

## 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:** Together, our work illustrates the advantage of mechanistic, population genetic models like ROC SEMPPR and the quantitative estimates they provide when analyzing sequence data.

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

The conclusion is too  
focused on our method.  
Emphasize the biological  
insight more.

## 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

codon usage evolution that accounts for and, in turn, can estimate the contribution of mutation bias  $\Delta M$ , selection bias. The mathematics of ROC SEMPPR are derived on a mechanistic description of ribosome movement along an mRNA, although the approximation of other biological mechanisms could also be consistent with the model. Broadly speaking, ROC SEMPPR 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 usageDE between synonymous codons and protein synthesis rate  $\phi$  to the patterns of codon usage observed within a set of genes. Briefly, the set of  $\Delta M$  for an amino acid quantifies the relative differences in mutational stability or bias between the synonymous codons of the amino acid  $S$ . In the absence of selection bias (or equivalently when gene expression  $\phi = 0$ ), the equilibrium frequency of synonymous codon  $i$  is simply  $\exp[-\Delta M_i] / \left( \sum_{j \in S} \exp[-\Delta M_j] \right)$ . Because the time units of protein production rate have no intrinsic time scale, we define the average protein production rate for a set of genes to be one, i.e.  $\bar{\phi} = 1$  by definition [10]. In order to facilitate comparisions between gene sets, we express both,  $\Delta M$  and  $\Delta \eta$ , as deviation from the mean of each synonymous codon family (see Materials and Methods for details). Nevertheless, the difference  $\Delta \eta$  describes the difference in fitness between two synonymous codons relative to drift for a gene whose protein production rate  $\phi$  is equal to the the average rate of protein production  $\bar{\phi}$  across the set of genes. In other words, for a gene whose protein is expressed at the average rate, for any two given synonymous codons  $i$  and  $j$ ,  $\Delta \eta_i - \Delta \eta_j = N_e s$ .

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

mikeg: Is the Lachancea  
de synonymous with the  
achancea genus? If so use

*Lachancea*

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

## 28 Results

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

30 We used our software package AnaCoDa [22] to compare model fits of ROC  
31 SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two  
32 sets of 4,864 endogenous and 497 exogenous genes. These two set where initially  
33 identified based on their striking difference in GC content [15], with very little over-

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

16 Bayes factor strongly support the hypothesis that the *L. kluyveri* genome consists  
17 of genes with two different and distinct patterns of codon usage bias rather than a  
18 single ( $K = \exp(42,294)$ ; Table 1). We find additional support for this hypothesis  
19 when we compare our predictions of protein synthesis rate to empirically observed  
20 mRNA expression values as a proxy for protein synthesis. Specifically, we improve  
21 the variance explained by our predicted protein synthesis rates by  $\sim 42\%$ , from  
22  $R^2 = 0.33$  ( $p \approx 0$ ) to 0.46 ( $p \approx 0$ ) (Figure 1). While the implicit consideration of GC  
23 content in this analysis certainly plays a roll, it does not explain the improvement  
24 in  $R^2$  (Figure S1b).

#### 25

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

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

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  $\log(\mathcal{L}_M)$ , Bayes Factor K, and the p-value of the 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

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 exogenous genes were negatively correlated ( $\rho = -0.49, p = 3.56 \times 10^{-5}$ ), indicating weak similarity with only  $\sim 5\%$  of the codons share the same sign between the two mutation environments (Figure 2a). Overall, mutation bias favors codons ending in the purines A and T over codons ending in the pyrimidines G or C, respectively, as indicated by where the endogenous model fit curves intercept the left axis in Figure 3. One exception is Lys, L, where mutation appears to favor the codon TTG over TTA; however, this difference is very small and not statistically meaningful. The exogenous genes display behavior, favoring G over A and C over G and doing so more strongly, as indicated by shift in order and the greater distance between the exogenous model fit curves where they intercept the left axis in Figure 3).

We find that the signature of selection bias  $\Delta\eta$  also differs substantially between the endogenous and exogenous gene sets. In terms of their magnitude relative to mutation bias  $\Delta M$ ,  $\Delta\eta$  for the endogenous gene set is substantially greater than for the exogenous gene set as indicated by the fact that the protein production rate  $\phi$  where the effect of selection fails to override the effect of The difference in codon usage between endogenous and exogenous genes is striking as some amino acids have opposite codon preferences. As a result, our estimates of the optimal codon differ in nine cases between endogenous and exogenous genes (Figure 3, Table S2). For example, the usage of the Asparagine (Asn, N) codon AAC is increased in highly expressed endogenous genes but the same codon is depleted in highly expressed exogenous genes. For Aspartic acid (Asp, D), the combined genome shows the same codon preference in highly expressed genes as the exogenous gene set. Generally,

1  
2  
3  
4  
5  
6

mikeg: Note that Alex has implemented DIC routines into AnaCoDa which is better than BF because it is insensitive to the width of flat priors, which I erroneously thought didn't matter previously.

12  
13  
14  
15  
16  
17

18 mikeg: Cedric, can you verify that the 95% PI of  $\Delta M$  for TTG and TTA revised the  $\Delta M$  results because they weren't actually consistent with the curves illustrated in Figure 3. Please check and, if changes are correct, please revise  $\Delta\eta$  results to mikro: structure of  $\Delta M$  section. make it more consistent with previous discussions of these terms?

32 For example, 'opposite codon preferences' is confusing, is it that  $\Delta\eta$  values have opposite signs between gene sets?

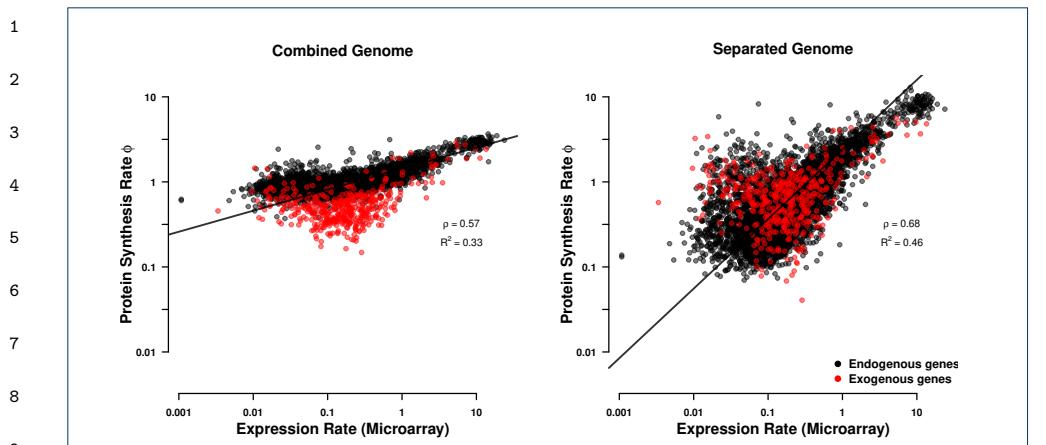


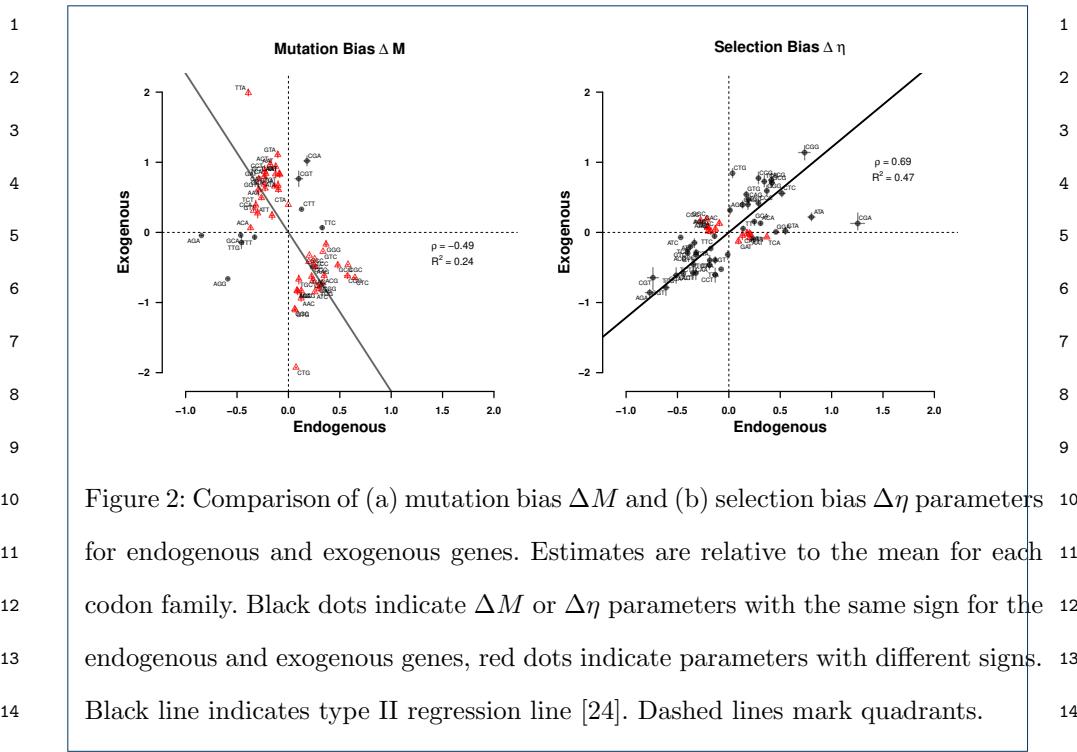
Figure 1: Comparison of predicted protein synthesis rate  $\phi$  to mRNA abundance from [23] for (a) the combined genome where mutation bias and selection bias parameters  $\Delta M$  and  $\Delta \eta$  are estimated for the combined endogenous and exogenous gene sets, and (b) where  $\Delta M$  and  $\Delta \eta$  are estimated separately for the endogenous and exogenous gene sets. Endogenous genes are displayed in black and exogenous genes in red. Black line indicates type II regression line [24].

mikeg: Type II is also referred to as ‘Deming regression’ or ‘errors-in-variables model’

fits to the complete *L. kluyveri* genome reveal that the relatively small exogenous gene set ( $\sim 10\%$  of genes) has a disproportionate effect on the model fit (Figure S2, S3).

