

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^{1,2,3*}, Brian C O'Meara^{1,2}, Russell Zaretzki^{2,4} and Michael A Gilchrist^{1,2}

10
Submitted to BMC Evolutionary Biology on Nov 14, 2019

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

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

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

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.

Results: 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. 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 endogenous codon preferences. In addition, using our estimates of mutation bias and selection, we identify *Eremothecium gossypii* as the closest relative to the exogenous genes, providing an alternative hypothesis about the origin of the exogenous genes, estimate 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

Background

Synonymous codon usage patterns varies within a genome and between taxa, reflecting differences in mutation bias, selection, and genetic drift. The signature of

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

17 Most studies implicitly assume that the CUB of a genome is shaped by a single
18 cellular environment. As genes are horizontally transferred, introgress, or combined
19 to form novel hybrid species, one would expect to see the influence of multiple cellu-
20 lar environments on a genomes codon usage pattern [11, 12]. Given that transferred
21 genes are likely to be less adapted than endogenous genes to their new cellular en-
22 vironment, we expect a greater selection against mismatched codon usage in trans-
23 ferred genes if donor and recipient environment differ greatly in their selection bias,
24 making such transfers less likely. More practically, if differences in codon usage of
25 transferred genes are not taken into account for, they may distort the interpretation
26 of codon usage patterns. Such distortion could lead to the wrong inference of codon
27 preference for an amino acid [8, 10], underestimate the variation in protein synthesis
28 rate, or influence mutation estimates when analyzing a genome. While such gene
29 transfer events may be rare, this study aims to provide a general approach to study
30 the evolution of codon usage that could as well be applied between species.

31 To illustrate these ideas, we analyze the CUB of the genome of the yeast *Lachancea*
32 *kluyveri*, which is sister to all other Lachancea species. The Lachancea clade diverged
33 from the *Saccharomyces* clade, prior to its whole genome duplication ~ 100 Mya

1 ago [13, 14]. Since that time, *L. kluyveri* has experienced a large introgression of
2 exogenous genes (1 Mb, 457 genes) which is found in all of its populations [15, 16],
3 but in no other known Lachancea species [17]. The introgression replaced the left
4 arm of the C chromosome and displays a 13% higher GC content than the en-
5 doogenous *L. kluyveri* genome [15, 16]. Previous studies suggest that the source of
6 the introgression is probably a currently unknown or potentially extinct Lachancea
7 lineage based on gene concatenation or synteny relationships [15–18]. These char-
8 acteristics make *L. kluyveri* an ideal model to study the effects of an introgressed
9 cellular environment and the resulting mismatch in codon usage.

10 Using ROC SEMPPR, a Bayesian population genetics model based on a mech-
11 anistic description of ribosome movement along an mRNA, allows us to quantify
12 the cellular environment in which genes have evolved by separately estimating the
13 effects of mutation bias and selection bias on codon usage. While previous studies
14 have used information on gene expression to separate the effects of mutation and
15 selection on codon usage, ROC SEMPPR does not need such information but can
16 provide it. ROC SEMPPR's resulting predictions of protein synthesis rates have
17 been shown to be on par with laboratory measurements [8, 10]. In contrast to often
18 used heuristic approaches to study codon usage [5, 6, 19], ROC SEMPPR explic-
19 itly incorporates and distinguishes between mutation and selection effects on codon
20 usage and properly weights by amino acid usage [20]. We use ROC SEMPPR to in-
21 dependently describe two cellular environments reflected in the *L. kluyveri* genome;
22 the signature of the current environment in the endogenous genes and the decaying
23 signature of the exogenous environment in the introgressed genes. Our results in-
24 dicate that the difference in GC content between endogenous and exogenous genes
25 is mostly due to the differences in mutation bias of their ancestral environments.
26 Correcting for these different signatures of mutation bias and selection bias of the
27 endogenous and exogenous sets of genes substantially improves our ability to pre-
28 dict present day protein synthesis rates. These endogenous and exogenous gene set
29 specific estimates of mutation bias and selection bias, in turn, allow us to address
30 more refined questions of biological importance. For example, they allow us to pro-
31 vide an alternative hypothesis about the origin of the introgression and identify *E.*
32 *gossypii* as the nearest sampled relative of the source of the introgressed genes out
33 of the 332 budding yeast lineages with sequenced genomes [21]. While this hypoth-

1 esis is in contrast to previous work [15–18], we find support for it in gene trees and
2 synteny. We also estimate the age of the introgression to be on the order of 0.2 - 1.7
3 Mya, estimate the selection against these genes, both at the time of introgression
4 and now, and predict a detectable signature of CUB to persist in the introgressed
5 genes for another 0.3 - 2.8 Mya, highlighting the sensitivity of our approach.
6

7 Results

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

9 We used our software package AnaCoDa [22] to compare model fits of ROC
10 SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two
11 sets of 4,864 endogenous and 497 exogenous genes. These two set where initially
12 identified based on their striking difference in GC content [15], with very little over-
13 lap in GC content between the two sets (Figure S1a). ROC SEMPPR is a statistical
14 model that relates the effects of mutation bias ΔM , selection bias $\Delta \eta$ between syn-
15 onymous codons and protein synthesis rate ϕ , to explain the observed codon usage
16 patterns. Thus, the probability of observing a synonymous codon is proportional
17 to $p \propto \exp(-\Delta M - \Delta \eta \phi)$ [10]. Briefly, ΔM describes the mutation bias between
18 two synonymous codons at stationarity under a time reversible mutation model.
19 Because ROC SEMPPR only considers the stationary probabilities, only variation
20 in mutation bias, not absolute mutation rates can be detected. $\Delta \eta$ describes the
21 fitness difference between two synonymous codons relative to drift [10]. Since $\Delta \eta$ is
22 scaled by protein synthesis rate ϕ , this term is dominant in highly expressed genes
23 and tends towards 0 in low expression genes, allowing us to separate the effect of
24 mutation bias and selection bias on codon usage. We express both, ΔM and $\Delta \eta$,
25 as deviation from the mean of each synonymous codon family which prevents that
26 the choice of the reference codon affects our results (see Materials and Methods for
27 details).

28 Bayes factor strongly support the hypothesis that the *L. kluyveri* genome consists
29 of genes with two different and distinct patterns of codon usage bias rather than a
30 single ($K = \exp(42,294)$; Table 1). We find additional support for this hypothesis
31 when we compare our predictions of protein synthesis rate to empirically observed
32 mRNA expression values as a proxy for protein synthesis. Specifically, we improve
33 the variance explained by our predicted protein synthesis rates by $\sim 42\%$, from

