assuming additive fitness effects  
 3 as

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

7 where  $s_g$  is the overall strength of selection for translational efficiency on gene,  $g$   
 8 in the exogenous gene set,  $\kappa$  is a constant, scaling the efficacy of selection between  
 9 the endogenous and exogenous cellular environments,  $L_g$  is length of the protein in  
 10 codons,  $\phi_g$  is the estimated protein synthesis rate of the gene in the endogenous  
 11 environment, and  $\Delta \eta'_i$  is the  $\Delta \eta'$  for the codon at position  $i$ . As stated previously,  
 12 our  $\Delta \eta$  are relative to the mean of the codon family. We find that the selection  
 13 against the introgressed genes is minimized at  $\kappa \sim 5$  (Figure S7b). Thus, we expect  
 14 a five fold difference in the efficacy of selection between *L. kluyveri* and *E. gossypii*,  
 15 due to differences in either protein synthesis rate  $\phi$ , effective population size  $N_e$ ,  
 16 and/or the selective cost of an ATP  $q$ . Therefore, we set  $\kappa = 1$  if we calculate the  $s_g$   
 17 for the endogenous and the current exogenous genes, and  $\kappa = 5$  for  $s_g$  for selection  
 18 calculations at the time of introgression.

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

$$\begin{aligned} E[n_{g,i}] &= P(c_i | \Delta M, \Delta \eta, \phi) \times m_{a_i} \\ &= \frac{\exp[-\Delta M_i - \Delta \eta_i \phi_g]}{\sum_j^C \exp[-\Delta M_j - \Delta \eta_j \phi_g]} \times m_{a_i} \end{aligned}$$

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

$$s_g = \sum_{i=1}^C -\kappa \phi_g \Delta \eta'_i E[n_{g,i}] \quad (7)$$

1 We report the selection due to mismatched codon usage of the introgression as  
2  $\Delta s_g = s_{\text{intro},g} - s_{\text{endo},g}$  where  $s_{\text{intro},g}$  is the selection against an introgressed gene  $g$   
3 either at the time of the introgression or presently.  
4

5 **Randomizing genes**

6 We randomized the codon content of the endogenous and exogenous genes while  
7 conserving the di-nucleotide distribution and GC content using the randomization  
8 algorithm from SPARCS [43]. We used the default settings of the randomization  
9 algorithm. The resulting gene sets were analyzed using the same scheme as described  
10 above.  
11

12 **Acknowledgments**

13 The authors would like to thank Alexander Cope for helpful criticisms and suggestions for this work.  
14

15 **Availability of data and materials**

16 Parameter estimates generated during this study are available from the corresponding author. All remaining data  
17 generated during this study are included in this published article as figures, tables.  
18

19 **Authors' contributions**

20 CL and MAG initiated the study. CL collected and analyzed the data and wrote the manuscript. MAG and BCO  
21 edited the manuscript. CL, MAG, BCO, and RZ contributed to the data analysis and acquiring of funding. All  
22 Authors approved the final manuscript.  
23

24 **Funding**

25 This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ), MCB-1546402 (A. Von Arnim and  
26 MAG), and DEB-1355033 (BCO, MAG, and RZ) with additional support from Department of Ecology &  
27 Evolutionary Biology (EEB) at the University of Tennessee Knoxville (UTK) and the National Institute for  
28 Mathematical and Biological Synthesis (NIMBioS), an Institute sponsored by the National Science Foundation  
29 through NSF Award DBI-1300426. CL received support as a Graduate Student Fellow from NIMBioS with  
30 additional support from Departments of Mathematics and EEB at UTK.  
31

32 **Ethics approval and consent to participate**

33 Not applicable  
34

35 **Consent for publication**

36 Not applicable  
37

38 **Competing interests**

39 The authors declare that they have no competing interests.  
40

41 **Author details**

42 <sup>1</sup>Department of Ecology & Evolutionary Biology, University of Tennessee, 37996, Knoxville, TN, USA. <sup>2</sup>National  
43 Institute for Mathematical and Biological Synthesis, 37996, Knoxville, TN, USA. <sup>3</sup>Max-Planck Institute of  
44 Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, 01307, Dresden, Germany. <sup>4</sup>Department of Business  
45 Analytics and Statistics, University of Tennessee, 37996, Knoxville, TN, USA.  
46

47 **References**

- 48 1. Gouy, M., Gautier, C.: Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Research* **10**,  
49 7055–7074 (1982)  
50

1. Ikemura, T.: Codon usage and tRNA content in unicellular and multicellular organisms. *Molecular Biology and Evolution* **2**, 13–34 (1985)

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

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

4. Sharp, P.M., Li, W.H.: The codon adaptation index - a measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Research* **15**, 1281–1295 (1987)

5. Wright, F.: The 'effective number of codons' used in a gene. *Genet* **87**, 23–29 (1990)

6. M, S.P., Stenico, M., Peden, J.F., Lloyd, A.T.: Codon usage: mutational bias, translational selection, or both? *Biochem Soc Trans.* **21**(4), 835–841 (1993)

7. Shah, P., Gilchrist, M.A.: Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences U.S.A* **108**(25), 10231–10236 (2011)

8. Wallace, E.W., Airoldi, E.M., Drummond, D.A.: Estimating selection on synonymous codon usage from noisy experimental data. *Molecular Biology and Evolution* **30**, 1438–1453 (2013)

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

10. Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., Danchin, A.: Evidence for horizontal gene transfer in *Escherichia coli* speciation. *Journal of Molecular Biology* **222**(4), 851–856 (1991)

11. Lawrence, J.G., Ochman, H.: Amelioration of bacterial genomes: Rates of change and exchange. *Journal of Molecular Biology* **44**, 383–397 (1997)

12. Marcet-Houben, M., Gabaldón, T.: Beyond the whole-genome duplication: Phylogenetic evidence for an ancient interspecies hybridization in the baker's yeast lineage. *PLoS Biology* **13**(8), 1002220 (2015)

13. Beimforde, C., Feldberg, K., Nylander, S., Rikkinen, J., Tuovila, H., Dörfelt, H., Gube, M., Jackson, D.J., Reitner, J., Seyfullah, L.J., Schmidt, A.R.: Estimating the phanerozoic history of the ascomycota lineages: combining fossil and molecular data. *Mol. Phylogenet. Evol.* **78**, 386–398 (2014)

14. Payen, C., Fischer, G., Marck, C., Proux, C., Sherman, D.J., Coppée, J.-Y., Johnston, M., Dujon, B., Neuvéglise, C.: Unusual composition of a yeast chromosome arm is associated with its delayed replication. *Genome Research* **19**(10), 1710–1721 (2009)

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

16. Vakirlis, N., Sarilar, V., Drillon, G., Fleiss, A., Agier, N., Meyniel, J.-P., Blanpain, L., Carbone, A., Devillers, H., Dubois, K., Gillet-Markowska, A., Graziani, S., Huu-Vang, N., Poirel, M., Reisser, C., Schott, J., Schacherer, J., Lafontaine, I., Llorente, B., Neuvéglise, C., Fischer, G.: Reconstruction of ancestral chromosome architecture and gene repertoire reveals principles of genome evolution in a model yeast genus. *Genome research* **26**(7), 918–32 (2016)

17. Brion, C., Legrand, S., Peter, J., Caradec, C., Pflieger, D., Hou, J., Friedrich, A., Llorente, B., Schacherer, J.: Variation of the meiotic recombination landscape and properties over a broad evolutionary distance in yeasts. *PLoS Genetics* **13**(8), 1006917 (2017)

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

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

20. Shen, X.X., Opulente, D.A., Kominek, J., Zhou, X., Steenwyk, J.L., Buh, K.V., Haase, M.A.B., Wisecaver, J.H., Wang, M., Doering, D.T., Boudouris, J.T., Schneider, R.M., Langdon, Q.K., Ohkuma, M., Endoh, R., Takashima, M., Manabe, R., Čadež, N., Libkind, D., Rosa, C., DeVirgilio, J., Hulfachor, A.B., Groenewald, M., Kurtzman, C., Hittinger, C.T., Rokas, A.: Tempo and mode of genome evolution in the budding yeast subphylum. *Cell* **175**(6), 1533–154520 (2018)