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

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

1 Can Codon Usage Help Determine the Source of the Exogenous Genes 1  
2

3 Since the origin of the exogenous genes is currently unknown, we explored if the 2  
4 information on codon usage extracted from the exogenous genes can be used to 3  
5 identify a potential source lineage. We combined our estimates of mutation bias 4  
6  $\Delta M$  and selection bias  $\Delta\eta$  with synteny information and searched for potential 5  
7 source lineages of the introgressed exogenous region. We used  $\Delta M$  to identify 6  
8 candidate lineages as the endogenous and exogenous genes show greater dissimilarity 7  
9 in mutation bias than in selection bias. We examined 332 budding yeasts [21] and, 8  
10 identified the ten lineages with the highest correlation to the exogenous  $\Delta M$  9  
11 parameters as potential source lineages (Figure 4, Table 2). Two of the ten candidate 10  
12 lineages utilize the alternative yeast nuclear code (NCBI codon table 12). In this 11  
13 case, the codon CTG codes for Serine instead of Leucine. We therefore excluded the 12  
14 Leucine codon family from our comparison of codon families; however, there was no 13  
15 need to exclude Serine as CTG is not a one step neighbor of the remaining Serine 14  
16 codons. A mutation between CTG and the remaining Serine codons would require 15  
17 two mutations with one of them being non-synonymous, which would violate the 16  
18 weak mutation assumption of ROC SEMPPR. 17

19 The endogenous *L. kluyveri* genome exhibits codon usage very similar to most 19  
20 (77 %) yeast lineages examined, indicating that most of the examined yeasts share 20  
21 a similar codon usage (Figure S4). Only ~ 17% of all examined yeast show a pos- 21  
22 itive correlation in both,  $\Delta M$  and  $\Delta\eta$  with the exogenous genes, whereas the vast 22  
23 majority of lineages (~ 83%) show a negative correlation for  $\Delta M$ , only 21 % show 23  
24 a negative correlation for  $\Delta\eta$ . 24

25 Comparing synteny between the exogenous genes, which are restricted to the left 25  
26 arm of chromosome C, and the candidate yeast species we find that *E. gossypii* 26  
27 is the only species that displays high synteny (Table 2). Furthermore, the synteny 27  
28 relationship between the exogenous region and other yeasts appears to be limited 28  
29 to Saccharomycetaceae clade. Given these results, we conclude that, of the 332 29  
30 examined yeast lineages the *E. gossypii* lineage is the most likely source of the 30  
31 introgressed exogenous genes. Previous studies which studied the exogenous genes and 31  
32 chromosome recombination in the Lachancea clade concluded that the exogenous 32  
33 region originated from within the Lachancea clade, from an unknown or potentially 33

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

extinct lineage [15–17]. While it is not possible for us to dispute this hypothesis, our results provide a novel hypothesis about the origin of the exogenous genes.

To further test the plausibility of *E. gossypii* as potential source lineage, we identified 127 genes in our dataset [21] with homologous genes in *E. gossypii* and other Lachancea and used IQTree [25] to infer the phylogenetic relationship of the exogenous genes. Our results show that at least ~ 45% of exogenous genes (57/127) are more closely related to *E. gossypii* than to other Lachancea S5. Interestingly, our results also indicate that codon usage does not necessarily correlate with phylogenetic distance (Table 2).

### Estimating Introgression Age

If we assume that the exogenous genes originated from the *E. gossypii* lineage, we can estimate the age of the introgression based on our estimates of mutation bias  $\Delta M$ . We modeled the change in codon frequency over time as exponential decay, and estimated the age of the introgression assuming that *E. gossypii* still represents the mutation bias of its ancestral source lineage at the time of the introgression and a constant mutation rate. We infer the age of the introgression to be on the order of  $6.2 \pm 1.2 \times 10^8$  generations. Assuming *L. kluyveri* experiences between one and

1 eight generations per day, we estimate the introgression to have occurred between  
2 212,000 to 1,700,000 years ago. Our estimate places the time of the introgression  
3 earlier than the previous estimate of 19,000 - 150,000 years by [16].  
4

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

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

12 We define the selection against inefficient codon usage as the difference between the  
13 fitness on the log scale of an expected, replaced endogenous gene and the exogenous  
14 gene,  $s \propto \phi \Delta \eta$  due to the mismatch in codon usage parameters (See Methods for  
15 details). As the introgression occurred before the diversification of *L. kluyveri* and  
16 has fixed throughout all populations [16], we can not observe the original endogenous  
17 sequences that have been replaced by the introgression. Overall, we predict that a  
18 small number of low expression genes ( $\phi < 1$ ) were weakly exapted at the time of the  
19 introgression (Figure 5a). Thus, they appear to provide a small fitness advantage  
20 due to the accordance of exogenous mutation bias with endogenous selection bias  
21 (compare Figure S2 and S3). High expression genes ( $\phi > 1$ ) are predicted to have  
22 faced the largest selection against their mismatched codon usage in the novel cellular  
23 environment. In order to account for differences in the efficacy of selection on codon  
24 usage either due to the cost of pausing, differences in the effective population size,  
25 or the decline in fitness with every ATP wasted between the donor lineage and *L.*  
26 *kluyveri* we added a linear scaling factor  $\kappa$  to scale our estimates of  $\Delta \eta$  between the  
27 donor lineage and *L. kluyveri* and searched for the value that minimized the cost of  
28 the introgression, thus giving us the best case scenario (See Methods for details).  
29

30 Using our estimates of  $\Delta M$  and  $\Delta \eta$  from the endogenous genes and assuming the  
31 current exogenous amino acid composition of genes is representative of the replaced  
32 endogenous genes, we estimate the strength of selection against the exogenous genes  
33 at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of  
selection bias for the exogenous genes show that, while well correlated with the

1 endogenous genes, only nine amino acids share the same selectively preferred codon.  
2 Exogenous genes are, therefore, expected to represent a significant reduction in  
3 fitness for *L. kluyveri* due to mismatch in codon usage. Since  $\Delta\eta$  is proportional  
4 to the difference in fitness between the wild type and a mutant, we can use our  
5 estimates of  $\Delta\eta$  to approximate the selection against the exogenous genes  $\Delta s$  [10,  
6 26]. We estimate that the selection against all exogenous genes due to mismatched  
7 codon usage to have been  $\Delta s \approx -0.0008$  at the time of the introgression and  
8  $\approx -0.0003$  today. This reduction in  $\Delta s$  is primarily due to adaptive changes to the  
9 codon usage of the most highly expressed, introgressed genes (Figures 5a & S8).  
10 Based on the selection against the codon mismatch at the time of the introgression  
11 and assuming an effective population size  $N_e$  on the order of  $10^7$  [27], we estimate  
12 a fixation probability of  $(1 - \exp[-\Delta s])/(1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$  [26] for the  
13 exogenous genes. Clearly, the possibility of fixation under this simple scenario is  
14 effectively zero. In order for the exogenous genes to have reached fixation one or  
15 more exogenous loci must have provided a selective advantage not considered in  
16 this study (See Discussion).

## 18 Discussion

19 In order to study the evolutionary effects of the large scale introgression of the left  
20 arm of chromosome C, we used ROC SEMPPR, a mechanistic model of ribosome  
21 movement along an mRNA. The usage of a mechanistic model rooted in popula-  
22 tion genetics allows us generate more nuanced quantitative parameter estimates  
23 and separate the effects of mutation and selection on the evolution of codon usage.  
24 This allowed us to calculate the selection against the introgression, and provides *E.*  
25 *gossypii* as a potential source lineage of the introgression which was previously not  
26 considered. Our parameter estimates indicate that the *L. kluyveri* genome contains  
27 distinct signatures of mutation and selection bias from both an endogenous and ex-  
28 ogenous cellular environment. By fitting ROC SEMPPR separately to *L. kluyveri*'s  
29 endogenous and exogenous sets of genes we generate a quantitative description of  
30 their signatures of mutation bias and natural selection for efficient protein transla-  
31 tion.

32 In contrast to other methods such as RSCU, CAI, or tAI, ROC SEMPPR does  
33 not rely on external information such as gene expression or tRNA gene copy number

[5, 19]. Instead, ROC SEMPPR allows for the estimation of protein synthesis rate  $\phi$  and separates the effects of mutation and selection on codon usage. In addition, [20] showed that approaches like CAI are sensitive to amino acid composition, another property that distinguishes the endogenous and exogenous genes [15].