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.**

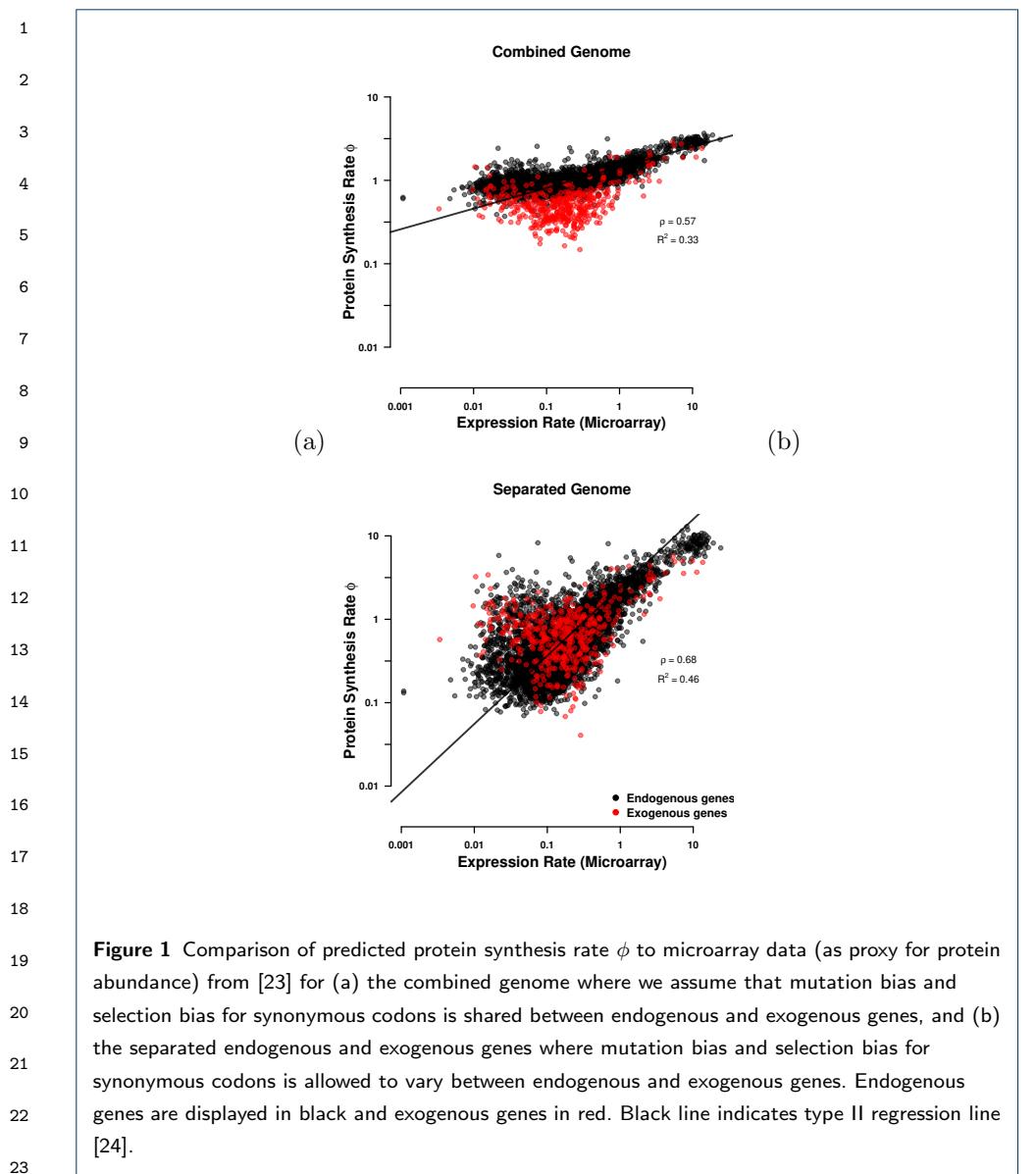
Hypothesis	$\log(\mathcal{L})$	n	$\log(\mathcal{L}_M)$	$\log(K)$	p
Combined	-2,650,047	5,483	-2,657,582	—	—
Separated	-2,612,397	5,402	-2,615,288	42,294	0

7
 8 $R^2 = 0.33 (p \approx 0)$ to $0.46 (p \approx 0)$ (Figure 1). While the implicit consideration of GC
 9 content in this analysis certainly plays a roll, it does not explain the improvement
 10 in R^2 (Figure S1b)

11
 12 Comparing Differences in the Endogenous and Exogenous Codon Usage

13 ROC SEMPPR constraints $E[\phi] = 1$, allowing us to interpret $\Delta\eta$ as selection on
 14 codon usage of the average gene with $\phi = 1$ and gives us the ability to compare the
 15 efficacy of selection sN_e across genomes. While it may be expected for the endoge-
 16 nous and exogenous genes to differ in their codon usage pattern due to the large
 17 difference in GC content it is not clear if this difference can be attributed to differ-
 18 ences in mutation or selection between endogenous genes. To better understand the
 19 differences in the endogenous and exogenous cellular environments, we compared
 20 our parameter estimates of mutation bias ΔM and selection $\Delta\eta$ for the two sets of
 21 genes. Our estimates of ΔM for the endogenous and exogenous genes were nega-
 22 tively correlated ($\rho = -0.49, p = 3.56 \times 10^{-5}$), indicating weak similarity with only
 23 $\sim 5\%$ of the codons share the same sign between the two mutation environments
 24 (Figure 2a). Overall, the endogenous genes only show a selection preference for C
 25 and G ending codons in $\sim 58\%$ of the codon families. In contrast, the exogenous
 26 genes display a strong preference for A and T ending codons in $\sim 89\%$ of the codon
 27 families.

28 For example, the endogenous genes show a mutational bias for A and T ending
 29 codons in $\sim 95\%$ of the codon families (the exception being Phe, F). The exogenous
 30 genes display an equally consistent mutational bias towards C and G ending codons
 31 (Table S1). In contrast to ΔM , our estimates of $\Delta\eta$ for the endogenous and exoge-
 32 nous genes were positively correlated ($\rho = 0.69, p = 9.76 \times 10^{-10}$) and showing the
 33 same sign in $\sim 53\%$ of codons between the two selection environments (Figure 2).



We find that the efficacy of selection within each codon family differs between sets of genes. 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, fits to the complete *L. kluyveri* genome reveal

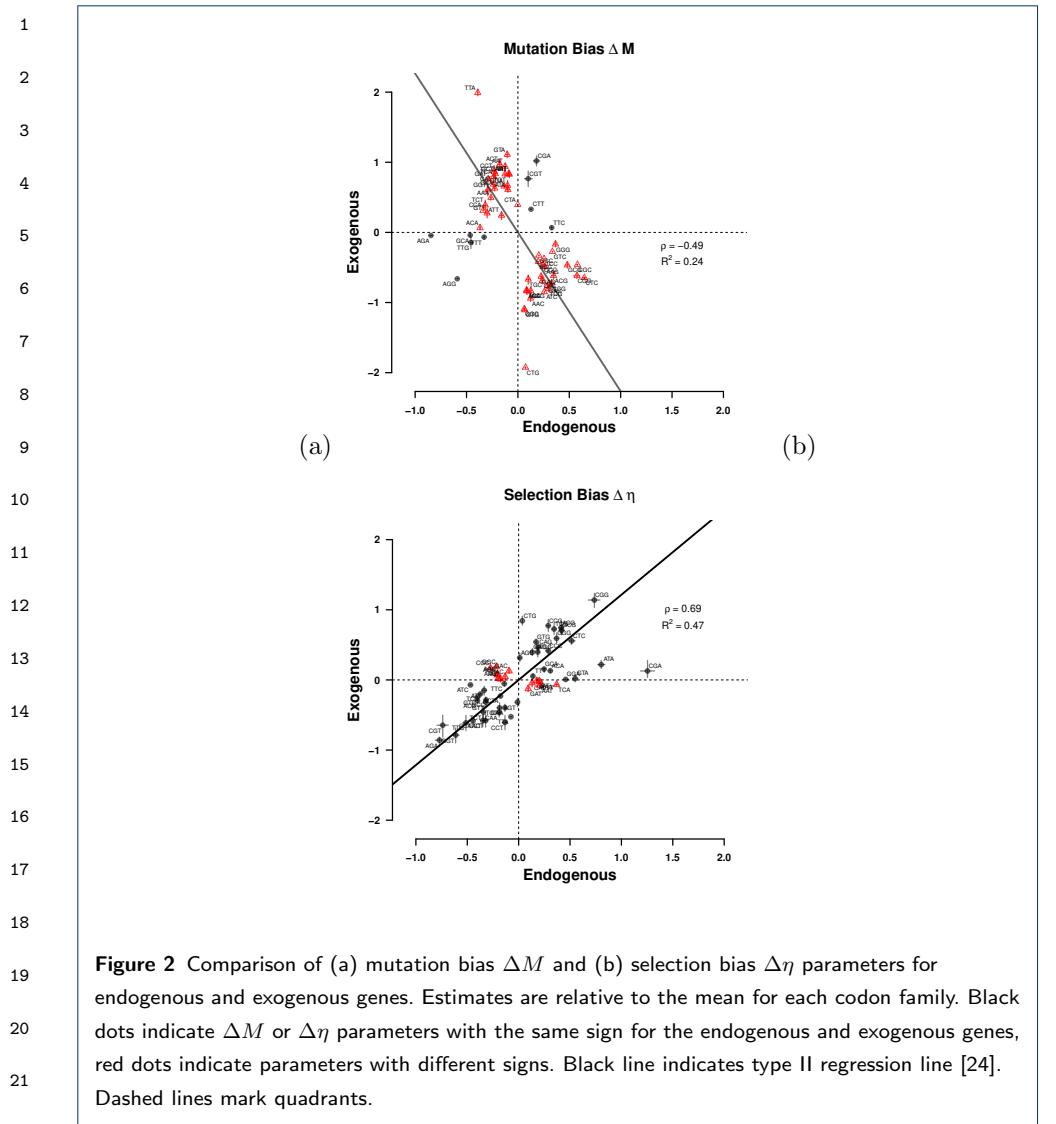


Figure 2 Comparison of (a) mutation bias ΔM and (b) selection bias $\Delta \eta$ parameters for endogenous and exogenous genes. Estimates are relative to the mean for each codon family. Black dots indicate ΔM or $\Delta \eta$ parameters with the same sign for the endogenous and exogenous genes, red dots indicate parameters with different signs. Black line indicates type II regression line [24]. Dashed lines mark quadrants.