21. Landerer, C., Cope, A., Zaretzki, R., Gilchrist, M.A.: AnaCoDa: analyzing codon data with bayesian mixture

- 1 models. *Bioinformatics* **34**(14), 2496–2498 (2018)
- 2 23. Tsankov, A.M., Thompson, D.A., Socha, A., Regev, A., Rando, O.J.: The role of nucleosome positioning in the  
3 evolution of gene regulation. *PLoS Biol* **8**(7), 1000414 (2010)
- 4 24. Sokal, R.R., Rohlf, F.J.: *Biometry - The principles and practice of statistics in biological*, pp. 547–555. W. H.  
5 Freeman, New York, NY (1981)
- 6 25. Nguyen, L.T., Schmidt, H.A., von Haeseler, A., Minh, B.Q.: Iq-tree: A fast and effective stochastic algorithm  
7 for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **32**(1), 268–274 (2015)
- 8 26. Sella, G., Hirsh, A.E.: The application of statistical physics to evolutionary biology. *Proceedings of the National  
9 Academy of Sciences of the United States of America* **102**, 9541–9546 (2005)
- 10 27. Wagner, A.: Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* **22**,  
11 1365–1374 (2005)
- 12 28. Nagylaki, T.: Evolution of a finite population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**,  
13 6278–6281 (1983)
- 14 29. Nagylaki, T.: Evolution of a large population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**,  
15 5941–5945 (1983)
- 16 30. Harrison, R.J., Charlesworth, B.: Biased gene conversion affects patterns of codon usage and amino acid usage  
17 in the *Saccharomyces sensu stricto* group of yeasts. *Molecular Biology and Evolution* **28**(1), 117–129 (2011)
- 18 31. Salichos, L., Rokas, A.: Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature*  
19 **497**, 327–331 (2013)
- 20 32. Medrano-Soto, A., Moreno-Hagelsieb, G., Vinuela, P., Christen, J.A., Collado-Vides, J.: Successful lateral  
21 transfer requires codon usage compatibility between foreign genes and recipient genomes. *Molecular Biology  
22 and Evolution* **21**(10), 1884–1894 (2004)
- 23 33. Tuller, T., Girshovich, Y., Sella, Y., Kreimer, A., Freilich, S., Kupiec, M., Gophna, U., Ruppin, E.: Association  
24 between translation efficiency and horizontal gene transfer within microbial communities. *Nucleic Acids  
25 Research* **39**(11), 4743–4755 (2011). doi:10.1093/nar/gkr054
- 26 34. Ruderfer, D.M., Pratt, S.C., Seidl, H.S., Kruglyak, L.: Population genomic analysis of outcrossing and  
27 recombination in yeast. *Nature Genetics* **38**(9), 1077–1081 (2006)
- 28 35. R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical  
29 Computing, Vienna, Austria (2013). R Foundation for Statistical Computing. <http://www.R-project.org/>
- 30 36. Gronau, Q.F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D.S., Forster, J.J.,  
31 Wagenmakers, E.J., Steingrover, H.: A tutorial on bridge sampling. *Journal of Mathematical Psychology* **81**,  
32 80–97 (2017)
- 33 37. Legendre, P.: Lmodel2: Model II Regression. (2018). R package version 1.7-3.  
34 <https://CRAN.R-project.org/package=lmodel2>
- 35 38. Soderlund, C., Nelson, W., Shoemaker, A., Paterson, A.: Symap A system for discovering and viewing syntenic  
36 regions of fpc maps. *Genome Research* **16**, 1159–1168 (2006)
- 37 39. Soderlund, C., Bomhoff, M., Nelson, W.: Symap v3.4: a turnkey synteny system with application to plant  
38 genomes. *Nucleic Acids Research* **39**(10), 68 (2011)
- 39 40. Marais, G., Charlesworth, B., Wright, S.I.: Recombination and base composition: the case of the highly  
40 self-fertilizing plant *Arabidopsis thaliana*. *Genome Biology* **5**, 45 (2004)
- 41 41. Lang, G.I., Murray, A.W.: Estimating the per-base-pair mutation rate in the yeast *Saccharomyces cerevisiae*.  
42 *Genetics* **178**(1), 67–82 (2008)
- 43 42. Wolfram Research Inc.: Mathematica 11. (2017). <http://www.wolfram.com>
- 44 43. Zhang, Y., Ponty, Y., Blanchette, M., Lécuyer, E., Waldspühl, J.: Sparcs: a web server to analyze  
45 (un)structured regions in coding RNA sequences. *Nucleic Acids Research* **41**, 480–485 (2013)
- 46 30
- 47 31
- 48 32
- 49 33