Previous work by [15] showed an increased bias towards GC rich codons in the exogenous genes but our results provide more nuanced insights by separating the effects of mutation bias and selection. We are able to show that the difference in GC content between endogenous and exogenous genes is mostly due to differences in mutation bias as 95% of exogenous codon families show a strong mutation bias towards GC ending codons (Table S1). However, the exogenous genes show a selective preference for AT ending codons for 90% of codon families (Table S2). Acknowledging the increased mutation bias towards GC ending codons and the difference in strength of selection between endogenous and exogenous genes by separating them also improves our estimates of protein synthesis rate  $\phi$  by 42% relative to the full genome estimate ( $R^2 = 0.46, p = 0$  vs.  $0.32, p = 0$ , respectively).

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

The estimates of mutation and selection bias parameters,  $\Delta M$  and  $\Delta \eta$ , are obtained under an equilibrium assumption. Given that the introgression is still adapting to its new environment, this assumption is clearly violated. However, the adaptation of the exogenous genes progresses very slowly as a quasi-static process as shown in this work as well as [16]. Therefore, the genome can be assumed to maintain an internal equilibrium at any given time. We see empirical evidence for this

1 behavior in our ability to predict gene expression and to correctly identify the low  
2 expression genes (Figure 1b).

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

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

32 In terms of understanding the spread of the introgression, we calculated the ex-  
33 pected selective cost of codon mismatch between the *L. kluyveri* and *E. gossypii*

lineages. Under our working hypothesis, the majority of the introgressed would have imposed a selective cost due to codon mismatch. Nevertheless,  $\sim 30\%$  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. Our estimate of the selective cost of codon mismatch on the order of  $-0.0008$ . While this selective cost may not seem very large, assuming *L. kluyveri* had a large  $N_e$ , the fixation probability of the introgression is the astronomically small value of  $\approx 10^{-6952} \approx 0$ . While this estimate heavily depends on the working hypothesis that the exogenous genes originated from the *E. gossypii* lineage, we can also calculate the hypothetical fixation probability if the current exogenous genes would introgress into *L. kluyveri*. Our estimate of the current selective cost of the mismatch of codon usage is on the order of  $-0.0003$ . The fixation probability of the current exogenous genes would still be astronomically small  $\approx 10^{-2609} \approx 0$ . These results are in accordance with previous work, highlighting the necessity of codon usage compatibility between endogenous and transferred exogenous genes [32, 33]. Thus, the basic scenario of an introgression between two yeast species with large  $N_e$  and where the introgression solely imposes a selective cost due to codon mismatch is clearly too simplistic.

One or more loci with a combined selective advantage on the order of 0.0008 or greater would have made the introgression change from disadvantageous to effectively neutral or advantageous. While this scenario seems plausible, it raises the question as to 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 [15, 18]. Estimates of the recombination rate as measured by crossovers (COs) for *L. kluyveri* are almost four times lower than for *S. cerevisiae* and about half that of *Schizosaccharomyces pombe* ( $\approx 1.6$  COs/Mb/meiosis,  $\approx 6$  COs/Mb/meiosis,  $\approx 3$  COs/Mb/meiosis) with no observed crossovers in the introgressed region [18], and no observed transposable elements [15]. 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. A population bottleneck reducing the  $N_e$  of the *L. kluyveri* lineage around the time of the introgression could also help explain the spread of the introgression. Compatible with these explanation is the possibility of several advantageous loci distributed

1 across the exogenous region drove a rapid selective sweep and/or the population  
2 through a bottleneck speciation process.

3 Assuming *E. gossypii* as potential source lineage of the exogenous region, we  
4 illustrated how information on codon usage can be used to infer the time since  
5 the introgression occurred using our estimates of mutation bias  $\Delta M$ . The  $\Delta M$   
6 estimates are well suited for this task as they are free of the influence of selection  
7 and unbiased by  $N_e$  and other scaling terms, which is in contrast to our estimates of  
8  $\Delta\eta$  [10]. Our estimated age of the introgression of  $6.2 \pm 1.2 \times 10^8$  generations is  $\sim 10$   
9 times longer than a previous minimum estimate by [16] of  $5.6 \times 10^7$  generations,  
10 which was based on the effective population recombination rate and the population  
11 mutation parameter [34]. Furthermore, these estimates assume that the current *E.*  
12 *gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the  
13 time of the introgression. Thus, if the ancestral mutation environments were more  
14 similar (dissimilar) at the time of the introgression then our result is an overestimate  
15 (underestimate).

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

## 21 Conclusion

22 Overall, our results show the usefulness of the separation of mutation bias and  
23 selection bias and the importance of recognizing the presence of multiple cellular  
24 environments in the study of codon usage. We also illustrate how a mechanistic  
25 model like ROC SEMPPR and the quantitative estimates it provides can be used for  
26 more sophisticated hypothesis testing in the future. In contrast to other approaches  
27 used to study codon usage like CAI [5] or tAI [19], ROC SEMPPR incorporates the  
28 effects of mutation bias and amino acid composition explicitly [20]. We highlight  
29 potential issues when estimating codon preferences, as estimates can be biased by  
30 the signature of a second, historical cellular environment. In addition, we show  
31 how quantitative estimates of mutation bias and selection relative to drift can be  
32 obtained from codon data and used to infer the fitness cost of an introgression as  
33 well as its history and potential future.

## Materials and Methods

## Separating Endogenous and Exogenous Genes

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

Model Fitting with ROC SEMPPR

ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) [22] and R (3.4.1) [35]. ROC SEMPPR was run from 10 different starting values for at least 250,000 iterations and thinned to every 50th iteration. After manual inspection to verify that the MCMC had converged, parameter posterior means, log posterior probability and log likelihood were estimated from the last 500 samples (last 10% of samples).

## Model selection

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

Comparing Codon Specific Parameter Estimates and Selecting Candidate lineages

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

$$\Delta M_i \equiv \Delta M_{i+1} - \overline{\Delta M_i} \quad (1)$$

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

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

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

### 14

### 15 Phylogenetic Analysis

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

### 27

### 28

### 29 Synteny Comparison

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

1      Estimating Age of Introgression

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

$$6 \quad \frac{dc_1}{dt} = -\mu_{1,2}c_1 - \mu_{2,1}(1 - c_1) \quad (3)$$

7      where  $\mu_{i,j}$  is the rate at which codon  $i$  mutates to codon  $j$  and  $c_1$  is the frequency  
 8      of the reference codon. Initial codon frequencies  $c_1(0)$  for each codon  
 9      family were taken from our mutation parameter estimates for *E. gossypii* where  
 10      $c_1(0) = \exp[\Delta M_{\text{gos}}]/(1 + \exp[\Delta M_{\text{gos}}])$ . Our estimates of  $\Delta M_{\text{endo}}$  can be used to  
 11     calculate the steady state of equation 3 were  $\frac{dc_1}{dt} = 0$  to obtain the equality  
 12      $\frac{\mu_{2,1}}{\mu_{1,2} + \mu_{2,1}} = \frac{1}{1 + \exp[\Delta M_{\text{endo}}]}$

$$14 \quad \frac{\mu_{2,1}}{\mu_{1,2} + \mu_{2,1}} = \frac{1}{1 + \exp[\Delta M_{\text{endo}}]} \quad (4)$$

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

$$17 \quad c_1(t) = \frac{1 + \exp[-X](K - 1)}{1 + \Delta M_{\text{endo}}} \quad (5)$$

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

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

26      Estimating Selection against Codon Mismatch

27      In order to estimate the selection against codon mismatch, we had to make three  
 28      key assumptions. First, we assumed that the current exogenous amino acid sequence  
 29      of a gene is representative of its ancestral state and the replaced endogenous gene  
 30      it replaced. Second, we assume that the currently observed cellular environment of

<sup>1</sup> *E. gossypii* reflects the cellular environment that the exogenous genes experienced  
<sup>2</sup> before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy  
<sup>3</sup> of selection between the cellular environments due to differences in either effective  
<sup>4</sup> population size  $N_e$  or the selective cost of an ATP  $q$  of the source lineage and *L.*  
<sup>5</sup> *kluyveri* can be expressed as a scaling constant and that protein synthesis rate  $\phi$   
<sup>6</sup> has not changed between the replaced endogenous and the introgressed exogenous  
<sup>7</sup> genes. Using estimates for  $N_e = 1.36 \times 10^7$  [27] for *Saccharomyces paradoxus* we  
<sup>8</sup> scale our estimates of  $\Delta\eta$  which explicitly contains the effective population size  $N_e$   
<sup>9</sup> [10] and define  $\Delta\eta' = \frac{\Delta\eta}{N_e}$ .