that the relatively small exogenous gene set ($\sim 10\%$ of genes) has a disproportionate effect on the 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; and
29 Pro, P) the endogenous genes favor the codon with the most abundant tRNA. For
30 the remaining four cases (Ile, I; Ser, S; Thr, T; and Val, V), there are no tRNA genes
31 for the wobble free cognate codon encoded in the *L. kluyveri* genome. However, the
32 codon preference of these four amino acids in the exogenous genes matches the most
33 abundant tRNA encoded in the *L. kluyveri* genome.

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

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

23 Can Codon Usage Help Determine the Source of the Exogenous Genes

24 Since the origin of the exogenous genes is currently unknown, we explored if the
25 information on codon usage extracted from the exogenous genes can be used to
26 identify a potential source lineage. We combined our estimates of mutation bias
27 ΔM and selection bias $\Delta\eta$ with synteny information and searched for potential
28 source lineages of the introgressed exogenous region. We used ΔM to identify can-
29 didate lineages as the endogenous and exogenous genes show greater dissimilarity
30 in mutation bias than in selection bias. We examined 332 budding yeasts [21] and,
31 identified the ten lineages with the highest correlation to the exogenous ΔM pa-
32 rameters as potential source lineages (Figure 4, Table 2). Two of the ten candidate
33 lineages utilize the alternative yeast nuclear code (NCBI codon table 12). In this

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 ΔM 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

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, ΔM and $\Delta \eta$ with the exogenous genes, whereas the vast majority of lineages (~ 83%) show a negative correlation for ΔM , 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 chromosome recombination in the Lachancea clade concluded that the exogenous region originated from within the Lachancea clade, from an unknown or potentially

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

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

10
11 **Estimating Introgression Age**

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

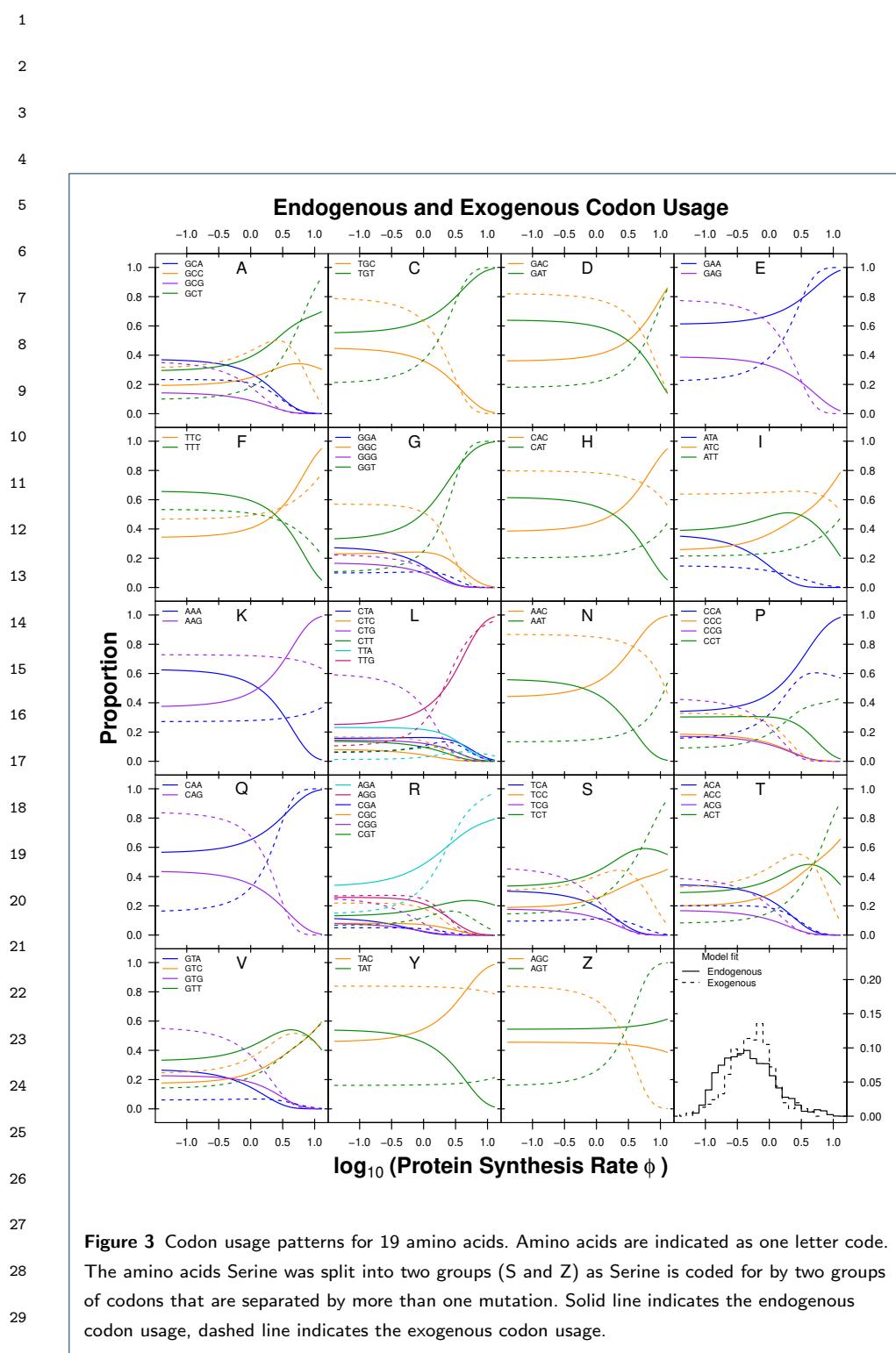
22 Using our model of exponential decay model, we also estimated the persistence of
23 the signal of the exogenous cellular environment. We predict that the ΔM signal of
24 the source cellular environment will have decayed to be within one percent of the
25 *L. kluyveri* environment in $\sim 5.4 \pm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