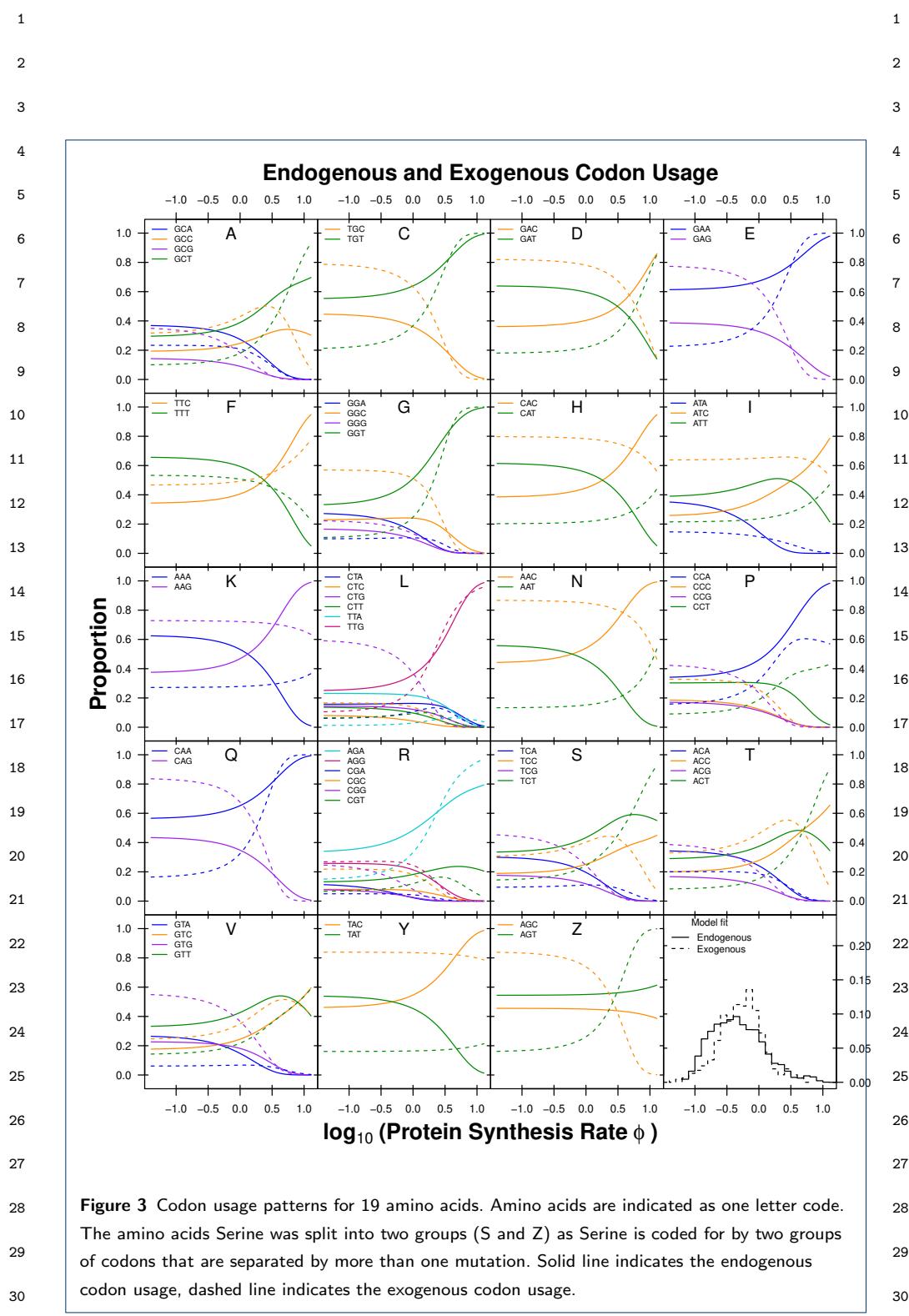
1      **Supplementary Material**

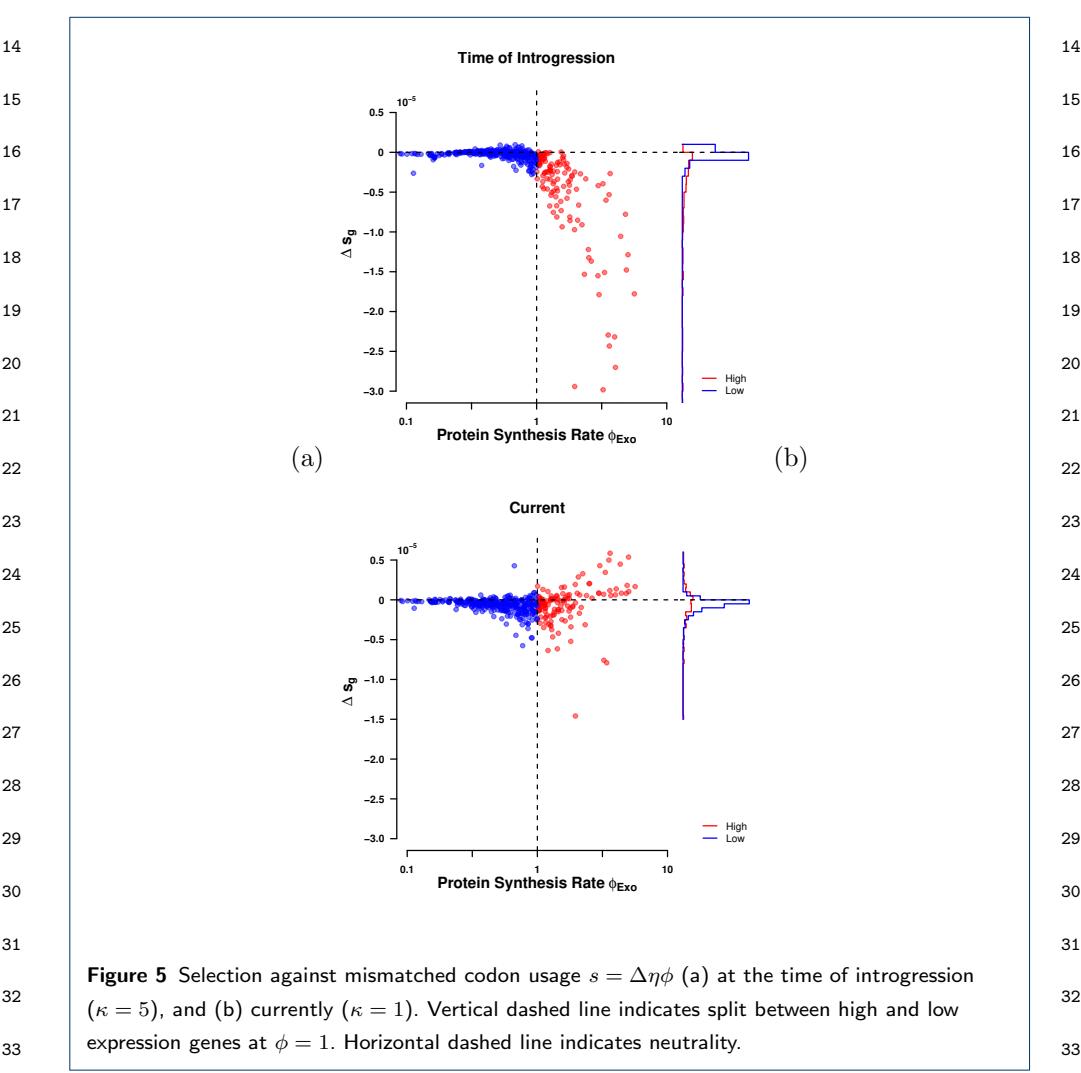
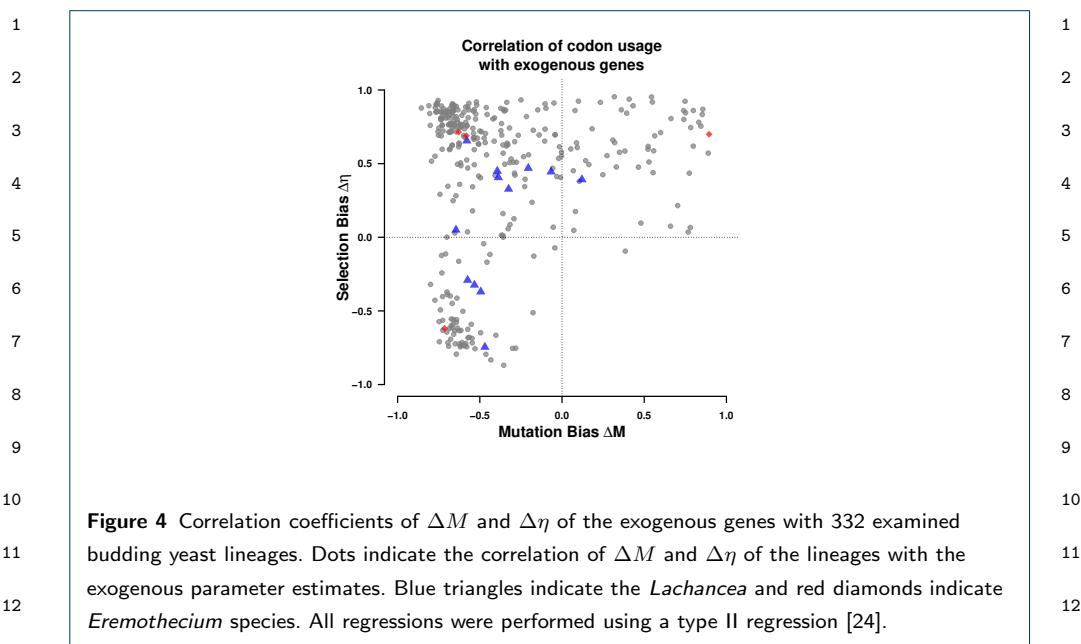
2      Supporting Materials for *Unlocking a signal of introgression from codons in Lachancea kluveri using a*  
*mutation-selection model* by Landerer et al..