<sup>10</sup> All of our genome parameter estimations are scaled by lineage specific effects such  
<sup>11</sup> as  $N_e$ , the average, absolute gene expression level, and/or the proportionate fitness  
<sup>12</sup> value of an ATP. In order to account for these genome specific differences in scaling,  
<sup>13</sup> we scale the difference in the efficacy of selection on codon usage between the donor  
<sup>14</sup> lineage and *L. kluyveri* using a linear scaling factor  $\kappa$ . As  $\Delta\eta$  is defined as  $\Delta\eta =$   
<sup>15</sup>  $2N_e q(\eta_i - \eta_j)$ , we cannot distinguish if  $\kappa$  is a scaling on protein synthesis rate  $\phi$ ,  
<sup>16</sup> effective population size  $N_e$ , or the selective cost of an ATP  $q$  [4, 10]. We calculated  
<sup>17</sup> the selection against each genes codon mismatch assuming additive fitness effects  
<sup>18</sup> as

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

<sup>22</sup> where  $s_g$  is the overall strength of selection for translational efficiency on gene,  $g$   
<sup>23</sup> in the exogenous gene set,  $\kappa$  is a constant, scaling the efficacy of selection between  
<sup>24</sup> the endogenous and exogenous cellular environments,  $L_g$  is length of the protein in  
<sup>25</sup> codons,  $\phi_g$  is the estimated protein synthesis rate of the gene in the endogenous  
<sup>26</sup> environment, and  $\Delta\eta'_i$  is the  $\Delta\eta'$  for the codon at position  $i$ . As stated previously,  
<sup>27</sup> our  $\Delta\eta$  are relative to the mean of the codon family. We find that the selection  
<sup>28</sup> against the introgressed genes is minimized at  $\kappa \sim 5$  (Figure S7b). Thus, we expect  
<sup>29</sup> a five fold difference in the efficacy of selection between *L. kluyveri* and *E. gossypii*,  
<sup>30</sup> due to differences in either protein synthesis rate  $\phi$ , effective population size  $N_e$ ,  
<sup>31</sup> and/or the selective cost of an ATP  $q$ . Therefore, we set  $\kappa = 1$  if we calculate the  $s_g$   
<sup>32</sup> for the endogenous and the current exogenous genes, and  $\kappa = 5$  for  $s_g$  for selection  
<sup>33</sup> calculations at the time of introgression.

1 However, since we are unable to observe codon sequences of the replaced en-  
2 dogenous genes and for the exogenous genes at the time of introgression, instead  
3 of summing over the sequence, we calculate the expected codon count  $E[n_{g,i}]$  for  
4 codon  $i$  in gene  $g$  simply as the probability of observing codon  $i$  multiplied by the  
5 number of times the corresponding amino acids is observed in gene  $g$ , yielding:  
6

$$7 E[n_{g,i}] = P(c_i | \Delta M, \Delta \eta, \phi) \times m_{a_i}$$
$$8 = \frac{\exp[-\Delta M_i - \Delta \eta_i \phi_g]}{\sum_j^C \exp[-\Delta M_j - \Delta \eta_j \phi_g]} \times m_{a_i}$$
$$9$$

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

$$13 s_g = \sum_{i=1}^C -\kappa \phi_g \Delta \eta'_i E[n_{g,i}] \quad 14$$
$$15$$

16 We report the selection due to mismatched codon usage of the introgression as  
17  $\Delta s_g = s_{\text{intro},g} - s_{\text{endo},g}$  where  $s_{\text{intro},g}$  is the selection against an introgressed gene  $g$   
18 either at the time of the introgression or presently.  
19

#### Acknowledgments

20 The authors would like to thank Alexander Cope for helpful criticisms and suggestions for this work.  
21

#### Availability of data and materials

22 Parameter estimates generated during this study are available from the corresponding author. All remaining data  
23 generated during this study are included in this published article as figures, tables.  
24

#### Authors' contributions

25 CL and MAG initiated the study. CL collected and analyzed the data and wrote the manuscript. MAG and BCO  
26 edited the manuscript. CL, MAG, BCO, and RZ contributed to the data analysis and acquiring of funding. All  
27 Authors approved the final manuscript.  
28

#### Funding

29 This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ), MCB-1546402 (A. Von Arnim and  
30 MAG), and DEB-1355033 (BCO, MAG, and RZ) with additional support from Department of Ecology &  
31 Evolutionary Biology (EEB) at the University of Tennessee Knoxville (UTK) and the National Institute for  
32 Mathematical and Biological Synthesis (NIMBioS), an Institute sponsored by the National Science Foundation  
33 through NSF Award DBI-1300426. CL received support as a Graduate Student Fellow from NIMBioS with  
additional support from Departments of Mathematics and EEB at UTK.  
34

#### Ethics approval and consent to participate

35 Not applicable  
36

#### Consent for publication

37 Not applicable  
38

1	<b>Competing interests</b>	1
2	The authors declare that they have no competing interests.	2
3	<b>Author details</b>	3
4	<sup>1</sup> Department of Ecology & Evolutionary Biology, University of Tennessee, 37996, Knoxville, TN, USA. <sup>2</sup> National	4
5	Institute for Mathematical and Biological Synthesis, 37996, Knoxville, TN, USA. <sup>3</sup> Max-Planck Institute of	5
6	Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, 01307, Dresden, Germany. <sup>4</sup> Department of Business	6
	Analytics and Statistics, University of Tennessee, 37996, Knoxville, TN, USA.	
7	<b>References</b>	7
8	1. Gouy, M., Gautier, C.: Codon usage in bacteria: correlation with gene expressivity. <i>Nucleic Acids Research</i> <b>10</b> , 7055–7074 (1982)	8
9	2. Ikemura, T.: Codon usage and tRNA content in unicellular and multicellular organisms. <i>Molecular Biology and Evolution</i> <b>2</b> , 13–34 (1985)	9
10	3. Bulmer, M.: The selection-mutation-drift theory of synonymous codon usage. <i>Genetics</i> <b>129</b> , 897–907 (1990)	10
11	4. Gilchrist, M.A.: Combining models of protein translation and population genetics to predict protein production rates from codon usage patterns. <i>Molecular Biology and Evolution</i> <b>24</b> (11), 2362–2372 (2007)	11
12	5. Sharp, P.M., Li, W.H.: The codon adaptation index - a measure of directional synonymous codon usage bias, and its potential applications. <i>Nucleic Acids Research</i> <b>15</b> , 1281–1295 (1987)	12
13	6. Wright, F.: The 'effective number of codons' used in a gene. <i>Genet</i> <b>87</b> , 23–29 (1990)	13
14	7. M, S.P., Stenico, M., Peden, J.F., Lloyd, A.T.: Codon usage: mutational bias, translational selection, or both? <i>Biochem Soc Trans.</i> <b>21</b> (4), 835–841 (1993)	14
15	8. Shah, P., Gilchrist, M.A.: Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. <i>Proceedings of the National Academy of Sciences U.S.A</i> <b>108</b> (25), 10231–10236 (2011)	15
16	9. Wallace, E.W., Airoldi, E.M., Drummond, D.A.: Estimating selection on synonymous codon usage from noisy experimental data. <i>Molecular Biology and Evolution</i> <b>30</b> , 1438–1453 (2013)	16
17	10. 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. <i>Genome Biology and Evolution</i> <b>7</b> , 1559–1579 (2015)	17
18	11. Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., Danchin, A.: Evidence for horizontal gene transfer in <i>Escherichia coli</i> speciation. <i>Journal of Molecular Biology</i> <b>222</b> (4), 851–856 (1991)	18
19	12. Lawrence, J.G., Ochman, H.: Amelioration of bacterial genomes: Rates of change and exchange. <i>Journal of Molecular Biology</i> <b>44</b> , 383–397 (1997)	19
20	13. Marcet-Houben, M., Gabaldón, T.: Beyond the whole-genome duplication: Phylogenetic evidence for an ancient interspecies hybridization in the baker's yeast lineage. <i>PLoS Biology</i> <b>13</b> (8), 1002220 (2015)	20
21	14. 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. <i>Mol. Phylogenet. Evol.</i> <b>78</b> , 386–398 (2014)	21
22	15. 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. <i>Genome Research</i> <b>19</b> (10), 1710–1721 (2009)	22
23	16. Friedrich, A., Reiser, C., Fischer, G., Schacherer, J.: Population genomics reveals chromosome-scale heterogeneous evolution in a protoploid yeast. <i>Molecular Biology and Evolution</i> <b>32</b> (1), 184–192 (2015)	23
24	17. 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. <i>Genome research</i> <b>26</b> (7), 918–32 (2016)	24
25	18. 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. <i>PLoS Genetics</i> <b>13</b> (8), 1006917 (2017)	25
26		26
27		27
28		28
29		29
30		30
31		31
32		32
33		33

