

RESEARCH

2
3
4
5
6
7
8
**Unlocking a signal of introgression from codons
in *Lachancea kluyveri* using a mutation-selection
model**

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

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

Correspondence:
anderer@mpi-cbg.de
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. In order to better understand the effects of introgression on codon usage patterns and vice versa, we examine the patterns of codon usage in *Lachancea kluyveri*, a yeast which has experienced a large introgression. We quantify the effects of mutation bias and selection for translation efficiency on the codon usage pattern of the endogenous and introgressed exogenous genes using a Bayesian mixture model, ROC SEMPPR, which is built on mechanistic assumptions about protein synthesis and grounded in population genetics.

Results: We find substantial differences in codon usage between the endogenous and exogenous genes, and show that these differences can be largely attributed to differences in mutation bias favoring A/T ending codons in the endogenous genes while favoring C/G ending codons in the exogenous genes. Recognizing the two different signatures of mutation bias and selection improves our ability to predict protein synthesis rate by 42% and allowed us to accurately assess the decaying signal of endogenous codon mutation and preferences. In addition, using our estimates of mutation bias and selection, we identify *Eremothecium gossypii* as the closest relative to the exogenous genes, providing an alternative hypothesis about the origin of the exogenous genes, estimate that the introgression occurred $\sim 6 \times 10^8$ generation ago, and estimate its historic and current selection against mismatched codon usage.

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

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

¹Background

²Synonymous codon usage patterns varies within a genome and between taxa, re-²
³flecting differences in mutation bias, selection, and genetic drift. The signature of³
⁴mutation bias is largely determined by the organism's internal or cellular environ-⁴
⁵ment, such as their DNA repair genes or UV exposure. While this mutation bias⁵
⁶is an omnipresent evolutionary force, its impact can be obscured or amplified by⁶
⁷selection. The signature of selection on codon usage is largely determined by an or-⁷
⁸ganism's cellular environment alone, such as, but not limited to, its tRNA species,⁸
⁹their copy number, and their post-transcriptional modifications. In general, the⁹
¹⁰strength of selection on codon usage is assumed to increase with its expression level¹⁰
¹¹[1–3], specifically its protein synthesis rate [4]. Thus as protein synthesis increases,¹¹
¹²codon usage shifts from a process dominated by mutation to a process dominated¹²
¹³by selection. The overall efficacy of mutation and selection on codon usage is a¹³
¹⁴function of the organism's effective population size N_e . ROC SEMPPR allows us¹⁴
¹⁵to disentangle the evolutionary forces responsible for the patterns of codon usage¹⁵
¹⁶bias [5–7] (CUB) encoded in an species' genome, by explicitly modeling the com-¹⁶
¹⁷bined evolutionary forces of mutation, selection, and drift [4, 8–10]. In turn, these¹⁷
¹⁸evolutionary parameters should provide biologically meaningful information about¹⁸
¹⁹the lineage's historical cellular and external environment.¹⁹

²⁰ Most studies implicitly assume that the CUB of a genome is shaped by a single²⁰
²¹cellular and external environment. However, this assumption is clearly violated to²¹
²²increasing degrees via horizontally gene transfer, large scale introgressions, and hy-²²
²³brid specie formation. In these scenarios, one would expect to see the signature of²³
²⁴multiple cellular environments in a genome's CUB [11, 12]. Indeed, differences in²⁴
²⁵CUB between linages have been proposed to have a major effect on their rates of²⁵
²⁶gene transfer with rates declining with differences in their CUB. On a more practical²⁶
²⁷level, if differences in codon usage of transferred genes are not taken into account²⁷
²⁸for, they may distort the interpretation of codon usage patterns. Such distortion²⁸
²⁹could lead to the wrong inference of codon preference for an amino acid [8, 10], un-²⁹
³⁰derestimate the variation in protein synthesis rate, or distort estimates of mutation³⁰
³¹bias when analyzing a genome.³¹

³² To illustrate these ideas, we analyze the CUB of the genome of the yeast *Lachancea*³²
³³*kluyveri* using ROC SEMPPR, a population genetics based model of synonymous³³