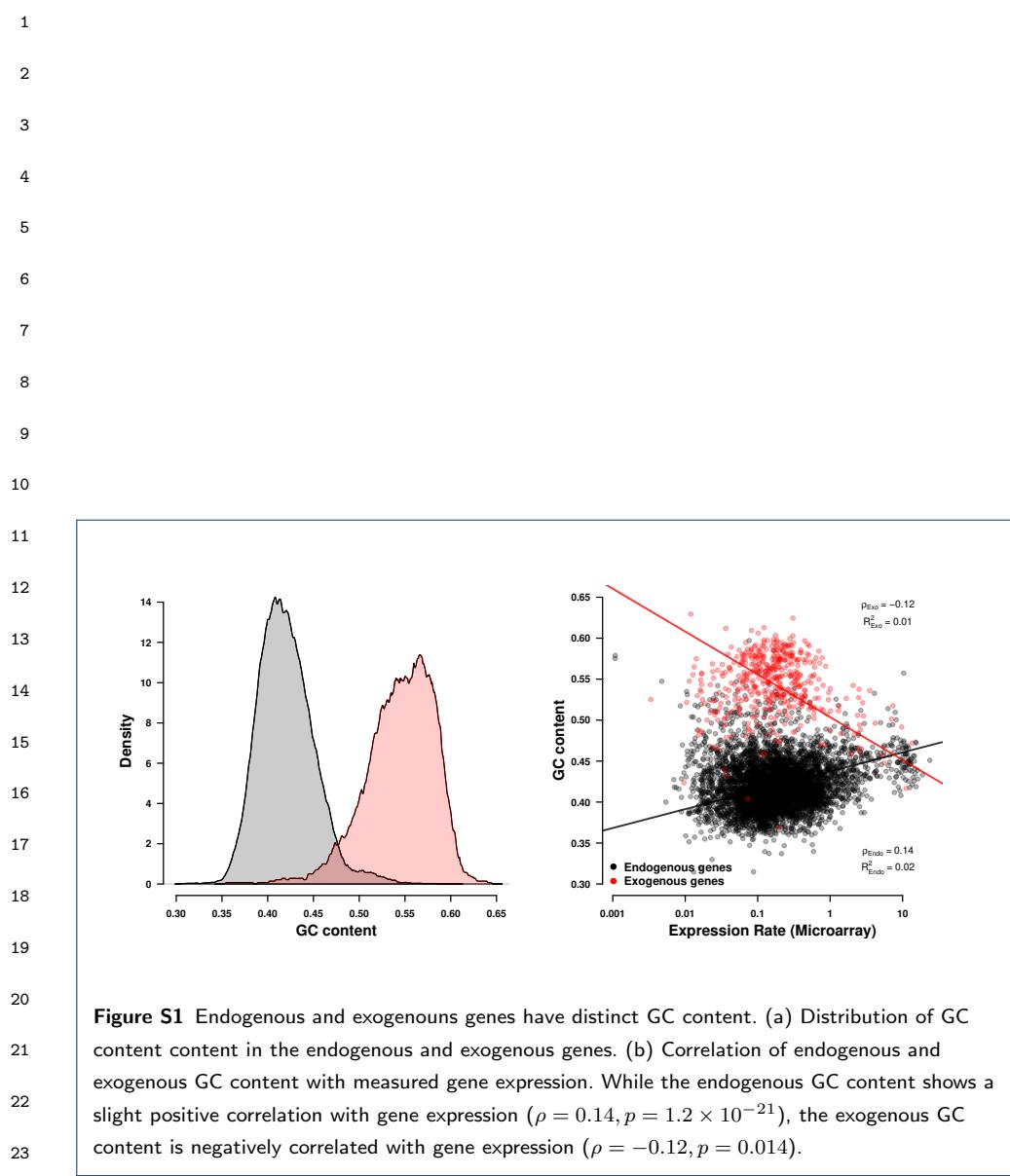
3      **Table S1** Synonymous mutation codon preference based on our estimates of  $\Delta M$ . Shown are the  
 4      most likely codon in low expression genes for each amino acid in: *E. gossypii*, in the endogenous and  
 5      exogenous genes of *L. kluyveri*, and in the combined *L. kluyveri*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

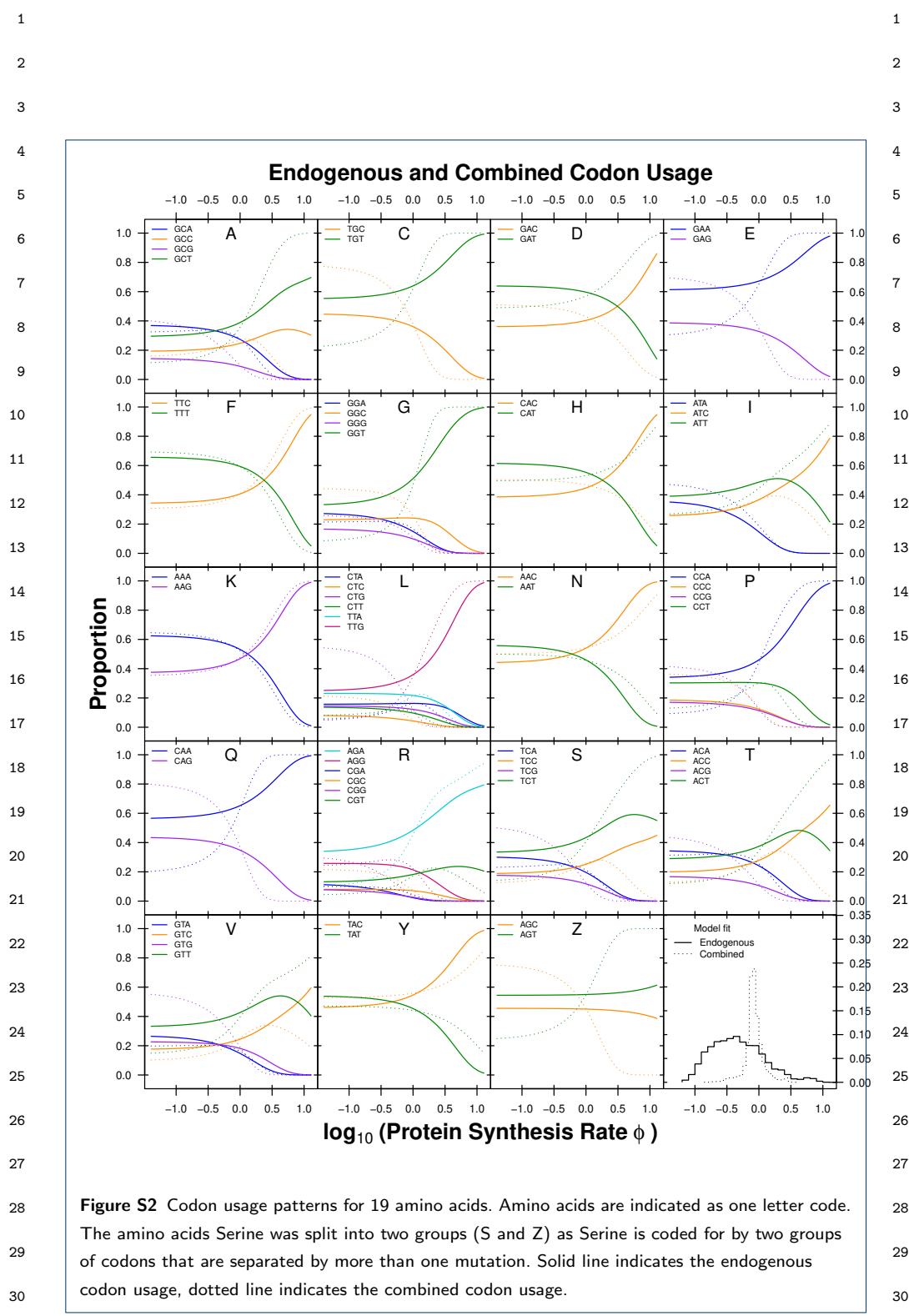
|    | Amino Acid         | <i>E. gossypii</i> | Endogenous | Exogenous | Combined |    |
|----|--------------------|--------------------|------------|-----------|----------|----|
| 7  | Ala A              | GCG                | GCA        | GCG       | GCG      | 7  |
| 8  | Cys C              | TGC                | TGT        | TGC       | TGC      | 8  |
| 9  | Asp D              | GAC                | GAT        | GAC       | GAC      | 9  |
| 10 | Glu E              | GAG                | GAA        | GAG       | GAG      | 10 |
| 11 | Phe F              | TTC                | TTT        | TTT       | TTT      | 11 |
| 12 | Gly G              | GGC                | GGT        | GGC       | GGC      | 12 |
| 13 | His H              | CAC                | CAT        | CAC       | CAC      | 13 |
| 14 | Ile I              | ATC                | ATT        | ATC       | ATA      | 14 |
| 15 | Lys K              | AAG                | AAA        | AAG       | AAA      | 15 |
| 16 | Leu L              | CTG                | TTG        | CTG       | CTG      | 16 |
| 17 | Asn N              | AAC                | AAT        | AAC       | AAT      | 17 |
| 18 | Pro P              | CCG                | CCA        | CCG       | CCG      | 18 |
| 19 | Gln Q              | CAG                | CAA        | CAG       | CAG      | 19 |
| 20 | Arg R              | CGC                | AGA        | AGG       | CGG      | 20 |
| 21 | Ser <sub>4</sub> S | TCG                | TCT        | TCG       | TCG      | 21 |
| 22 | Thr T              | ACG                | ACA        | ACG       | ACG      | 22 |
| 23 | Val V              | GTG                | GTT        | GTG       | GTG      | 23 |
| 24 | Tyr Y              | TAC                | TAT        | TAC       | TAC      | 24 |
| 25 | Ser <sub>2</sub> Z | AGC                | AGT        | AGC       | AGC      | 25 |
| 26 |                    |                    |            |           |          | 26 |
| 27 |                    |                    |            |           |          | 27 |
| 28 |                    |                    |            |           |          | 28 |
| 29 |                    |                    |            |           |          | 29 |
| 30 |                    |                    |            |           |          | 30 |
| 31 |                    |                    |            |           |          | 31 |
| 32 |                    |                    |            |           |          | 32 |
| 33 |                    |                    |            |           |          | 33 |