- 1 19. dos Reis, M., Savva, R., Wernisch, L.: Solving the riddle of codon usage preferences: a test for translational  
2 selection. *Nucleic Acids Research* **32**(17), 5036–5044 (2004) 1
- 3 20. Cope, A.L., Hettich, R.L., Gilchrist, M.A.: Quantifying codon usage in signal peptides: Gene expression and  
4 amino acid usage explain apparent selection for inefficient codons. *Biochimica et Biophysica Acta (BBA)* -  
5 *Biomembranes* **1860**(12), 2479–2485 (2018) 2
- 6 21. Shen, X.X., Opulente, D.A., Kominek, J., Zhou, X., Steenwyk, J.L., Buh, K.V., Haase, M.A.B., Wisecaver,  
7 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) 3
- 8 22. Landerer, C., Cope, A., Zaretzki, R., Gilchrist, M.A.: AnaCoDa: analyzing codon data with bayesian mixture  
models. *Bioinformatics* **34**(14), 2496–2498 (2018) 4
- 9 23. Tsankov, A.M., Thompson, D.A., Socha, A., Regev, A., Rando, O.J.: The role of nucleosome positioning in the  
evolution of gene regulation. *PLoS Biol* **8**(7), 1000414 (2010) 5
- 10 24. Sokal, R.R., Rohlf, F.J.: *Biometry - The principles and practice of statistics in biological*, pp. 547–555. W. H.  
Freeman, New York, NY (1981) 6
- 11 25. Nguyen, L.T., Schmidt, H.A., von Haeseler, A., Minh, B.Q.: Iq-tree: A fast and effective stochastic algorithm  
for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **32**(1), 268–274 (2015) 7
- 12 26. Sella, G., Hirsh, A.E.: The application of statistical physics to evolutionary biology. *Proceedings of the National  
Academy of Sciences of the United States of America* **102**, 9541–9546 (2005) 8
- 13 27. Wagner, A.: Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* **22**,  
1365–1374 (2005) 9
- 14 28. Nagylaki, T.: Evolution of a finite population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**,  
6278–6281 (1983) 10
- 15 29. Nagylaki, T.: Evolution of a large population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**,  
5941–5945 (1983) 11
- 16 30. Harrison, R.J., Charlesworth, B.: Biased gene conversion affects patterns of codon usage and amino acid usage  
in the *Saccharomyces sensu stricto* group of yeasts. *Molecular Biology and Evolution* **28**(1), 117–129 (2011) 12
- 17 31. Salichos, L., Rokas, A.: Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature*  
**497**, 327–331 (2013) 13
- 18 32. Medrano-Soto, A., Moreno-Hagelsieb, G., Vinuela, P., Christen, J.A., Collado-Vides, J.: Successful lateral  
transfer requires codon usage compatibility between foreign genes and recipient genomes. *Molecular Biology  
and Evolution* **21**(10), 1884–1894 (2004) 14
- 19 33. Tuller, T., Girshovich, Y., Sella, Y., Kreimer, A., Freilich, S., Kupiec, M., Gophna, U., Ruppin, E.: Association  
between translation efficiency and horizontal gene transfer within microbial communities. *Nucleic Acids  
Research* **39**(11), 4743–4755 (2011). doi:10.1093/nar/gkr054 15
- 20 34. Ruderfer, D.M., Pratt, S.C., Seidl, H.S., Kruglyak, L.: Population genomic analysis of outcrossing and  
recombination in yeast. *Nature Genetics* **38**(9), 1077–1081 (2006) 16
- 21 35. R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical  
Computing, Vienna, Austria (2013). R Foundation for Statistical Computing. <http://www.R-project.org/> 17
- 22 36. Gronau, Q.F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D.S., Forster, J.J.,  
Wagenmakers, E.J., Steingrover, H.: A tutorial on bridge sampling. *Journal of Mathematical Psychology* **81**,  
80–97 (2017) 18
- 23 37. Legendre, P.: Lmodel2: Model II Regression. (2018). R package version 1.7-3. 19
- 24 29 https://CRAN.R-project.org/package=lmodel2 20
- 25 38. Soderlund, C., Nelson, W., Shoemaker, A., Paterson, A.: Symap A system for discovering and viewing syntenic  
regions of fpc maps. *Genome Research* **16**, 1159–1168 (2006) 21
- 26 39. Soderlund, C., Bomhoff, M., Nelson, W.: Symap v3.4: a turnkey synteny system with application to plant  
genomes. *Nucleic Acids Research* **39**(10), 68 (2011) 22
- 27 40. Marais, G., Charlesworth, B., Wright, S.I.: Recombination and base composition: the case of the highly  
self-fertilizing plant *Arabidopsis thaliana*. *Genome Biology* **5**, 45 (2004) 23

1	41. Lang, G.I., Murray, A.W.: Estimating the per-base-pair mutation rate in the yeast <i>Saccharomyces cerevisiae</i> . <i>Genetics</i> <b>178</b> (1), 67–82 (2008)	1
2	42. Wolfram Research Inc.: Mathematica 11. (2017). <a href="http://www.wolfram.com">http://www.wolfram.com</a>	2
3		3
4		4
5		5
6		6
7		7
8		8
9		9
10		10
11		11
12		12
13		13
14		14
15		15
16		16
17		17
18		18
19		19
20		20
21		21
22		22
23		23
24		24
25		25
26		26
27		27
28		28
29		29
30		30
31		31
32		32
33		33

## Supplementary Material

Supporting Materials for *Unlocking a signal of introgression from codons in Lachancea kluveri using a mutation-selection model* 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	8
Cys C	TGC	TGT	TGC	TGC	9
Asp D	GAC	GAT	GAC	GAC	10
Glu E	GAG	GAA	GAG	GAG	
Phe F	TTC	TTT	TTT	TTT	11
Gly G	GGC	GGT	GGC	GGC	
His H	CAC	CAT	CAC	CAC	12
Ile I	ATC	ATT	ATC	ATA	13
Lys K	AAG	AAA	AAG	AAA	
Leu L	CTG	TTG	CTG	CTG	14
Asn N	AAC	AAT	AAC	AAT	
Pro P	CCG	CCA	CCG	CCG	15
Gln Q	CAG	CAA	CAG	CAG	
Arg R	CGC	AGA	AGG	CGG	16
Ser <sub>4</sub> S	TCG	TCT	TCG	TCG	17
Thr T	ACG	ACA	ACG	ACG	18
Val V	GTG	GTT	GTG	GTG	
Tyr Y	TAC	TAT	TAC	TAC	19
Ser <sub>2</sub> Z	AGC	AGT	AGC	AGC	