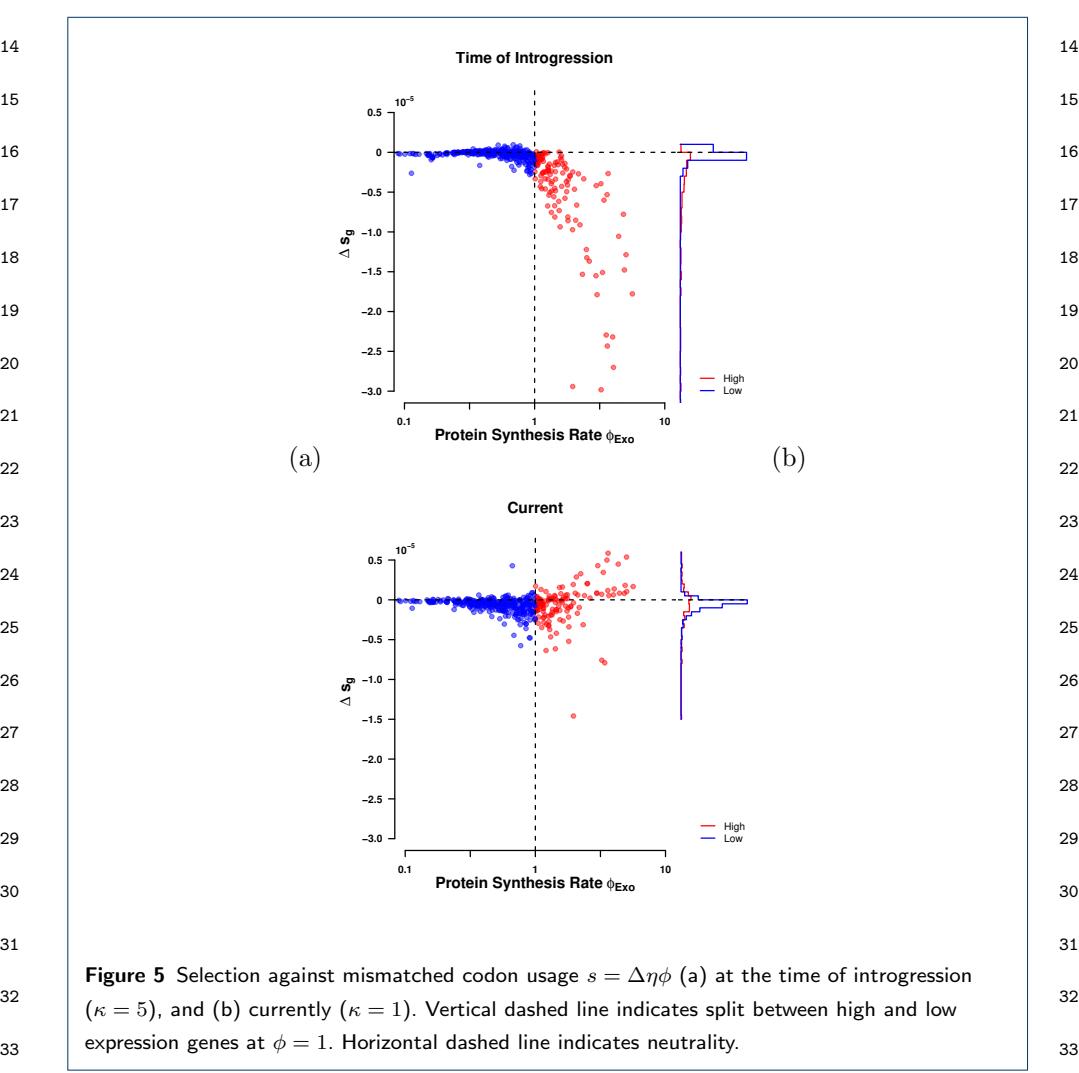
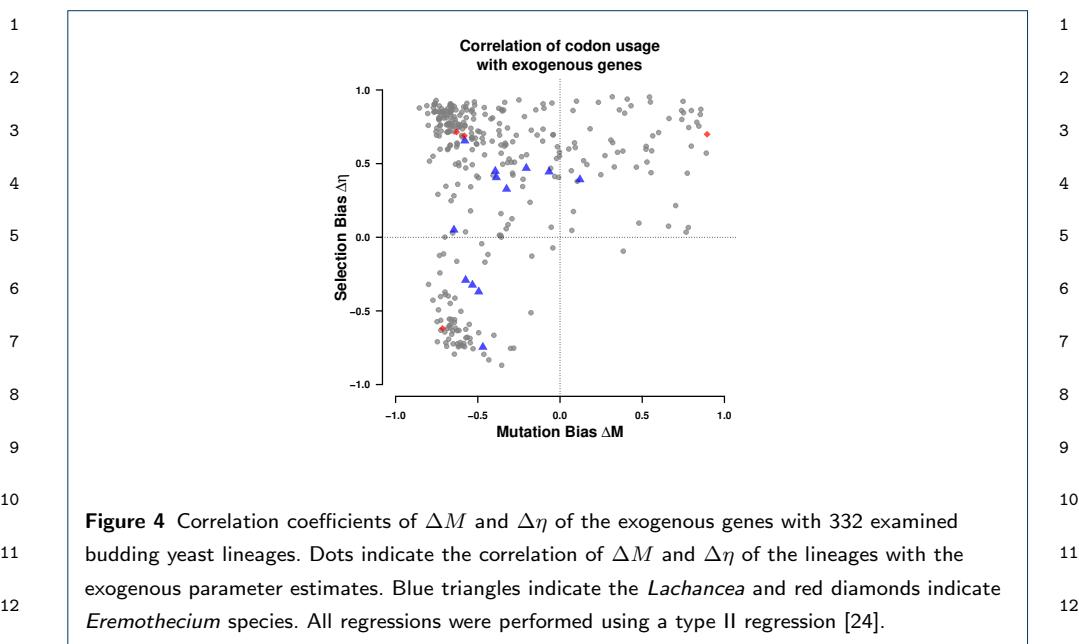
28
29 **Estimating Selection against Codon Mismatch of the Exogenous Genes**

30 We define the selection against inefficient codon usage as the difference between the
31 fitness on the log scale of an expected, replaced endogenous gene and the exogenous
32 gene, $s \propto \phi\Delta\eta$ due to the mismatch in codon usage parameters (See Methods for
33 details). As the introgression occurred before the diversification of *L. kluyveri* and
has fixed throughout all populations [16], we can not observe the original endogenous

1 sequences that have been replaced by the introgression. Overall, we predict that a
2 small number of low expression genes ($\phi < 1$) were weakly exapted at the time of the
3 introgression (Figure 5a). Thus, they appear to provide a small fitness advantage
4 due to the accordance of exogenous mutation bias with endogenous selection bias
5 (compare Figure S2 and S3). High expression genes ($\phi > 1$) are predicted to have
6 faced the largest selection against their mismatched codon usage in the novel cellular
7 environment. In order to account for differences in the efficacy of selection on codon
8 usage either due to the cost of pausing, differences in the effective population size,
9 or the decline in fitness with every ATP wasted between the donor lineage and *L.*
10 *kluyveri* we added a linear scaling factor κ to scale our estimates of $\Delta\eta$ between the
11 donor lineage and *L. kluyveri* and searched for the value that minimized the cost of
12 the introgression, thus giving us the best case scenario (See Methods for details).

13 Using our estimates of ΔM and $\Delta\eta$ from the endogenous genes and assuming the
14 current exogenous amino acid composition of genes is representative of the replaced
15 endogenous genes, we estimate the strength of selection against the exogenous genes
16 at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of
17 selection bias for the exogenous genes show that, while well correlated with the
18 endogenous genes, only nine amino acids share the same selectively preferred codon.
19 Exogenous genes are, therefore, expected to represent a significant reduction in
20 fitness for *L. kluyveri* due to mismatch in codon usage. Since $\Delta\eta$ is proportional
21 to the difference in fitness between the wild type and a mutant, we can use our
22 estimates of $\Delta\eta$ to approximate the selection against the exogenous genes Δs [10,
23 26]. We estimate that the selection against all exogenous genes due to mismatched
24 codon usage to have been $\Delta s \approx -0.0008$ at the time of the introgression and
25 ≈ -0.0003 today. This reduction in Δs is primarily due to adaptive changes to the
26 codon usage of the most highly expressed, introgressed genes (Figures 5a & S8).
27 Based on the selection against the codon mismatch at the time of the introgression
28 and assuming an effective population size N_e on the order of 10^7 [27], we estimate
29 a fixation probability of $(1 - \exp[-\Delta s])/(1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$ [26] for the
30 exogenous genes. Clearly, the possibility of fixation under this simple scenario is
31 effectively zero. In order for the exogenous genes to have reached fixation one or
32 more exogenous loci must have provided a selective advantage not considered in
33 this study (See Discussion).





1 Discussion

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

15 In contrast to other methods such as RSCU, CAI, or tAI, ROC SEMPPR does
16 not rely on external information such as gene expression or tRNA gene copy number
17 [5, 19]. Instead, ROC SEMPPR allows for the estimation of protein synthesis rate ϕ
18 and separates the effects of mutation and selection on codon usage. In addition, [20]
19 showed that approaches like CAI are sensitive to amino acid composition, another
20 property that distinguishes the endogenous and exogenous genes [15].