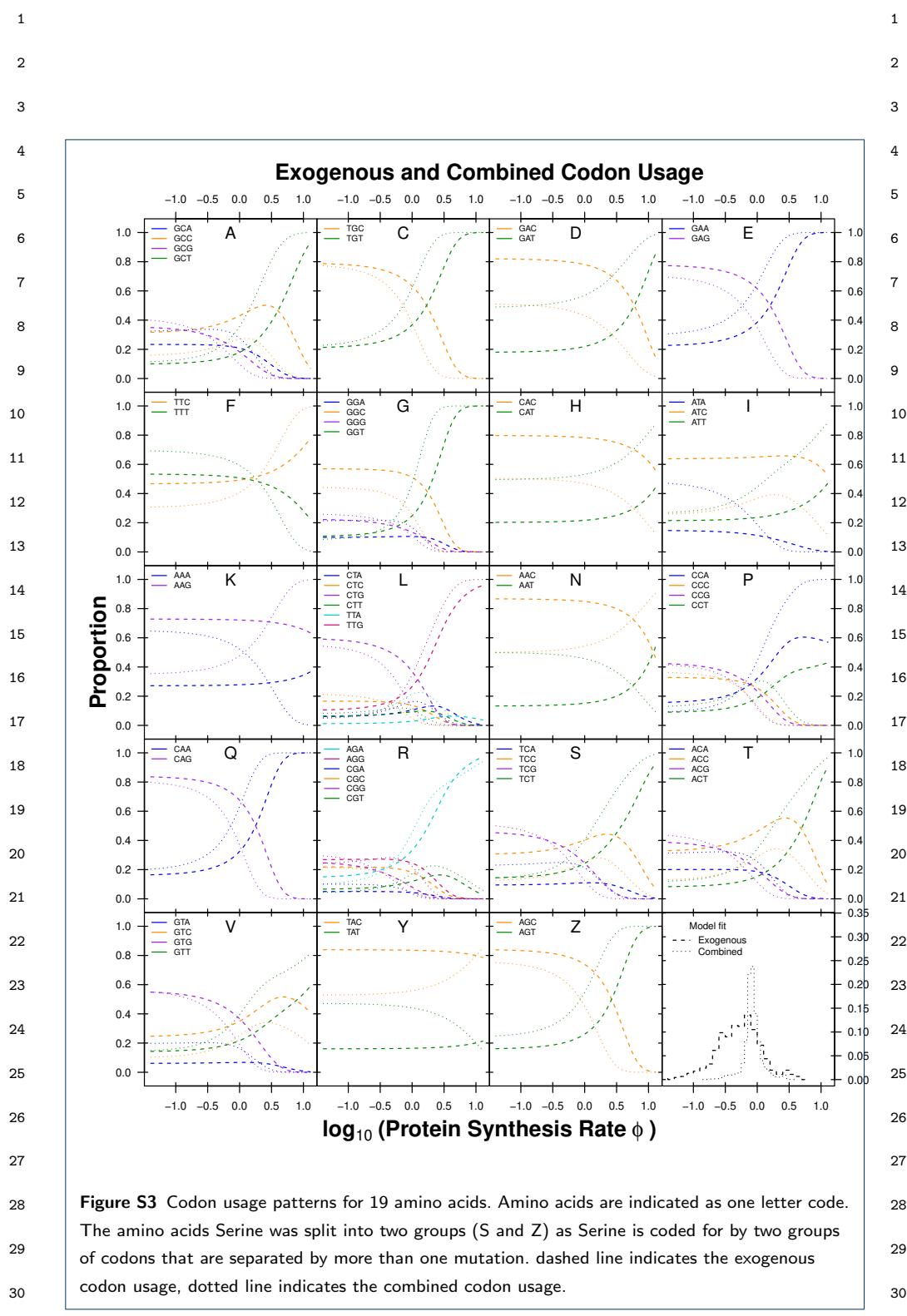
|    |  |                    |            |           |          |    |
|----|--|--------------------|------------|-----------|----------|----|
| 1  |  | 1                  |            |           |          |    |
| 2  |  | 2                  |            |           |          |    |
| 3  |  | 3                  |            |           |          |    |
| 4  |  | 4                  |            |           |          |    |
| 5  |  | 5                  |            |           |          |    |
| 6  |  | 6                  |            |           |          |    |
| 7  |  | 7                  |            |           |          |    |
| 8  |  | 8                  |            |           |          |    |
| 9  |  | 9                  |            |           |          |    |
| 10 | <b>Table S2</b> 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: <i>E. gossypii</i> , in the endogenous and exogenous genes of <i>L. kluyveri</i> , and in the combined <i>L. kluyveri</i> genome without accounting for the two cellular environments. | 10                 |            |           |          |    |
| 11 |  | 11                 |            |           |          |    |
| 12 |  | 12                 |            |           |          |    |
| 13 | Amino Acid   | <i>E. gossypii</i> | Endogenous | Exogenous | Combined | 13 |
| 14 | Ala A  | GCT                | GCT        | GCT       | GCT      |    |
| 15 | Cys C  | TGT                | TGT        | TGT       | TGT      |    |
| 16 | Asp D  | GAT                | GAC        | GAT       | GAT      |    |
| 17 | Glu E  | GAA                | GAA        | GAA       | GAA      |    |
| 18 | Phe F  | TTT                | TTC        | TTC       | TTC      |    |
| 19 | Gly G  | GGA                | GGT        | GGT       | GGT      |    |
| 20 | His H  | CAT                | CAC        | CAT       | CAT      |    |
| 21 | Ile I  | ATA                | ATC        | ATT       | ATT      |    |
| 22 | Lys K  | AAA                | AAG        | AAA       | AAG      |    |
| 23 | Leu L  | TTA                | TTG        | TTG       | TTG      |    |
| 24 | Asn N  | AAT                | AAC        | AAT       | AAC      |    |
| 25 | Pro P  | CCA                | CCA        | CCT       | CCA      |    |
| 26 | Gln Q  | CAA                | CAA        | CAA       | CAA      |    |
| 27 | Arg R  | AGA                | AGA        | AGA       | AGA      |    |
| 28 | Ser <sub>4</sub> S   | TCA                | TCC        | TCT       | TCT      |    |
| 29 | Thr T  | ACT                | ACC        | ACT       | ACT      |    |
| 30 | Val V  | GTT                | GTC        | GTT       | GTT      |    |
| 31 | Tyr Y  | TAT                | TAC        | TAT       | TAC      |    |
| 32 | Ser <sub>2</sub> Z   | AGT                | AGT        | AGT       | AGT      |    |
| 33 |  |                    |            |           |          | 33 |