¹codon usage evolution that accounts for and, in turn, can estimate the contribution¹
²of mutation bias ΔM , selection bias. The mathematics of ROC SEMPPR are de-²
³rived on a mechanistic description of ribosome movement along an mRNA, although³
⁴the approximation of other biological mechanisms could also be consistent with the⁴
⁵model. Broadly speaking, ROC SEMPPR allows us to quantify the cellular environ-⁵
⁶ment in which genes have evolved by separately estimating the effects of mutation⁶
⁷bias and selection bias on codon usageDE between synonymous codons and pro-⁷
⁸tein synthesis rate ϕ to the patterns of codon usage observed within a set of genes.⁸
⁹Briefly, the set of ΔM for an amino acid quantifies the relative differences in muta-⁹
¹⁰tional stability or bias between the synonymous codons of the amino acid S . In the¹⁰
¹¹absence of selection bias (or equivalently when gene expression $\phi = 0$), the equilib-¹¹
¹²rium frequency of synonymous codon i is simply $\exp[-\Delta M_i] / \left(\sum_{j \in S} \exp[-\Delta M_j] \right)$.¹²
¹³Because the time units of protein production rate have no intrinsic time scale, we¹³
¹⁴define the average protein production rate for a set of genes to be one, i.e. $\bar{\phi} = 1$ ¹⁴
¹⁵by definition [10]. In order to facilitate comparisons between gene sets, we express¹⁵
¹⁶both, ΔM and $\Delta \eta$, as deviation from the mean of each synonymous codon family¹⁶
¹⁷(see Materials and Methods for details). Nevertheless, the difference $\Delta \eta$ describes¹⁷
¹⁸the difference in fitness between two synonymous codons relative to drift for a gene¹⁸
¹⁹whose protein production rate ϕ is equal to the the average rate of protein produc-¹⁹
²⁰tion $\bar{\phi}$ across the set of genes. In other words, for a gene whose protein is expressed²⁰
²¹at the average rate, for any two given synonymous codons i and j , $\Delta \eta_i - \Delta \eta_j = N_e s$.²¹
22

²³ The *Lachancea* clade diverged from the *Saccharomyces* clade, prior to its whole²³
²⁴genome duplication ~ 100 Mya ago [13, 14]. Since that time, *L. kluyveri*, which is²⁴
²⁵sister species to all other *Lachancea* spp., has experienced a large introgression of²⁵
²⁶exogenous genes (1 Mb, 457 genes) which is found in all of its populations [15, 16],²⁶
²⁷but in no other known *Lachancea* species [17]. The introgression replaced the left²⁷
²⁸arm of the C chromosome and displays a 13% higher GC content than the en-²⁸
²⁹dogenous *L. kluyveri* genome [15, 16]. Previous studies suggest that the source of²⁹
³⁰the introgression is probably a currently unknown or potentially extinct *Lachancea*³⁰
³¹lineage based on gene concatenation or synteny relationships [15–18]. These char-³¹
³²acteristics make *L. kluyveri* an ideal model to study the effects of an introgressed³²
³³cellular environment and the resulting mismatch in codon usage.³³