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

32 Previous studies showed that nucleotide composition can be strongly affected by
33 biased gene conversion, which, in turn would affect codon usage. Biased gene conver-

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

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

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

30 Previous work indicated that the donor lineage of the exogenous genes has to be
31 a, potentially unknown, Lachancea lineage [15–18]. These previous results, however,
32 are based on species rather than gene trees, ignoring the differential adaptation rate
33 to their novel cellular environment between genes or do not consider lineages outside

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

15
16 In terms of understanding the spread of the introgression, we calculated the ex-
17 pected selective cost of codon mismatch between the *L. kluyveri* and *E. gossypii*
18 lineages. Under our working hypothesis, the majority of the introgressed would have
19 imposed a selective cost due to codon mismatch. Nevertheless, $\sim 30\%$ of low expres-
20 sion exogenous genes ($\phi < 1$) appeared to be exapted at the time of the introgres-
21 sion. This exaptation is due to the mutation bias in the endogenous genes matching
22 the selection bias in the exogenous genes for GC ending codons. Our estimate of
23 the selective cost of codon mismatch on the order of -0.0008 . While this selective
24 cost may not seem very large, assuming *L. kluyveri* had a large N_e , the fixation
25 probability of the introgression is the astronomically small value of $\approx 10^{-6952} \approx 0$.
26 While this estimate heavily depends on the working hypothesis that the exogenous
27 genes originated from the *E. gossypii* lineage, we can also calculate the hypothetical
28 fixation probability if the current exogenous genes would introgress into *L. kluyveri*.
29 Our estimate of the current selective cost of the mismatch of codon usage is on the
30 order of -0.0003 . The fixation probability of the current exogenous genes would
31 still be astronomically small $\approx 10^{-2609} \approx 0$. These results are in accordance with
32 previous work, highlighting the necessity of codon usage compatibility between en-
33 dogenous and transferred exogenous genes [32, 33]. Thus, the basic scenario of an

1 introgression between two yeast species with large N_e and where the introgression
2 solely imposes a selective cost due to codon mismatch is clearly too simplistic.
3

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

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

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

5 Conclusion

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

19 Materials and Methods

20 Separating Endogenous and Exogenous Genes

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

30 Model Fitting with ROC SEMPPR

31 ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) [22] and R (3.4.1)
32 [35]. ROC SEMPPR was run from 10 different starting values for at least 250,000
33 iterations and thinned to every 50th iteration. After manual inspection to verify that

1 the MCMC had converged, parameter posterior means, log posterior probability and
2 log likelihood were estimated from the last 500 samples (last 10% of samples).
3

4 **Model selection**
5

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

13 **Comparing Codon Specific Parameter Estimates and Selecting Candidate lineages**
14

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

$$\Delta M_i = \Delta M_{i,1} - \overline{\Delta M_i} \quad (1) \quad 18$$

$$\Delta \eta_i = \Delta \eta_{i,1} - \overline{\Delta \eta_i} \quad (2) \quad 19$$

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

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

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 assess 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 ΔM as

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

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

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

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

$$3 \\ 4 c_1(t) = \frac{1 + \exp[-X](K - 1)}{1 + \Delta M_{\text{endo}}} \quad (5) \\ 5$$

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

7 Equation 5 was solved with a mutation rate $\mu_{2,1}$ of 3.8×10^{-10} per nucleotide per
8 generation [41]. Current codon frequencies for each codon family were taken from
9 our estimates of ΔM from the exogenous genes. Mathematica (11.3) [42] was used
10 to calculate the time t_{intro} it takes for the initial codon frequencies $c_1(0)$ for each
11 codon family to equal the current exogenous codon frequencies. The same equation
12 was used to determine the time t_{decay} at which the signal of the exogenous cellular
13 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 ϕ
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 κ . As $\Delta\eta$ is defined as $\Delta\eta =$
33 $2N_e q(\eta_i - \eta_j)$, we cannot distinguish if κ is a scaling on protein synthesis rate ϕ ,

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, κ 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, ϕ_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 ϕ , 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 **Acknowledgments**
6 The authors would like to thank Alexander Cope for helpful criticisms and suggestions for this work.
7
- 8 **Availability of data and materials**
9 Parameter estimates generated during this study are available from the corresponding author. All remaining data
10 generated during this study are included in this published article as figures, tables.
11
- 12 **Authors' contributions**
13 CL and MAG initiated the study. CL collected and analyzed the data and wrote the manuscript. MAG and BCO
14 edited the manuscript. CL, MAG, BCO, and RZ contributed to the data analysis and acquiring of funding. All
15 Authors approved the final manuscript.
16
- 17 **Funding**
18 This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ), MCB-1546402 (A. Von Arnim and
19 MAG), and DEB-1355033 (BCO, MAG, and RZ) with additional support from Department of Ecology &
20 Evolutionary Biology (EEB) at the University of Tennessee Knoxville (UTK) and the National Institute for
21 Mathematical and Biological Synthesis (NIMBioS), an Institute sponsored by the National Science Foundation
22 through NSF Award DBI-1300426. CL received support as a Graduate Student Fellow from NIMBioS with
23 additional support from Departments of Mathematics and EEB at UTK.
24
- 25 **Ethics approval and consent to participate**
26 Not applicable
27
- 28 **Consent for publication**
29 Not applicable
30
- 31 **Competing interests**
32 The authors declare that they have no competing interests.
33
- 34 **Author details**
35 ¹Department of Ecology & Evolutionary Biology, University of Tennessee, 37996, Knoxville, TN, USA. ²National
36 Institute for Mathematical and Biological Synthesis, 37996, Knoxville, TN, USA. ³Max-Planck Institute of
37 Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, 01307, Dresden, Germany. ⁴Department of Business
38 Analytics and Statistics, University of Tennessee, 37996, Knoxville, TN, USA.
39
- 40 **References**
1. Gouy, M., Gautier, C.: Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Research* **10**, 7055–7074 (1982)
 2. Ikemura, T.: Codon usage and tRNA content in unicellular and multicellular organisms. *Molecular Biology and Evolution* **2**, 13–34 (1985)
 3. Bulmer, M.: The selection-mutation-drift theory of synonymous codon usage. *Genetics* **129**, 897–907 (1990)
 4. 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)
 5. 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)
 6. Wright, F.: The 'effective number of codons' used in a gene. *Genet* **87**, 23–29 (1990)
 7. 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)

- 1 1365–1374 (2005)

2 28. Nagylaki, T.: Evolution of a finite population under gene conversion. Proc. Natl. Acad. Sci. U. S. A. **80**,
6278–6281 (1983)

3 29. Nagylaki, T.: Evolution of a large population under gene conversion. Proc. Natl. Acad. Sci. U. S. A. **80**,
5941–5945 (1983)

4 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)

5 31. Salichos, L., Rokas, A.: Inferring ancient divergences requires genes with strong phylogenetic signals. Nature
497, 327–331 (2013)

6 32. Medrano-Soto, A., Moreno-Hagelsieb, G., Vinuesa, 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)

7 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

8 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)

9 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/>

10 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)

11 37. Legendre, P.: Lmodel2: Model II Regression. (2018). R package version 1.7-3.
<https://CRAN.R-project.org/package=lmodel2>

12 38. Soderlund, C., Nelson, W., Shoemaker, A., Paterson, A.: Symap: A system for discovering and viewing synteny
regions of fpc maps. Genome Research **16**, 1159–1168 (2006)

13 39. Soderlund, C., Bornhoff, M., Nelson, W.: Symap v3.4: a turnkey synteny system with application to plant
genomes. Nucleic Acids Research **39**(10), 68 (2011)

14 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)

15 41. Lang, G.I., Murray, A.W.: Estimating the per-base-pair mutation rate in the yeast *Saccharomyces cerevisiae*.
Genetics **178**(1), 67–82 (2008)