21 21

22 22

23 23

24 24

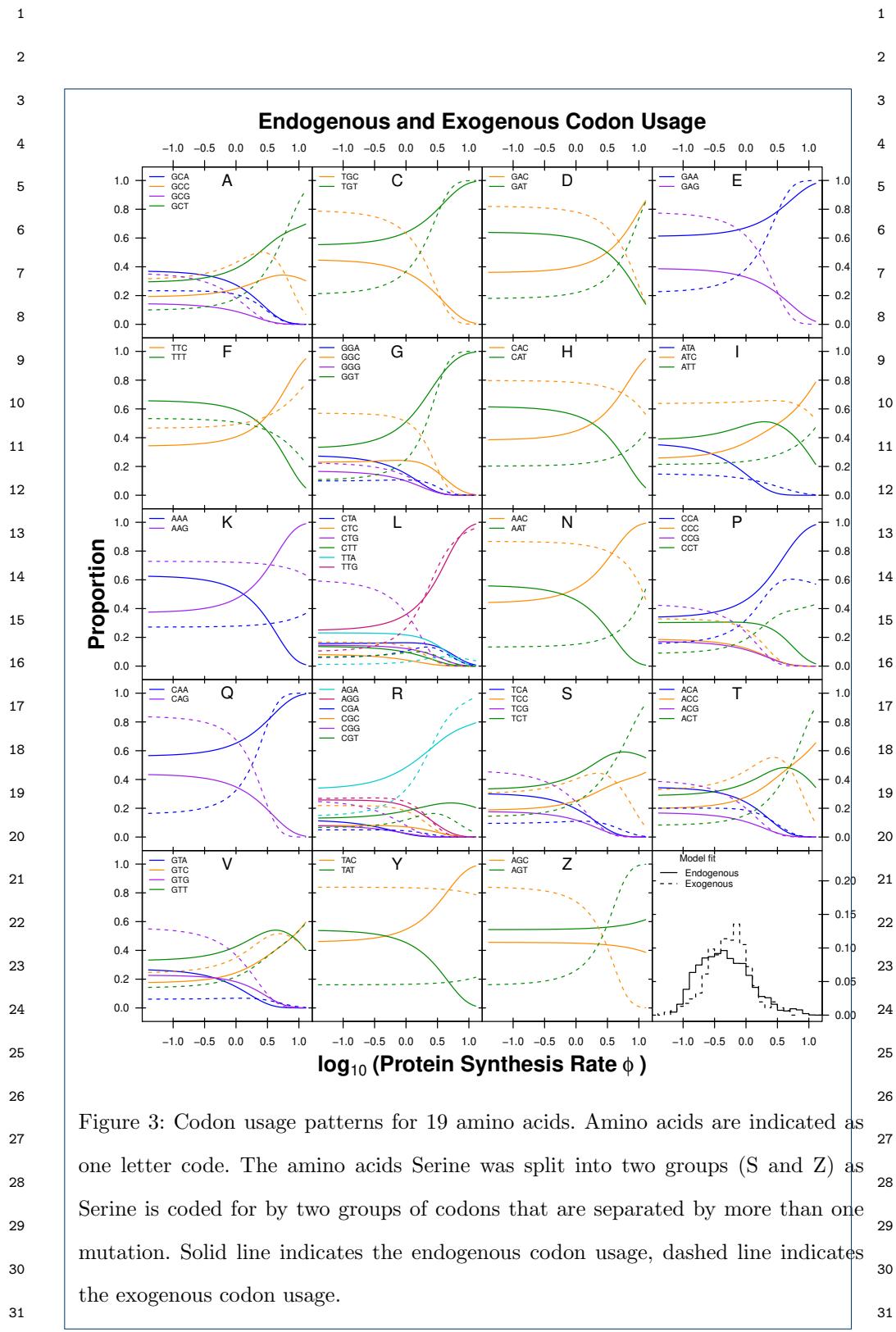
25 25

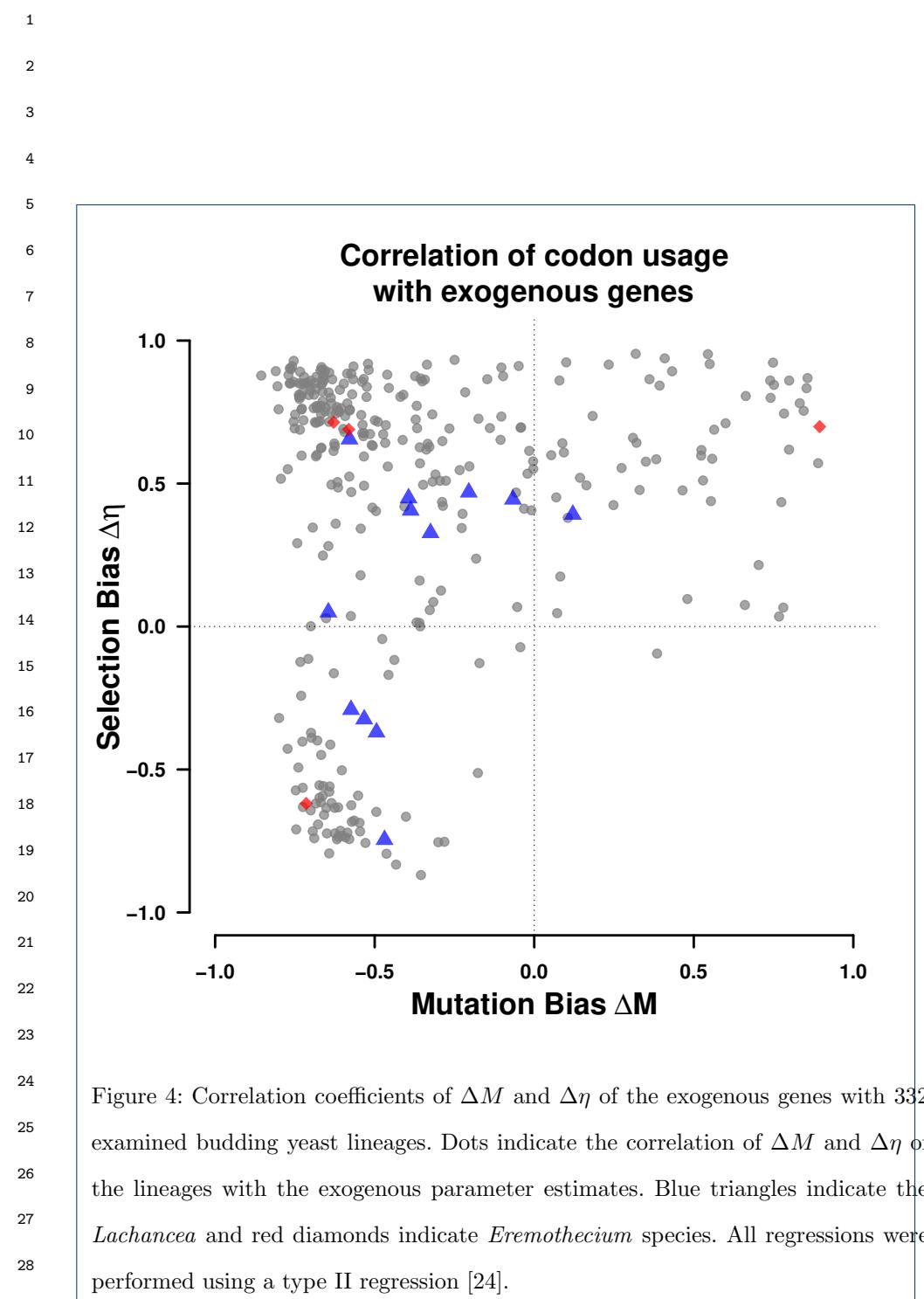
26

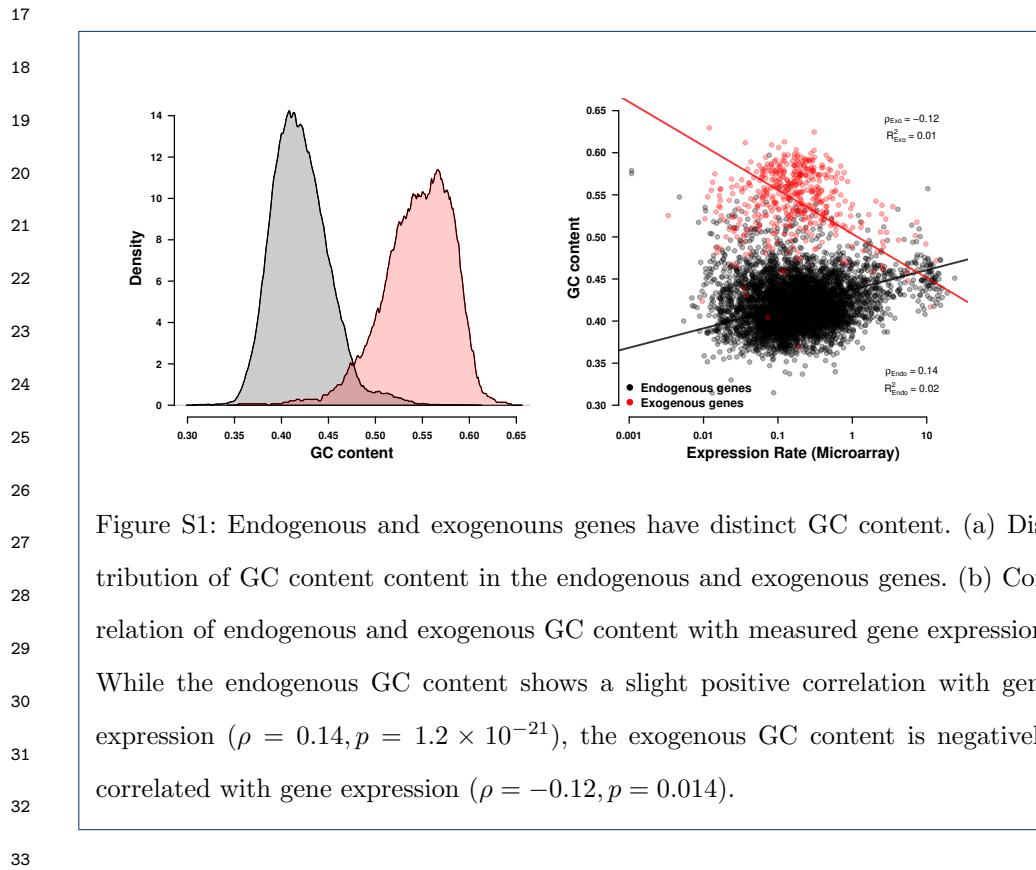
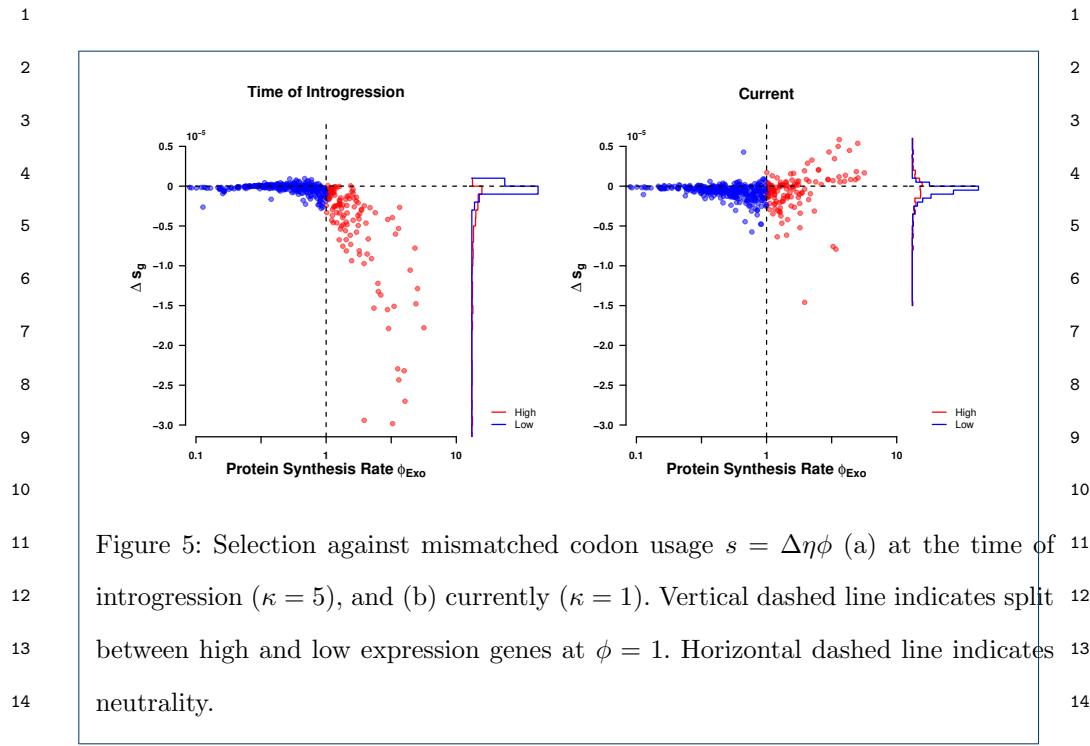
27