¹ While previous studies [8, 9] have used information on gene expression to separate¹
² the effects of mutation and selection on codon usage, ROC SEMPPR does not²
³ need such information but can provide it. ROC SEMPPR's resulting predictions³
⁴ of protein synthesis rates have been shown to be on par with laboratory measurements⁴
⁵ [8, 10]. In contrast to often used heuristic approaches to study codon usage⁵
⁶ [5, 6, 19], ROC SEMPPR explicitly incorporates and distinguishes between mu-⁶
⁷tation and selection effects on codon usage and properly weights its estimates by⁷
⁸ amino acid usage [20]. We use ROC SEMPPR to separately describe the two cellular⁸
⁹ environments reflected in the *L. kluyveri* genome; the signature of the endogenous⁹
¹⁰ environment reflected in the larger set of non-introgressed genes and the decaying¹⁰
¹¹ signature of the ancestral, exogenous environment in the smaller set of introgressed¹¹
¹² genes. Our results indicate that the current difference in GC content between en-¹²
¹³dogenous and exogenous genes is mostly due to the differences in mutation bias¹³
¹⁴ ΔM of their respective cellular environments. Taking the different signatures of¹⁴
¹⁵ ΔM and selection bias $\Delta \eta$ of the endogenous and exogenous sets of genes substan-¹⁵
¹⁶tially improves our ability to predict present day protein synthesis rates ϕ . These¹⁶
¹⁷ endogenous and exogenous gene set specific estimates of ΔM and $\Delta \eta$, in turn, allow¹⁷
¹⁸ us to address more refined biological questions. For example, we find support for¹⁸
¹⁹ an alternative origin of the exogenous genes and identify *E. gossypii* as the nearest¹⁹
²⁰ sampled relative of the source of the introgressed genes out of the 332 budding yeast²⁰
²¹ lineages with sequenced genomes [21]. While this inference is in contrast to previous²¹
²² work [15–18], we find additional phylogenetic support for via gene tree reconstruc-²²
²³ tion and gene synteny. We also estimate the age of the introgression to be on the²³
²⁴ order of 0.2 - 1.7 Mya, estimate the selection against these genes, both at the time²⁴
²⁵ of introgression and now, and predict a detectable signature of CUB to persist in²⁵
²⁶ the introgressed genes for another 0.3 - 2.8 Mya, highlighting the sensitivity of our²⁶
²⁷ approach.²⁷

28

28

²⁹Results

29

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

30

³¹We used our software package AnaCoDa [22] to compare model fits of ROC³¹
³² SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two³²
³³ sets of 4,864 endogenous and 497 exogenous genes. These two set where initially³³

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

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

7

7

8

8

9 identified based on their striking difference in GC content [15], with very little over-
⁹lap in GC content between the two sets (Figure S1a). ROC SEMPPR is a statistical
¹⁰model that relates the effects of mutation bias ΔM , selection bias $\Delta\eta$ between syn-
¹¹onymous codons and protein synthesis rate ϕ , to explain the observed codon usage₁₂
¹²patterns. Thus, the probability of observing a synonymous codon is proportional₁₃
¹³to $p \propto \exp(-\Delta M - \Delta\eta\phi)$ [10]. Briefly, ΔM describes the mutation bias between₁₄
¹⁴two synonymous codons at stationarity under a time reversible mutation model.₁₅
¹⁵Because ROC SEMPPR only considers the stationary probabilities, only variation₁₆
¹⁶in mutation bias, not absolute mutation rates can be detected. $\Delta\eta$ describes the₁₇
¹⁷fitness difference between two synonymous codons relative to drift [10]. Since $\Delta\eta$ is₁₈
¹⁸scaled by protein synthesis rate ϕ , this term is dominant in highly expressed genes₁₉
¹⁹and tends towards 0 in low expression genes, allowing us to separate the effect of₂₀
²⁰mutation bias and selection bias on codon usage. We express both, ΔM and $\Delta\eta$,₂₁
²¹as deviation from the mean of each synonymous codon family which prevents that₂₂
²²the choice of the reference codon affects our results (see Materials and Methods for₂₃
²³details).₂₄