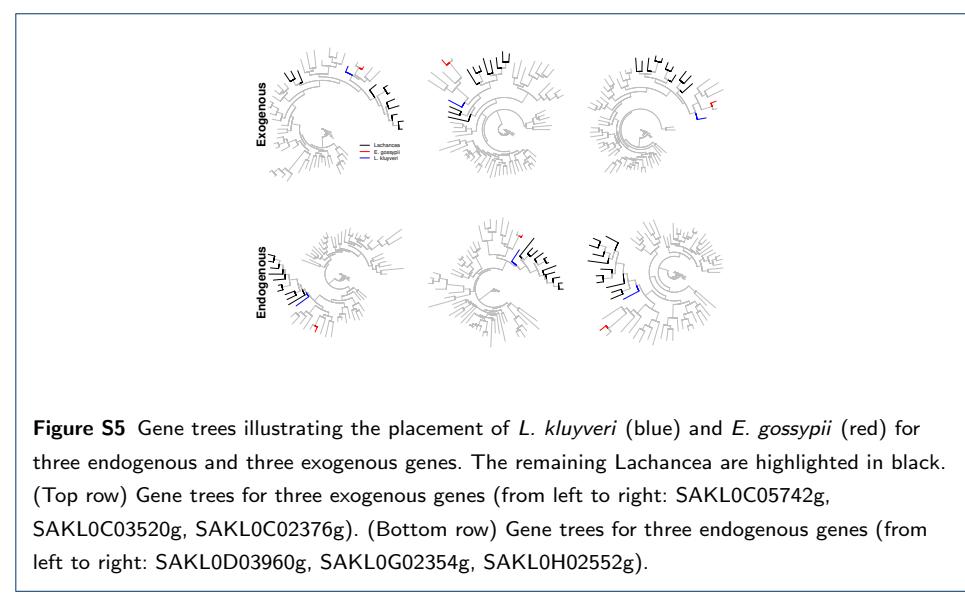
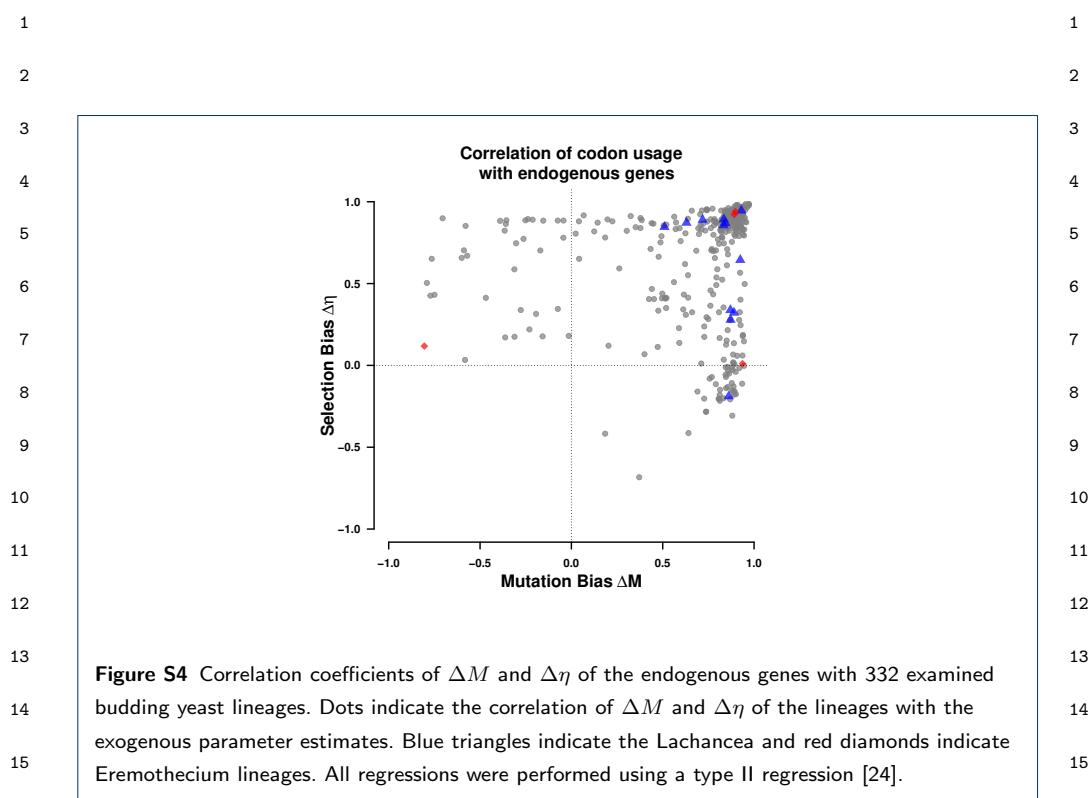


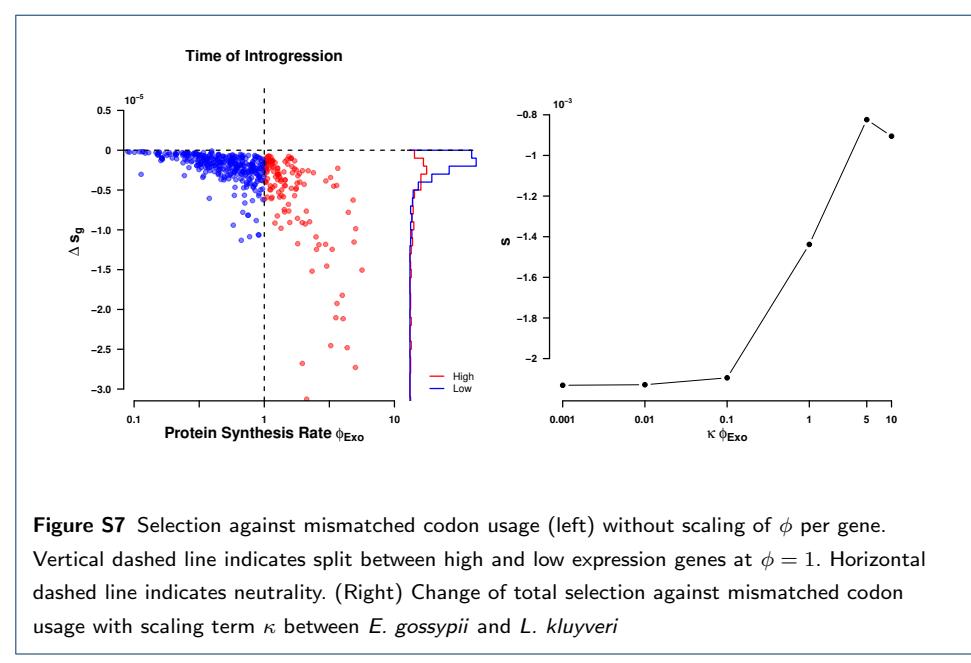
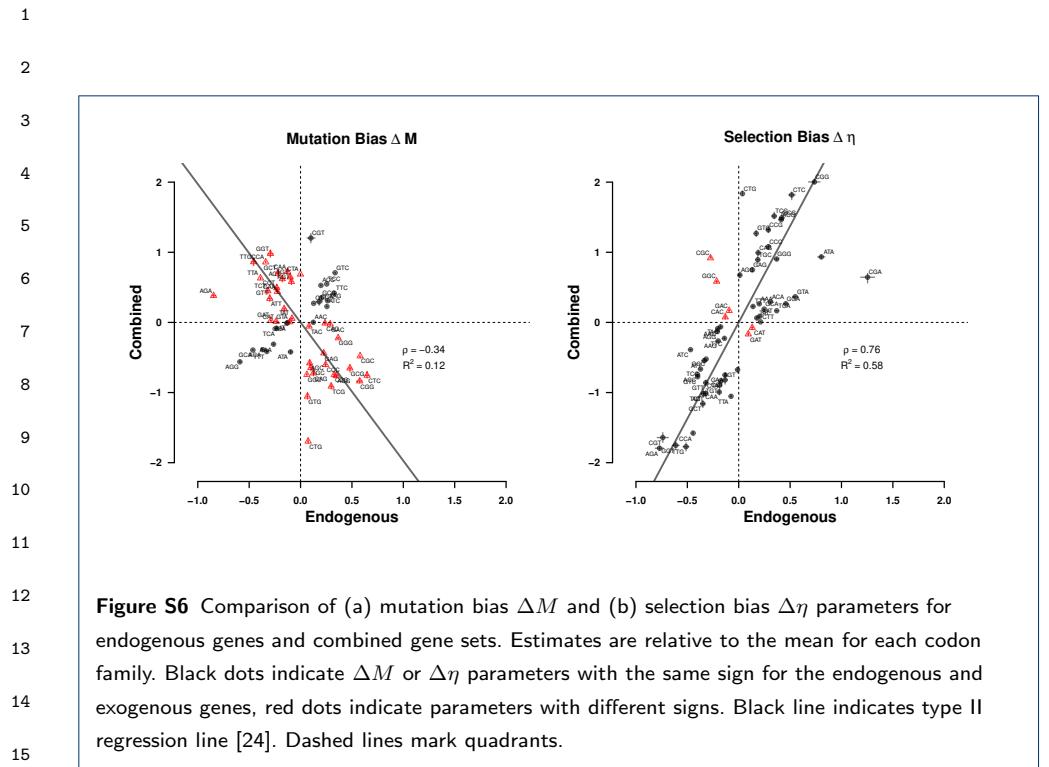