29 29

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







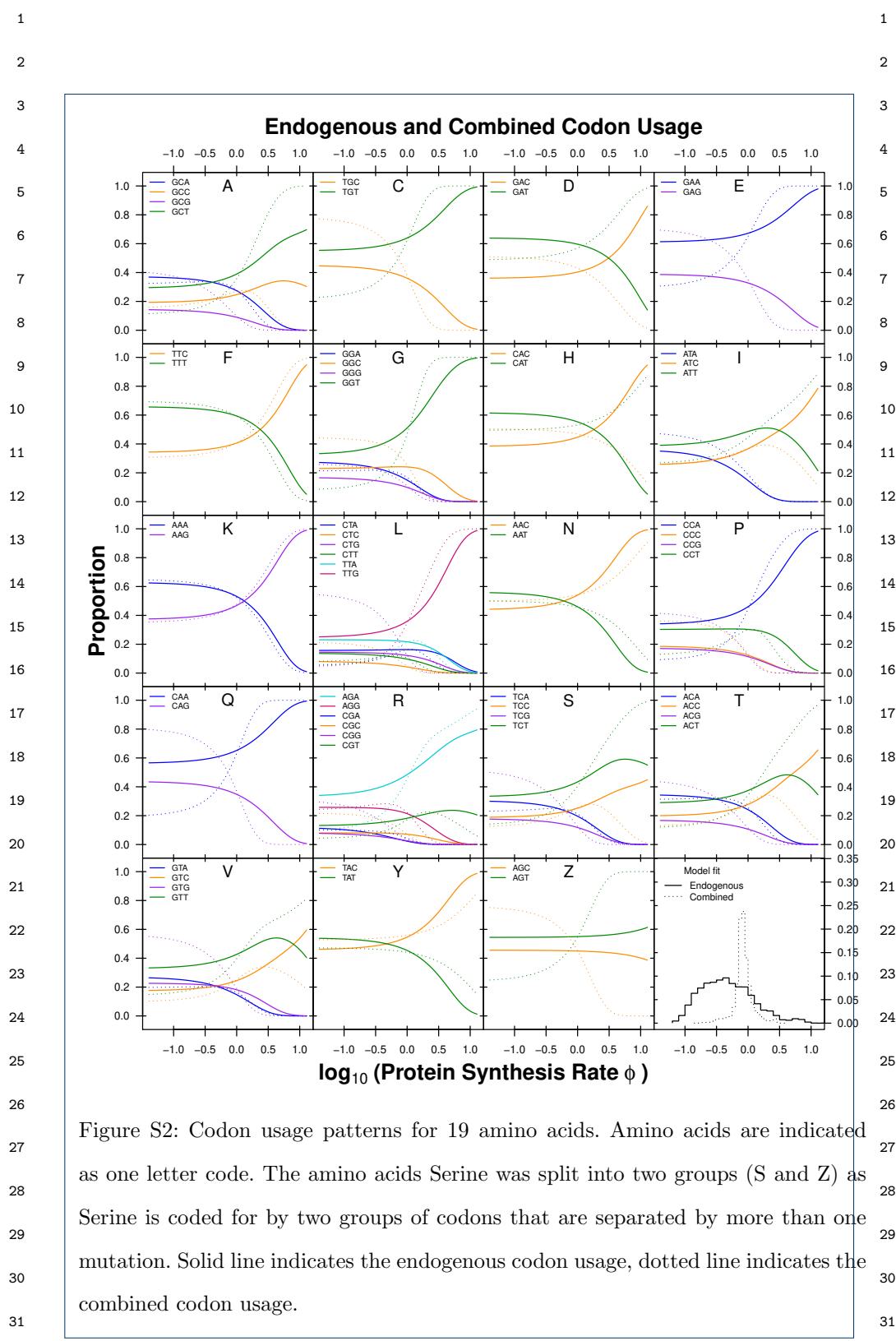


Figure S2: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dotted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

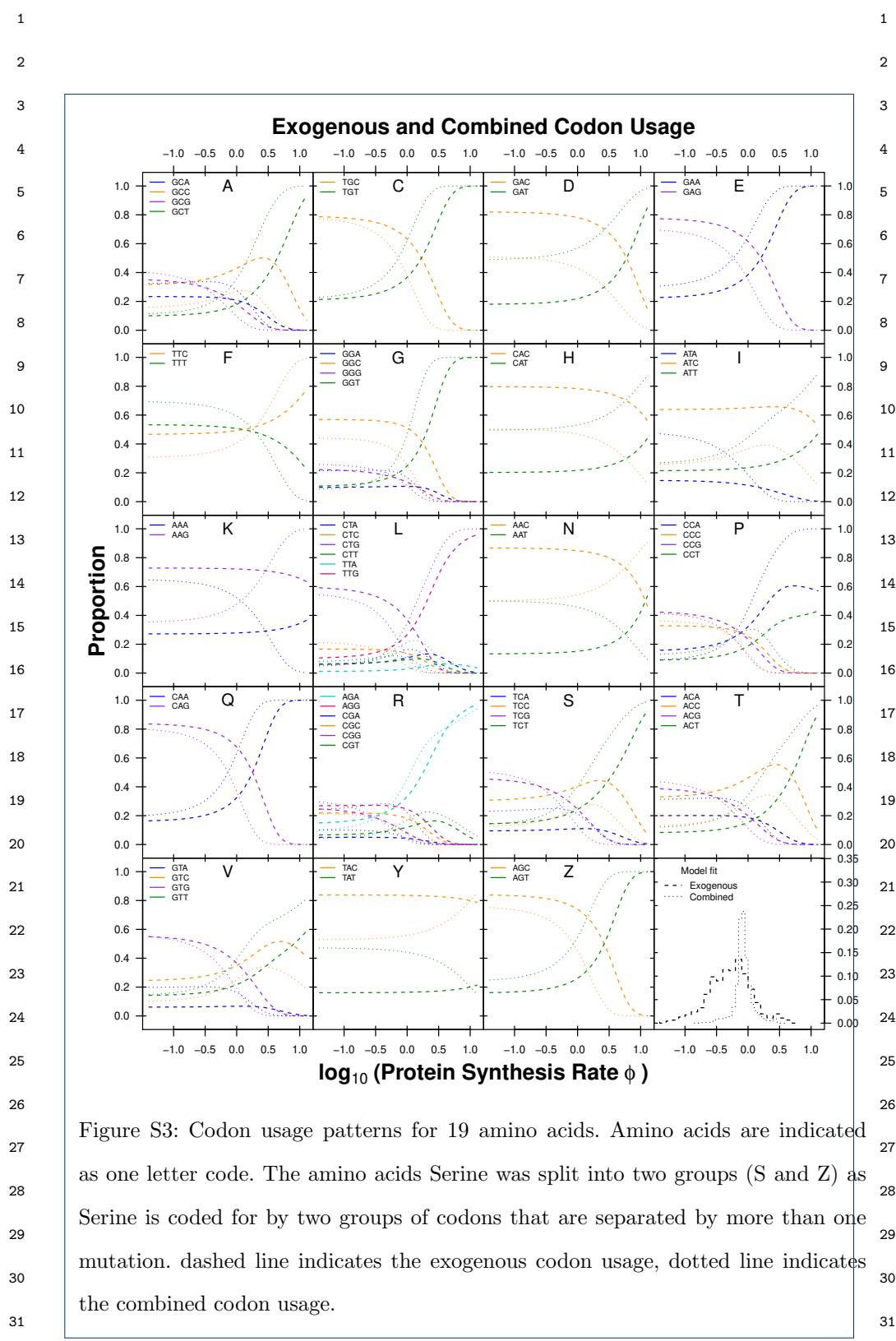
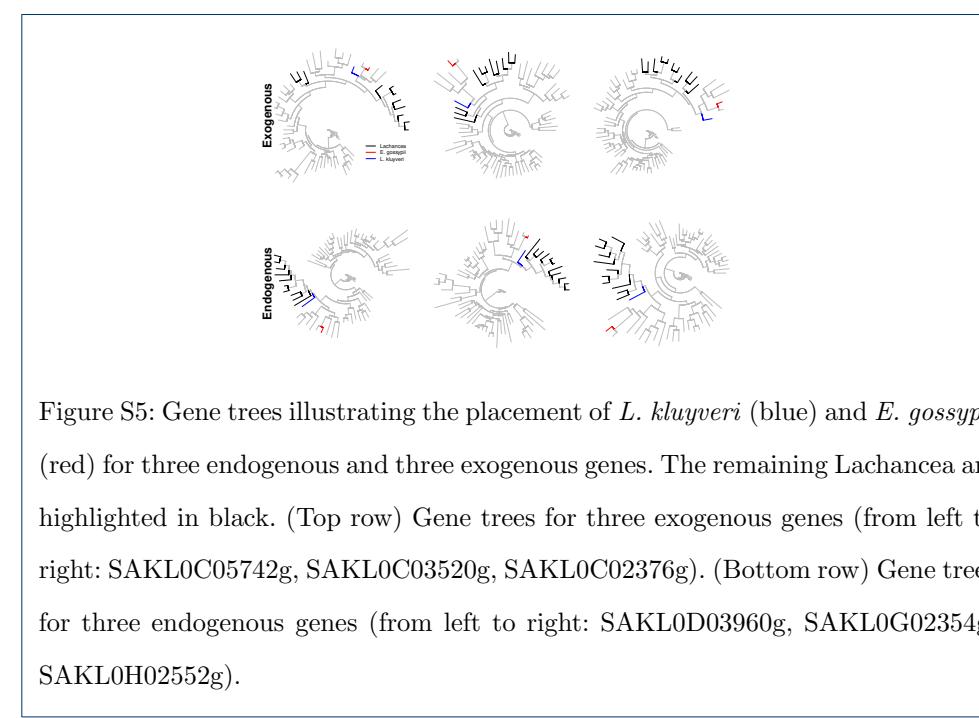
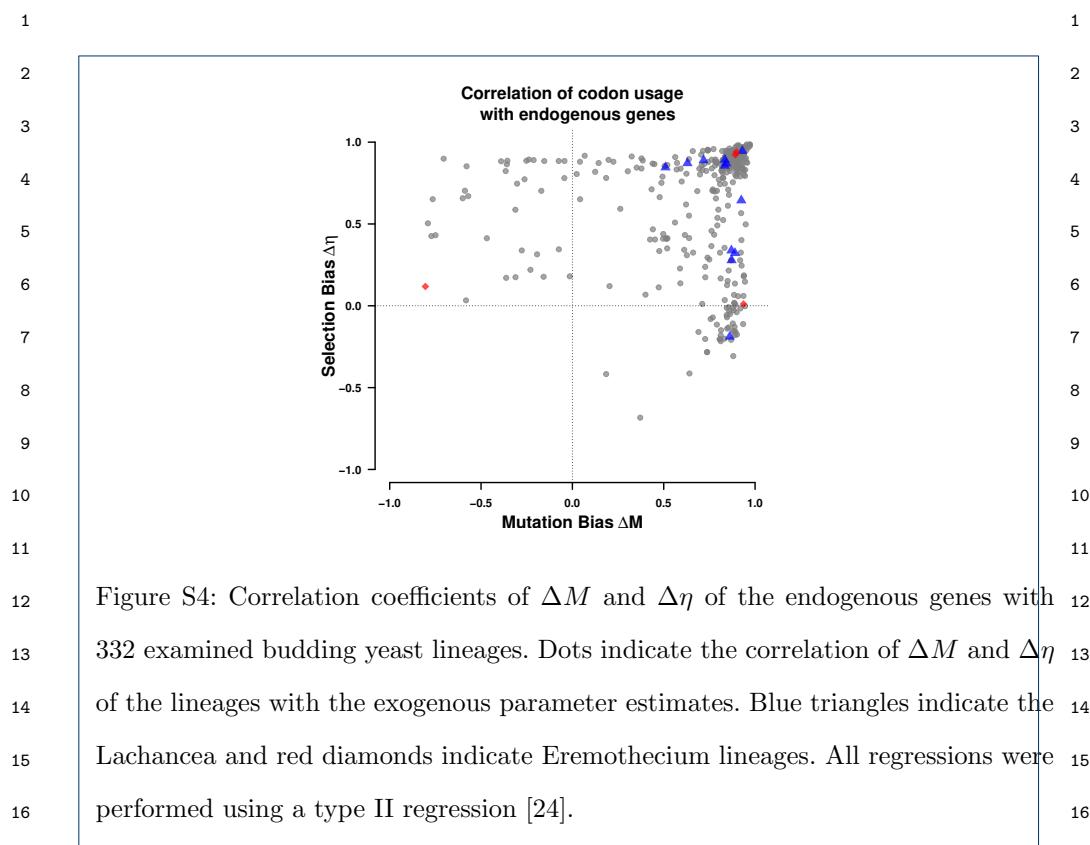