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

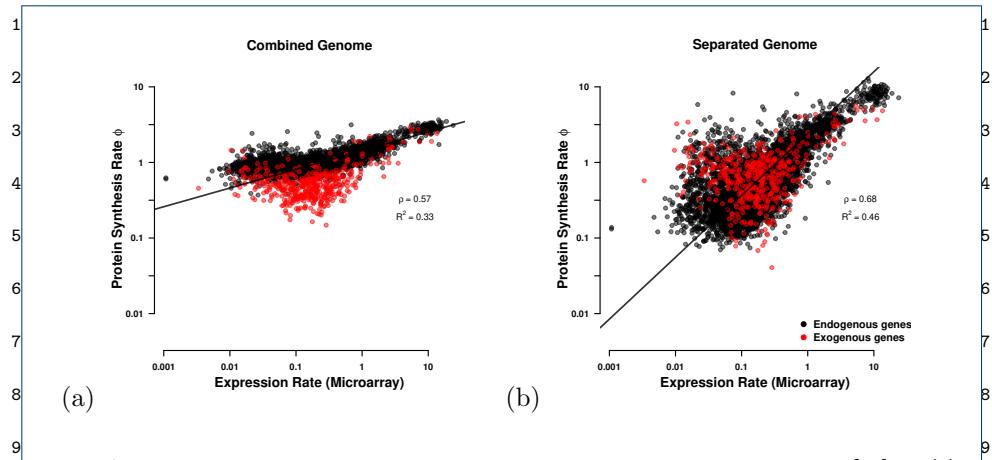


Figure 1 Comparison of predicted protein synthesis rate ϕ to mRNA abundance from [23] for (a) the combined genome where mutation bias and selection bias parameters ΔM and $\Delta \eta$ are estimated for the combined endogenous and exogenous gene sets, and (b) where ΔM and $\Delta \eta$ are estimated 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].

13

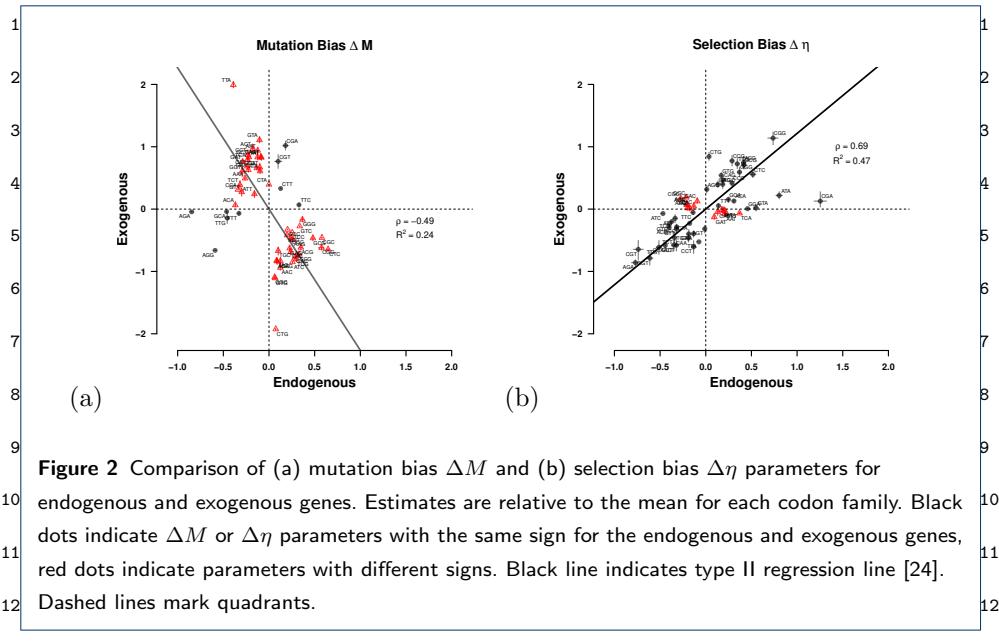
13

14 Comparing Differences in the Endogenous and Exogenous Codon Usage

14

15 Because ROC SEMPPR defines $\bar{\phi} = 1$, it makes the interpretation of $\Delta \eta$ as selection¹⁵
 16 on codon usage of the average gene with $\phi = 1$ straightforward and gives us the¹⁶
 17 ability to compare the efficacy of selection sN_e across genomes. While it may be¹⁷
 18 expected for the endogenous and exogenous genes to differ in their codon usage¹⁸
 19 pattern due to the large difference in GC content it is not clear how much of¹⁹
 20 this difference is due to differences in the mutation bias ΔM or selection bias $\Delta \eta$ ²⁰
 21 between the gene sets. To better understand the differences in the endogenous and²¹
 22 exogenous cellular environments, we compared our parameter estimates of ΔM ²²
 23 and $\Delta \eta$ for the two sets of genes. Our estimates of ΔM for the endogenous and²³
 24 exogenous genes were negatively correlated ($\rho = -0.49$, $p = 3.56 \times 10^{-5}$), indicating²⁴
 25 weak similarity with only $\sim 5\%$ of the codons share the same sign between the two²⁵
 26 mutation environments (Figure 2a). Overall, the endogenous genes only show a²⁶
 27 selection preference for C and G ending codons in $\sim 58\%$ of the codon families.²⁷
 28 In contrast, the exogenous genes display a strong preference for A and T ending²⁸
 29 codons in $\sim 89\%$ of the codon families.²⁹