16 42. Wolfram Research Inc.: Mathematica 11. (2017). <http://www.wolfram.com>

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

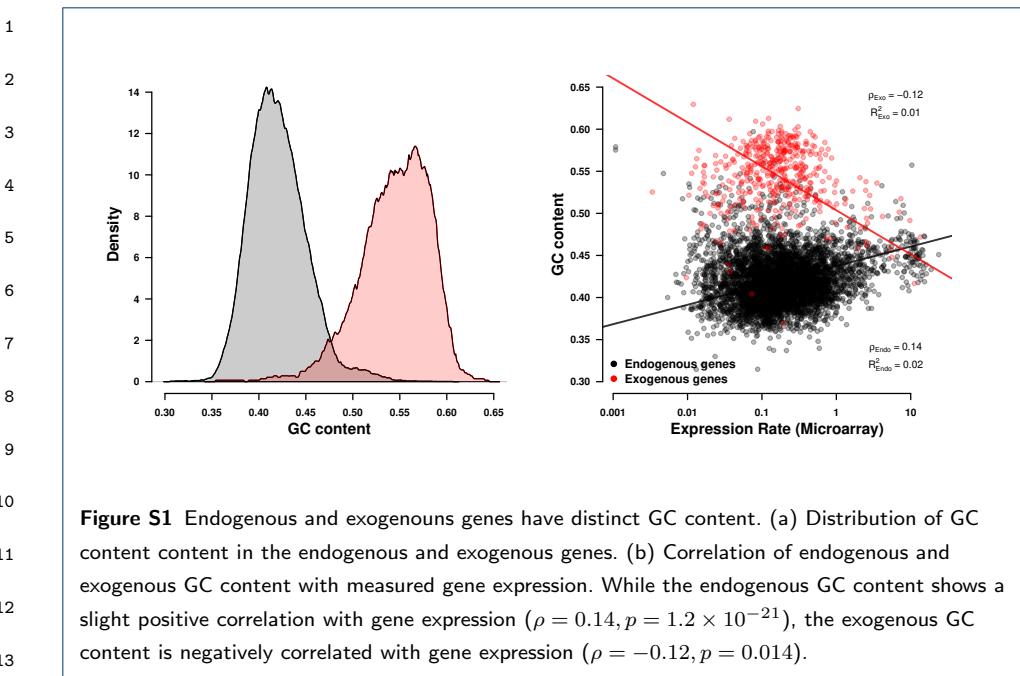
1 **Supplementary Material**

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 ΔM . 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 ₄ 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 ₂ 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	Table S2 Synonymous selection codon preference based on our estimates of $\Delta\eta$. Shown are the most likely codon in high expression genes for each amino acid in: <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 ₄ 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 ₂ Z	AGT	AGT	AGT	AGT	
33						33



1	1
2	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

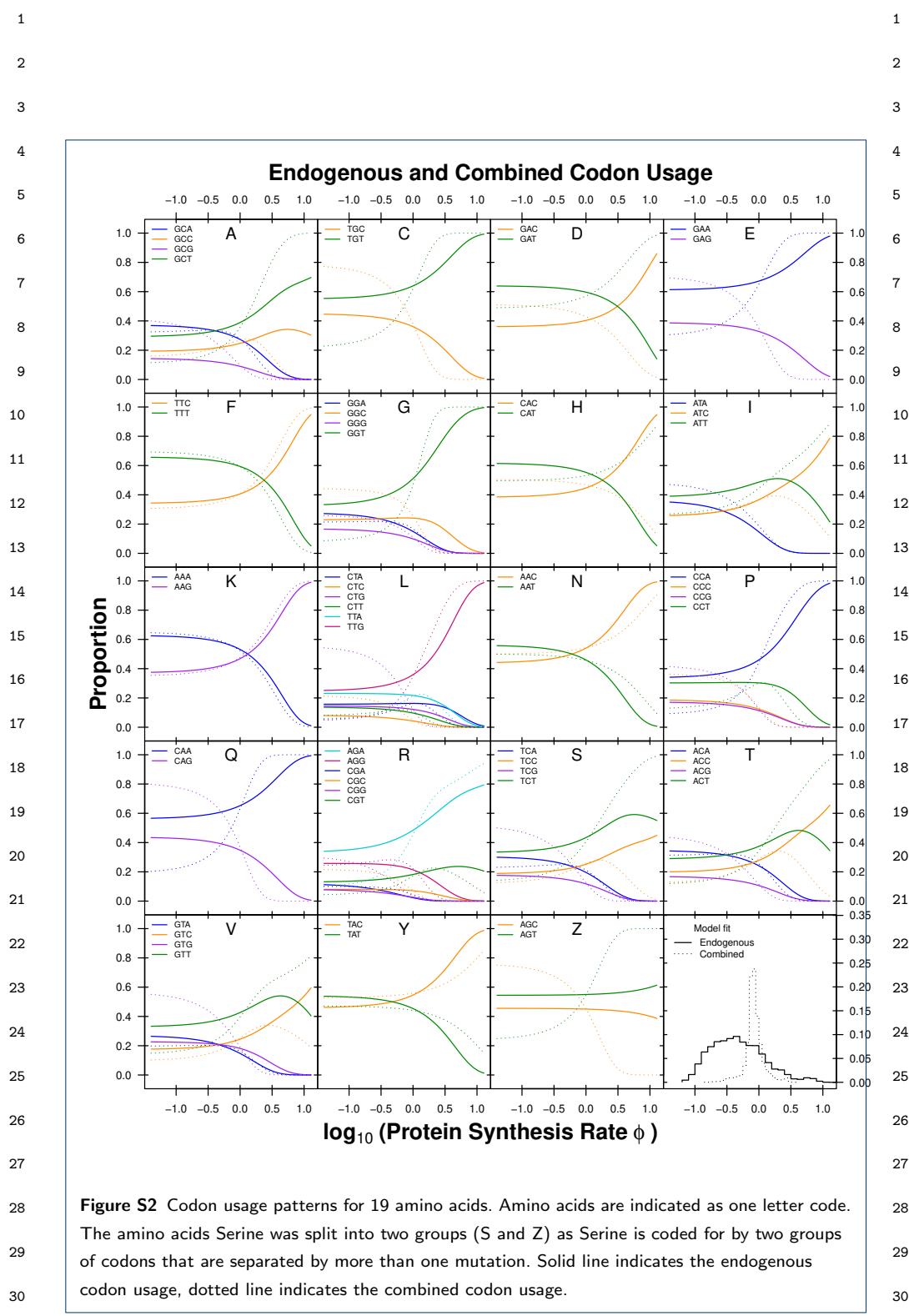


Figure S2 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dotted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