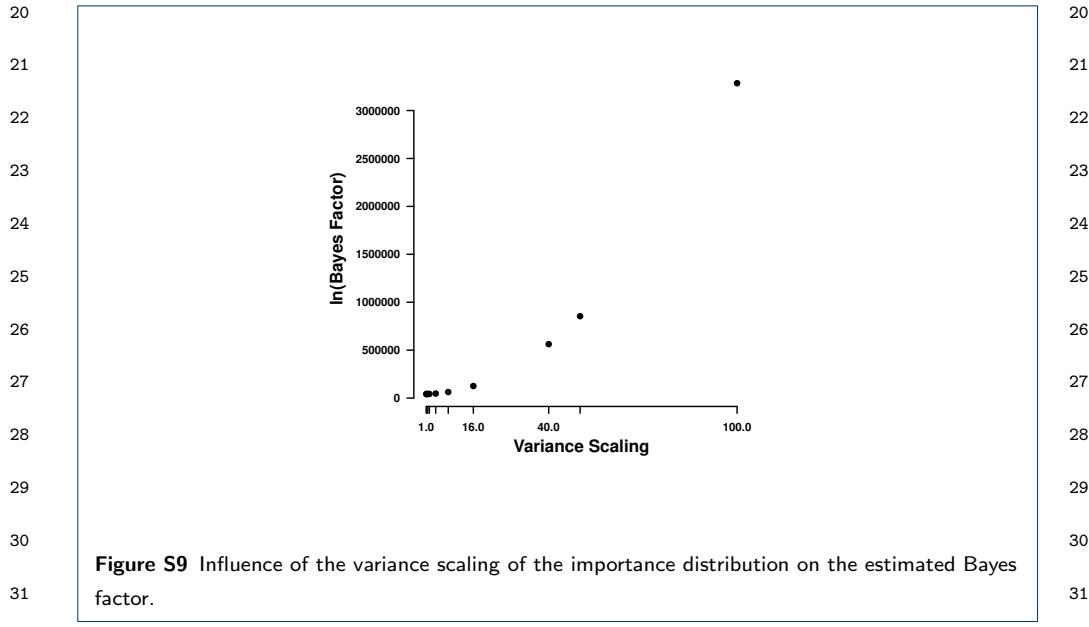
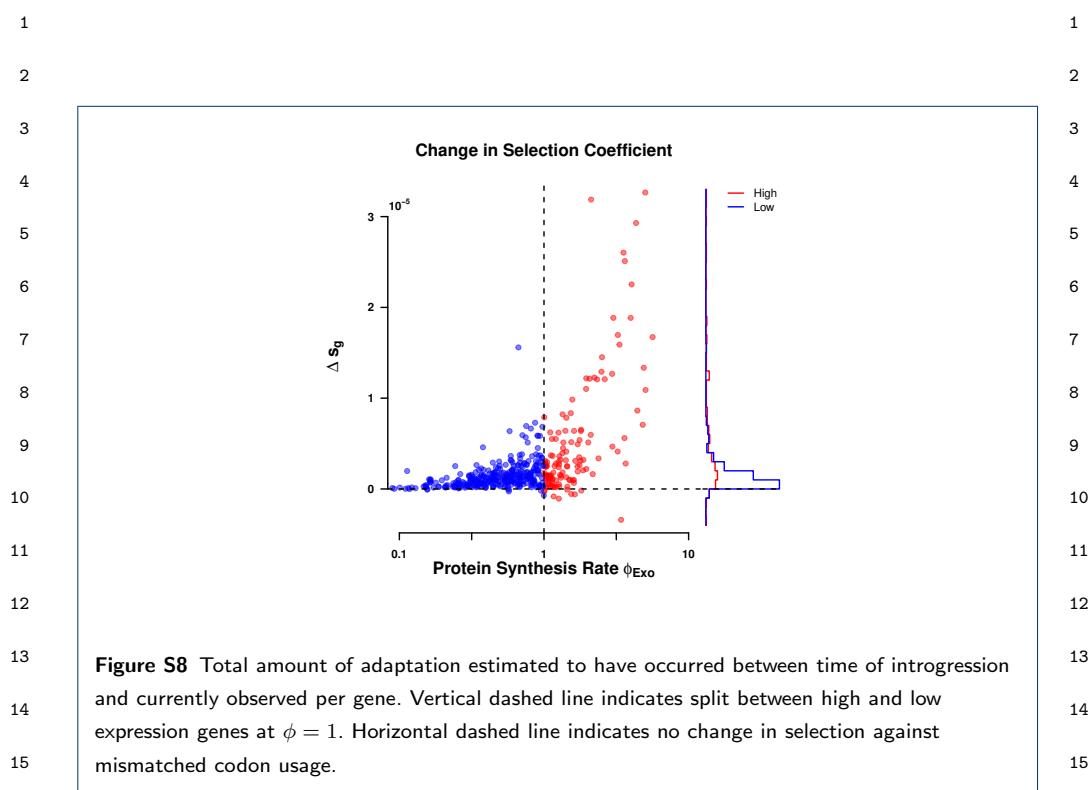