30 For example, the endogenous genes show a mutational bias for A and T ending³⁰
 31 codons in $\sim 95\%$ of the codon families. The exception is Leucine (Leu, L), where³¹
 32 mutation appears to favor the codon TTG over TTA (Figure 3, Table S1). The³²
 33 exogenous genes display an equally consistent mutational bias towards C and G³³



¹⁴ending codons (the exception being Phe, F). In contrast to ΔM , our estimates of $\Delta \eta^{14}$ ¹⁵for the endogenous and exogenous genes were positively correlated ($\rho = 0.69$) and¹⁵ ¹⁶showing the same sign in ~ 53% of codons between the two selection environments¹⁶ ¹⁷(Figure 2).

¹⁸ We find that the signature of selection bias $\Delta \eta$ also differs substantially between¹⁸ ¹⁹the endogenous and exogenous gene sets. The difference in codon usage between¹⁹ ²⁰endogenous and exogenous genes is striking as the sign for $\Delta \eta$ changes, indicating a²⁰ ²¹change in codon preference for some amino acids. As a result, our estimates of the²¹ ²²optimal codon differ in nine cases between endogenous and exogenous genes (Figure²² ²³, Table S2). For example, the usage of the Asparagine (Asn, N) codon AAC is²³ ²⁴increased in highly expressed endogenous genes but the same codon is depleted in²⁴ ²⁵highly expressed exogenous genes. For Aspartic acid (Asp, D), the combined genome²⁵ ²⁶shows the same codon preference in highly expressed genes as the exogenous gene²⁶ ²⁷set. Generally, fits to the complete *L. kluyveri* genome reveal that the relatively²⁷ ²⁸small exogenous gene set (~ 10% of genes) has a disproportionate effect on the²⁸ ²⁹model fit (Figure S2, S3).

³⁰ Of the nine cases in which the endogenous and exogenous genes show differences³⁰ ³¹in the selectively most favored codon five cases (Asp, D; His, H; Lys, K; Asn, N;³¹ ³²and Pro, P) the endogenous genes favor the codon with the most abundant tRNA.³² ³³For the remaining four cases (Ile, I; Ser, S; Thr, T; and Val, V), there are no³³

¹tRNA genes for the wobble free cognate codon encoded in the *L. kluyveri* genome.¹
²However, the codon preference of these four amino acids in the exogenous genes²
³matches the most abundant tRNA encoded in the *L. kluyveri* genome. In contrast³
⁴to ΔM , our estimates of selection bias $\Delta\eta$ for the endogenous and exogenous genes⁴
⁵are positively correlated ($\rho = 0.69$, $p = 9.76 \times 10^{-10}$) and show the same sign in⁵
⁶~53% of the cases (Figure 2).⁶

⁷ This striking difference in codon usage was noted previously. For example, using⁷
⁸RSCU [5], GAA (coding for Glu, E) was identified as the optimal synonymous codon⁸
⁹in the whole genome and GAG as the optimal codon in the exogenous genes [15].⁹
¹⁰Our results, however, indicate that GAA is the optimal codon in both, endogenous¹⁰
¹¹and exogenous genes, and that the high RSCU in the exogenous genes of GAG is¹¹
¹²driven by mutation bias (Table S1 and S2). Similar effects are observed for other¹²
¹³amino acids.¹³