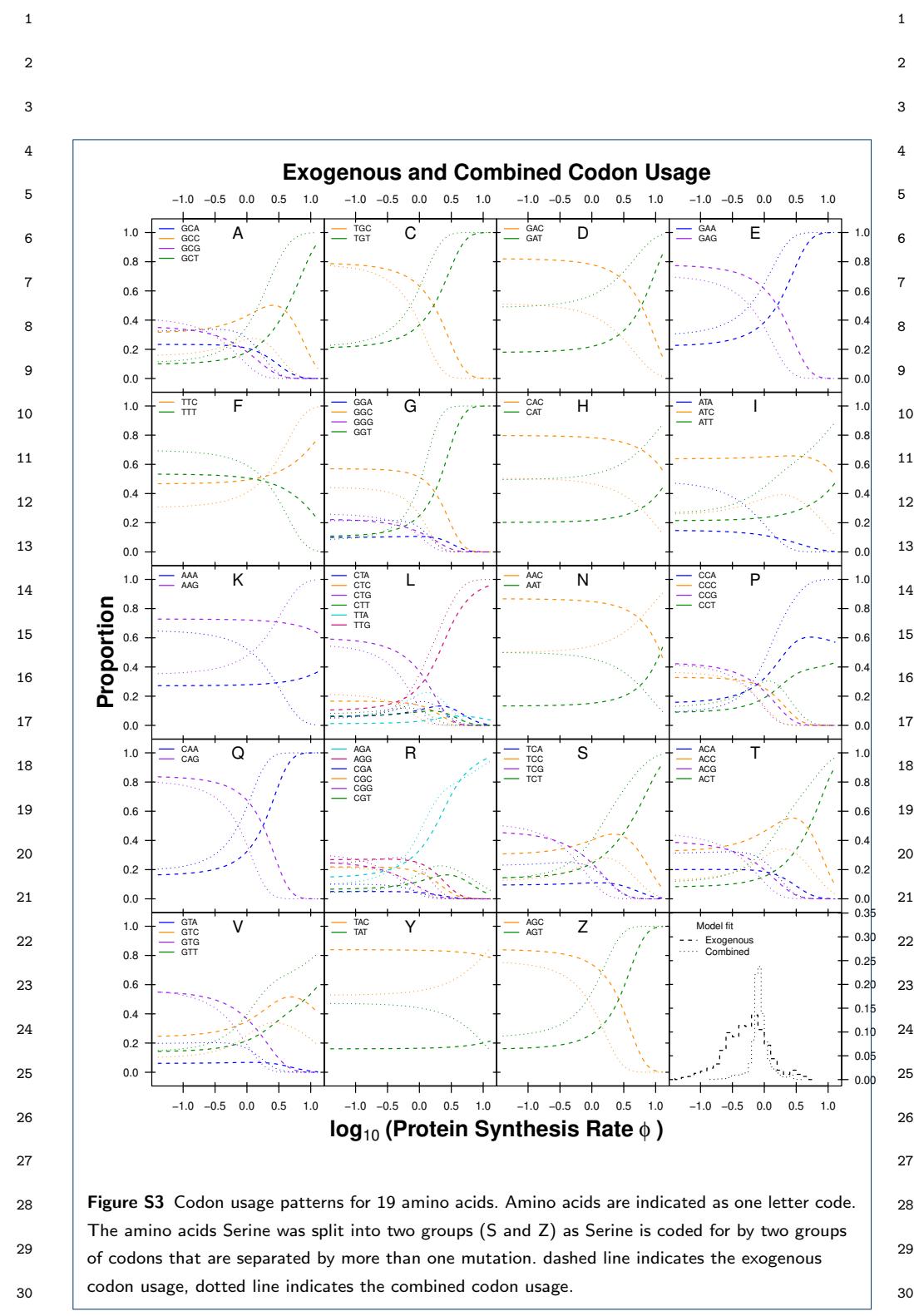
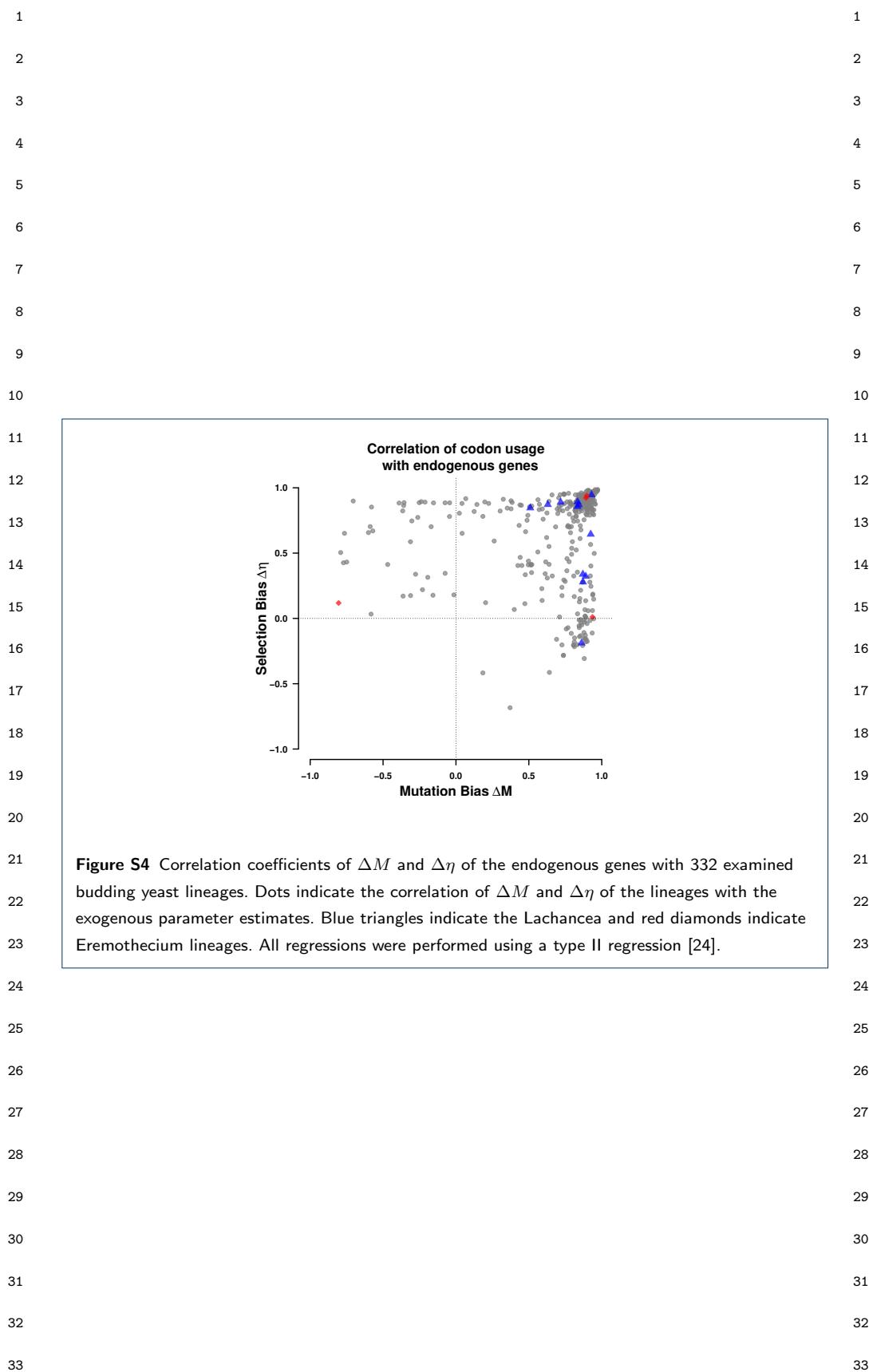
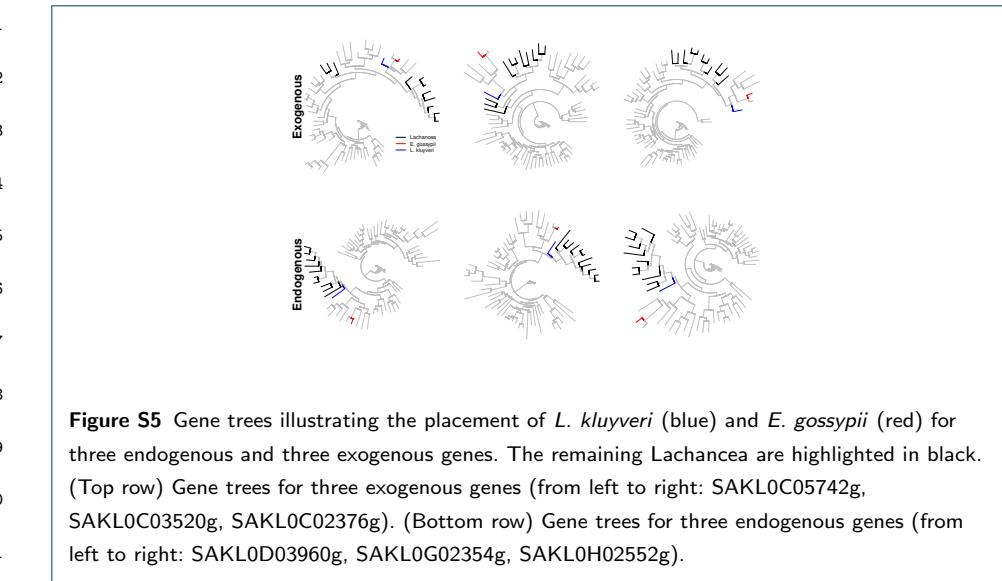
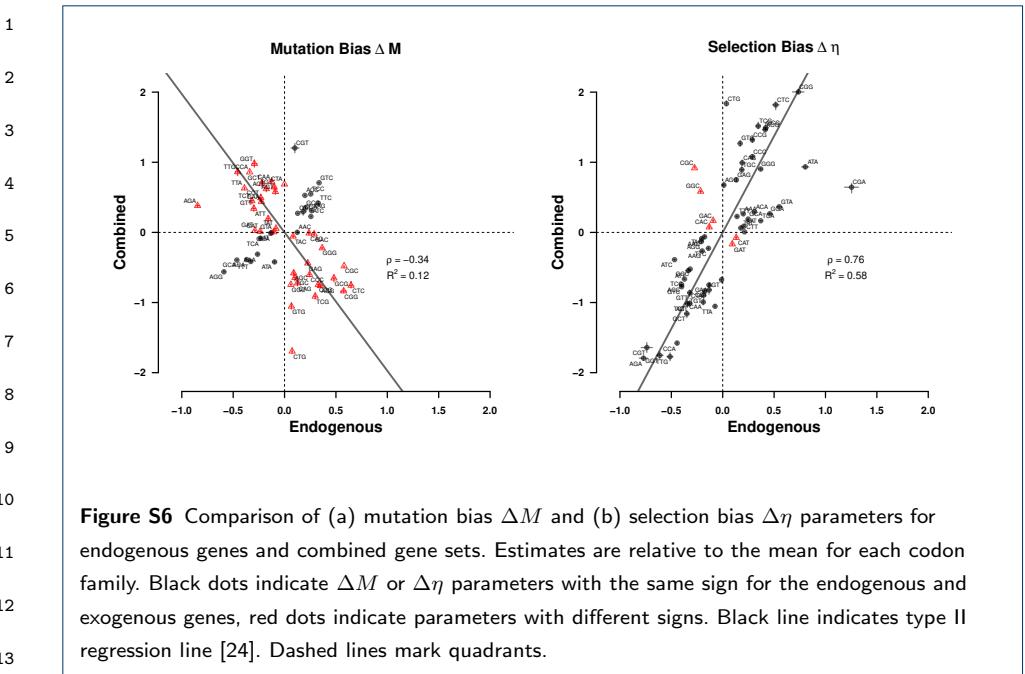
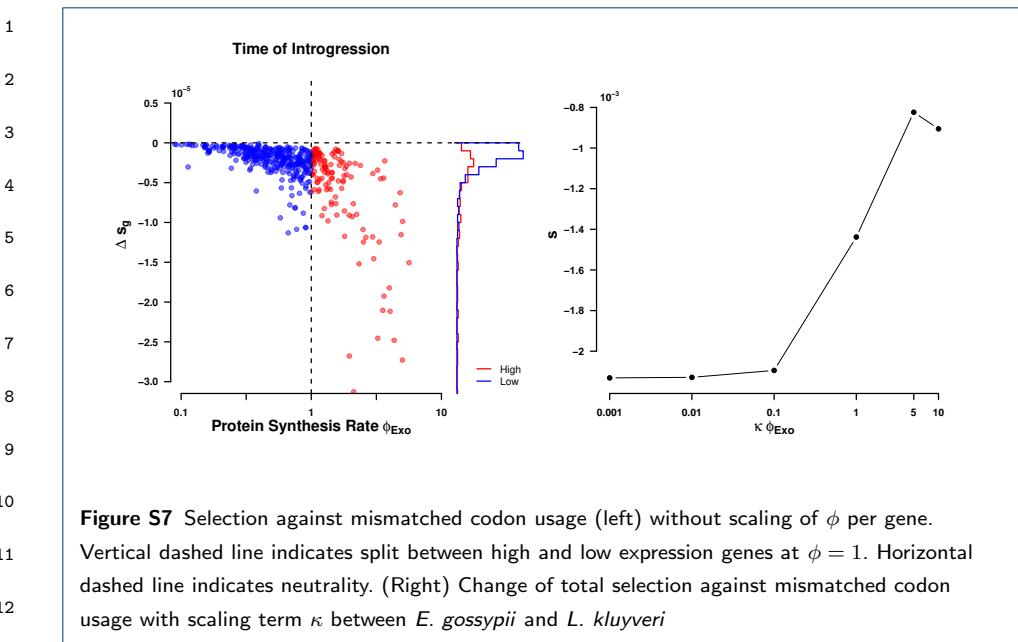