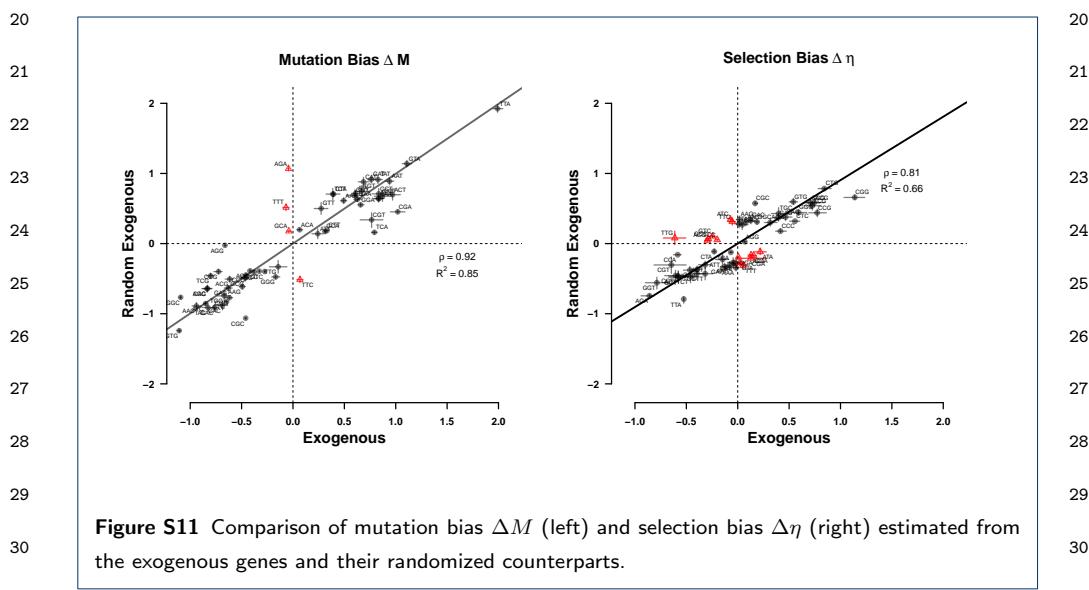
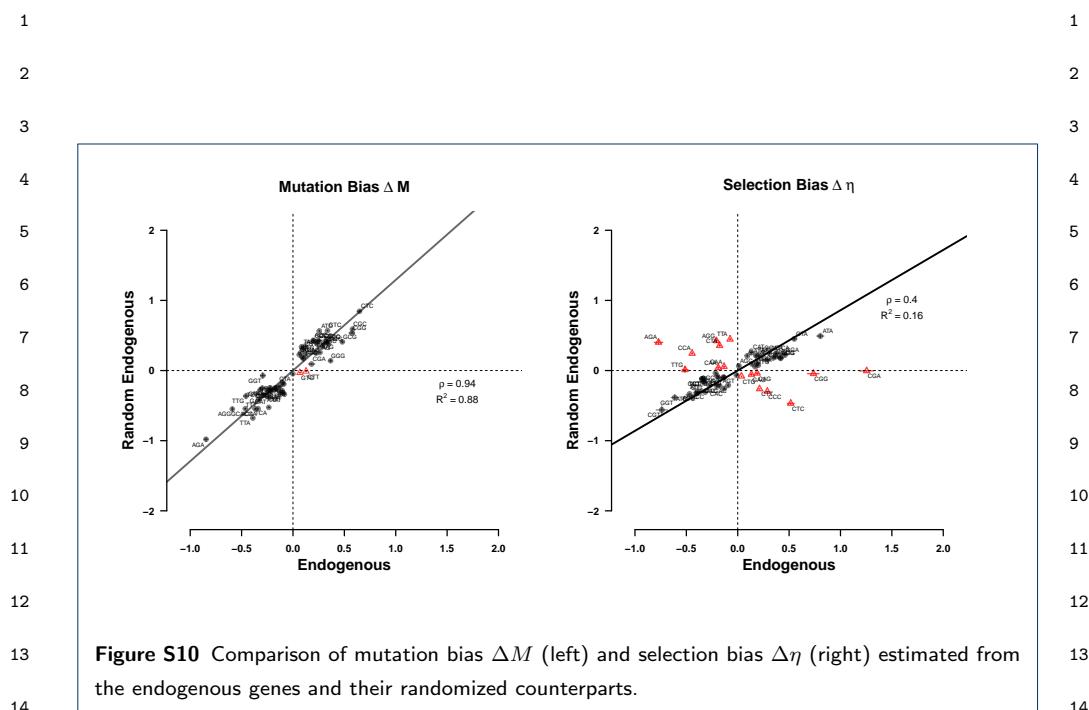


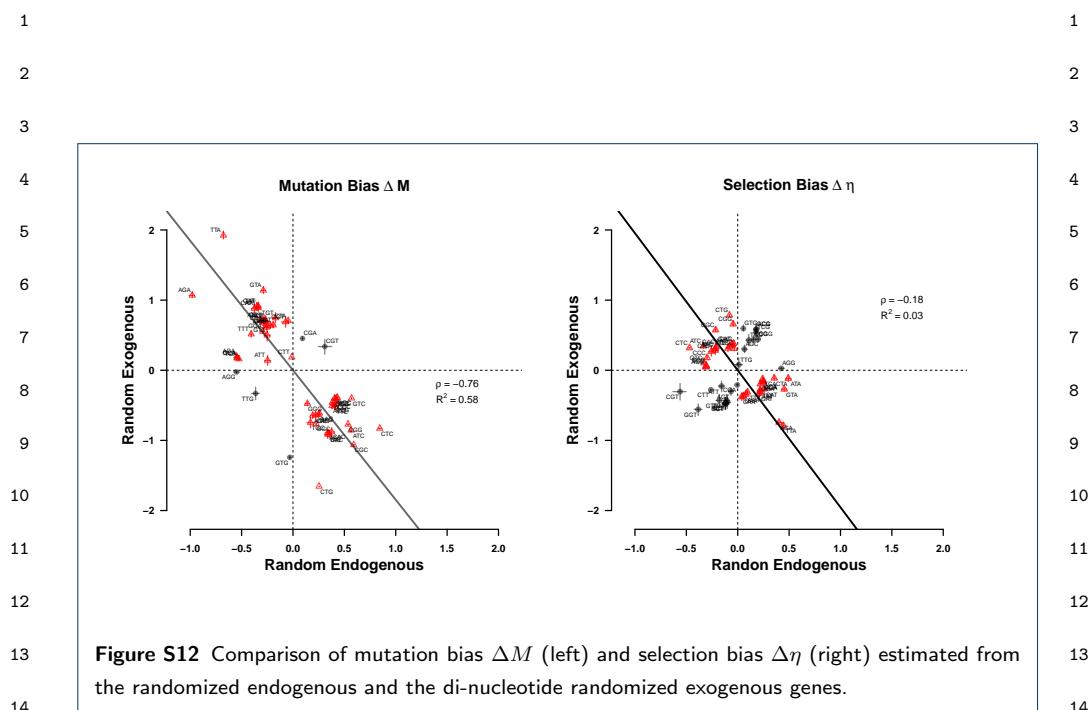




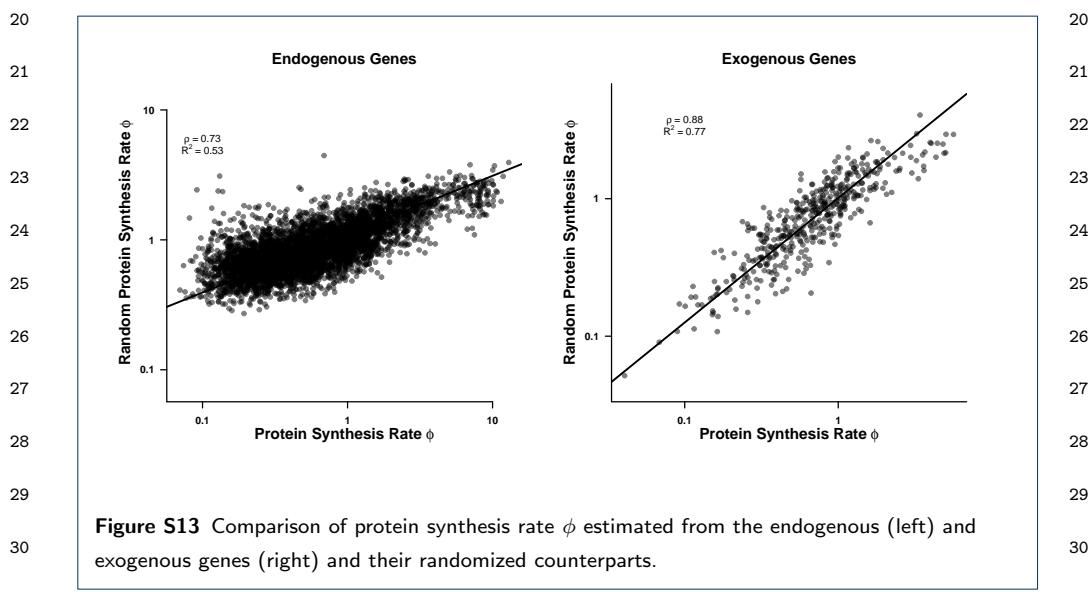








**Figure S12** Comparison of mutation bias  $\Delta M$  (left) and selection bias  $\Delta \eta$  (right) estimated from the randomized endogenous and the di-nucleotide randomized exogenous genes.



**Figure S13** Comparison of protein synthesis rate  $\phi$  estimated from the endogenous (left) and exogenous genes (right) and their randomized counterparts.