¹⁴ The effect of the small exogenous gene set on the fit to the complete *L. kluyveri*¹⁴
¹⁵genome is smaller for our estimates of selection bias $\Delta\eta$ than ΔM , but still large.¹⁵
¹⁶ We find that the complete *L. kluyveri* genome is estimated to share the selectively¹⁶
¹⁷preferred codon with the exogenous genes in ~60% of codon families that show dis-¹⁷
¹⁸similarity between endogenous and exogenous genes. We also find that the complete¹⁸
¹⁹*L. kluyveri* genome fit shares mutationally preferred codons with the exogenous¹⁹
²⁰genes in ~78% of the 19 codon families showing a difference in mutational codon²⁰
²¹preference between the endogenous and exogenous genes. In two cases, Isoleucine²¹
²²(Ile, I) and Arginine (Arg, R), the strong dissimilarity in mutation preference results²²
²³in an estimated codon preference in the complete *L. kluyveri* genome that differs²³
²⁴from both the endogenous, and the exogenous genes. These results clearly show that²⁴
²⁵it is important to recognize the difference in endogenous and exogenous genes and²⁵
²⁶treat these genes as separate sets to avoid the inference of incorrect synonymous²⁶
²⁷codon preferences and better predict protein synthesis.²⁷

²⁸

²⁹Can Codon Usage Help Determine the Source of the Exogenous Genes²⁹

³⁰ Since the origin of the exogenous genes is currently unknown, we explored if the³⁰
³¹information on codon usage extracted from the exogenous genes can be used to³¹
³²identify a potential source lineage. We combined our estimates of mutation bias³²
³³ ΔM and selection bias $\Delta\eta$ with synteny information and searched for potential³³

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

source lineages of the introgressed exogenous region. We used ΔM to identify candidate lineages as the endogenous and exogenous genes show greater dissimilarity in mutation bias than in selection bias. We examined 332 budding yeasts [21] and, identified the ten lineages with the highest correlation to the exogenous ΔM parameters as potential source lineages (Figure 4, Table 2). Two of the ten candidate lineages utilize the alternative yeast nuclear code (NCBI codon table 12). In this case, the codon CTG codes for Serine instead of Leucine. We therefore excluded the Leucine codon family from our comparison of codon families; however, there was no need to exclude Serine as CTG is not a one step neighbor of the remaining Serine codons. A mutation between CTG and the remaining Serine codons would require two mutations with one of them being non-synonymous, which would violate the weak mutation assumption of ROC SEMPPr.

The endogenous *L. kluyveri* genome exhibits codon usage very similar to most 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

¹relationship between the exogenous region and other yeasts appears to be limited¹
²to Saccharomycetaceae clade. Given these results, we conclude that, of the 332²
³examined yeast lineages the *E. gossypii* lineage is the most likely source of the in-³
⁴trogressed exogenous genes. Previous studies which studied the exogenous genes and⁴
⁵chromosome recombination in the Lachancea clade concluded that the exogenous⁵
⁶region originated from within the Lachancea clade, from an unknown or potentially⁶
⁷extinct lineage [15–17]. While it is not possible for us to dispute this hypothesis,⁷
⁸our results provide a novel hypothesis about the origin of the exogenous genes.⁸

⁹ To further test the plausibility of *E. gossypii* as potential source lineage, we iden-⁹
¹⁰tified 127 genes in our dataset [21] with homologous genes in *E. gossypii* and other¹⁰
¹¹Lachancea and used IQTree [25] to infer the phylogenetic relationship of the exoge-¹¹
¹²nous genes. Our results show that at least ~ 45% of exogenous genes (57/127) are¹²
¹³more closely related to *E. gossypii* than to other Lachancea S5. Interestingly, our re-¹³
¹⁴sults also indicate that codon usage does not necessarily correlate with phylogenetic¹⁴
¹⁵distance (Table 2).¹⁵