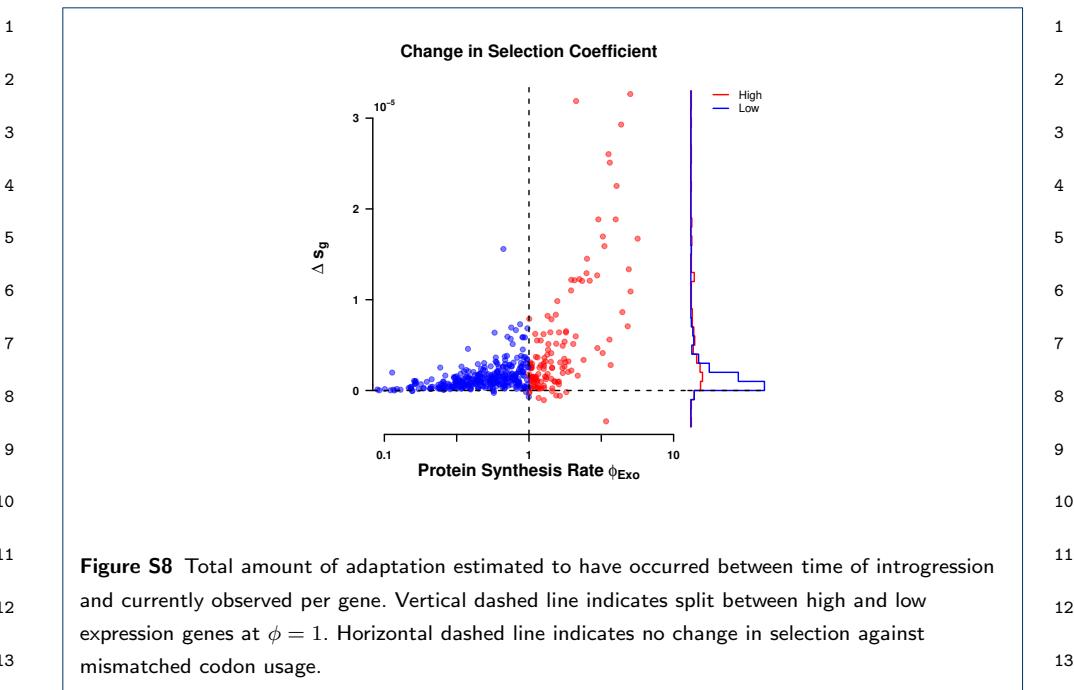
Figure S3 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. dashed line indicates the exogenous codon usage, dotted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











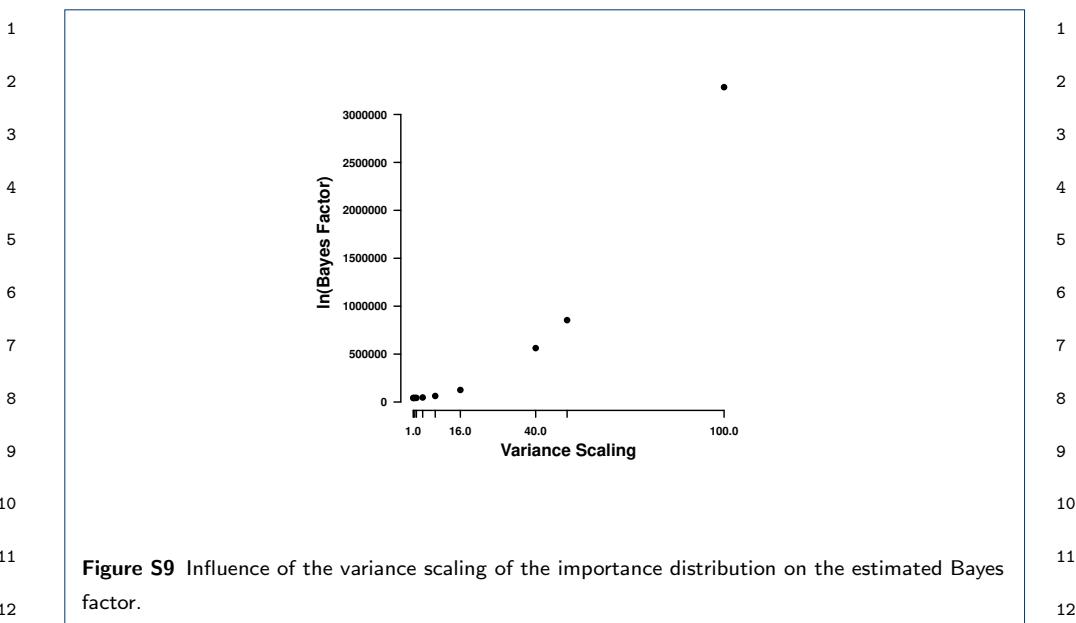


Figure S9 Influence of the variance scaling of the importance distribution on the estimated Bayes factor.