Figure S3: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. dashed line indicates the exogenous codon usage, dotted line indicates the combined codon usage.



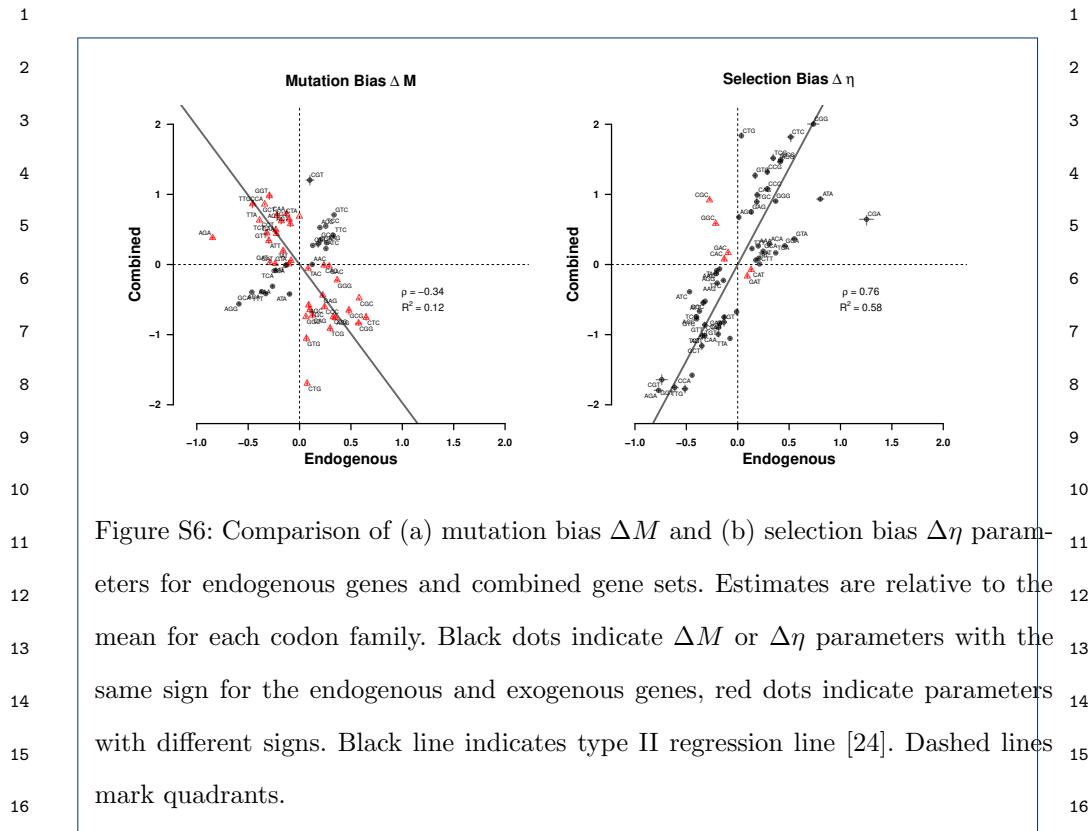


Figure S6: Comparison of (a) mutation bias  $\Delta M$  and (b) selection bias  $\Delta \eta$  parameters for endogenous genes and combined gene sets. Estimates are relative to the mean for each codon family. Black dots indicate  $\Delta M$  or  $\Delta \eta$  parameters with the same sign for the endogenous and exogenous genes, red dots indicate parameters with different signs. Black line indicates type II regression line [24]. Dashed lines mark quadrants.

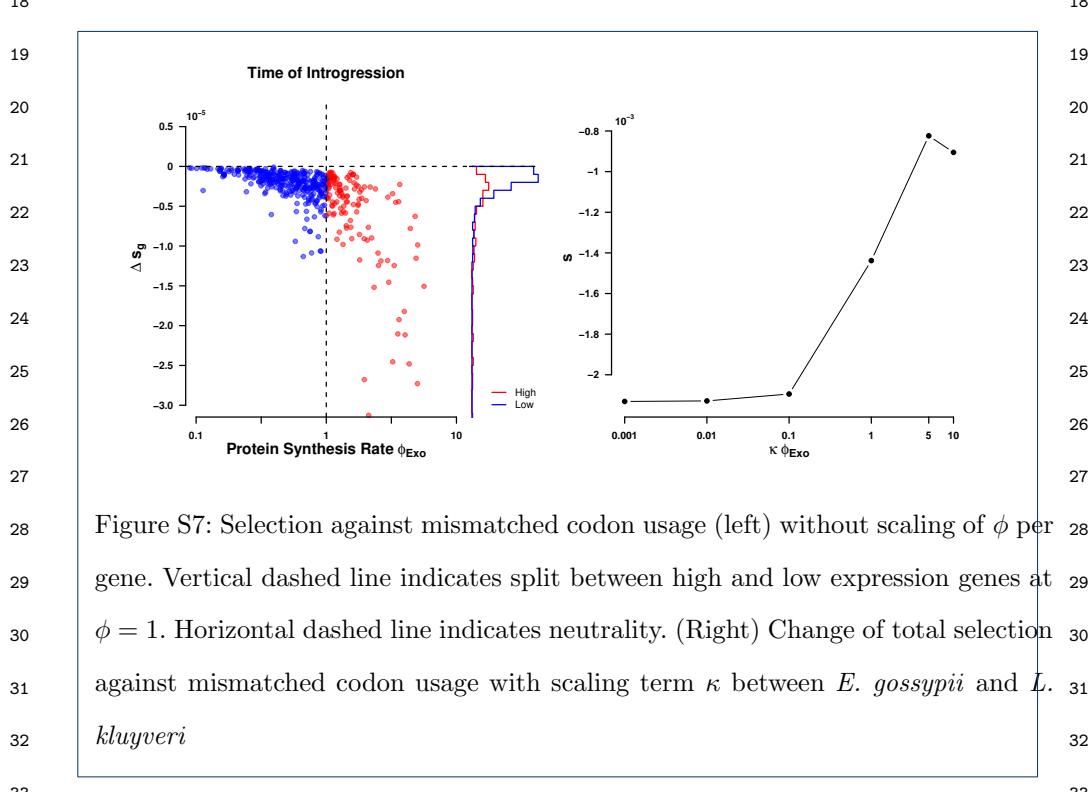


Figure S7: Selection against mismatched codon usage (left) without scaling of  $\phi$  per gene. Vertical dashed line indicates split between high and low expression genes at  $\phi = 1$ . Horizontal dashed line indicates neutrality. (Right) Change of total selection against mismatched codon usage with scaling term  $\kappa$  between *E. gossypii* and *L. kluyveri*