16

16

¹⁷Estimating Introgression Age

¹⁷

¹⁸If we assume that the exogenous genes originated from the *E. gossypii* lineage, we¹⁸
¹⁹can estimate the age of the introgression based on our estimates of mutation bias¹⁹
²⁰ ΔM . We modeled the change in codon frequency over time as exponential decay,²⁰
²¹and estimated the age of the introgression assuming that *E. gossypii* still represents²¹
²²the mutation bias of its ancestral source lineage at the time of the introgression and²²
²³a constant mutation rate. We infer the age of the introgression to be on the order²³
²⁴of $6.2 \pm 1.2 \times 10^8$ generations. Assuming *L. kluyveri* experiences between one and²⁴
²⁵eight generations per day, we estimate the introgression to have occurred between²⁵
²⁶212,000 to 1,700,000 years ago. Our estimate places the time of the introgression²⁶
²⁷earlier than the previous estimate of 19,000 - 150,000 years by [16].²⁷

²⁸Using our model of exponential decay model, we also estimated the persistence of²⁸
²⁹the signal of the exogenous cellular environment. We predict that the ΔM signal of²⁹
³⁰the source cellular environment will have decayed to be within one percent of the³⁰
³¹*L. kluyveri* environment in $\sim 5.4 \pm 0.2 \times 10^9$ generations, or between 1,800,000 and³¹
³²15,000,000 years. Together, these results indicate that the mutation signature of³²
³³the exogenous genes will persist for a very long time.³³

¹Estimating Selection against Codon Mismatch of the Exogenous Genes

²We define the selection against inefficient codon usage as the difference between the ² fitness on the log scale of an expected, replaced endogenous gene and the exogenous ³ gene, $s \propto \phi\Delta\eta$ due to the mismatch in codon usage parameters (See Methods for ⁴ details). As the introgression occurred before the diversification of *L. kluyveri* and ⁵ has fixed throughout all populations [16], we can not observe the original endogenous ⁶ sequences that have been replaced by the introgression. Overall, we predict that a ⁷ small number of low expression genes ($\phi < 1$) were weakly exapted at the time of the ⁸ introgression (Figure 5a). Thus, they appear to provide a small fitness advantage ⁹ due to the accordance of exogenous mutation bias with endogenous selection bias ¹⁰ (compare Figure S2 and S3). High expression genes ($\phi > 1$) are predicted to have ¹¹ faced the largest selection against their mismatched codon usage in the novel cellular ¹² environment. In order to account for differences in the efficacy of selection on codon ¹³ usage either due to the cost of pausing, differences in the effective population size, ¹⁴ or the decline in fitness with every ATP wasted between the donor lineage and *L.* ¹⁵ *kluyveri* we added a linear scaling factor κ to scale our estimates of $\Delta\eta$ between the ¹⁶ donor lineage and *L. kluyveri* and searched for the value that minimized the cost of ¹⁷ the introgression, thus giving us the best case scenario (See Methods for details). ¹⁸

¹⁹ Using our estimates of ΔM and $\Delta\eta$ from the endogenous genes and assuming the ¹⁹ current exogenous amino acid composition of genes is representative of the replaced ²⁰ endogenous genes, we estimate the strength of selection against the exogenous genes ²¹ at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of ²² selection bias for the exogenous genes show that, while well correlated with the ²³ endogenous genes, only nine amino acids share the same selectively preferred codon. ²⁴

²⁵ Exogenous genes are, therefore, expected to represent a significant reduction in ²⁶ fitness for *L. kluyveri* due to mismatch in codon usage. Since $\Delta\eta$ is proportional ²⁶ to the difference in fitness between the wild type and a mutant, we can use our ²⁷ estimates of $\Delta\eta$ to approximate the selection against the exogenous genes Δs [10, ²⁸ 26]. We estimate that the selection against all exogenous genes due to mismatched ²⁹ codon usage to have been $\Delta s \approx -0.0008$ at the time of the introgression and ³⁰ ≈ -0.0003 today. This reduction in Δs is primarily due to adaptive changes to the ³¹ codon usage of the most highly expressed, introgressed genes (Figures 5a & S8). ³²

³³ Based on the selection against the codon mismatch at the time of the introgression ³³

¹and assuming an effective population size N_e on the order of 10^7 [27], we estimate¹
²a fixation probability of $(1 - \exp[-\Delta s])/(1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$ [26] for the²
³exogenous genes. Clearly, the possibility of fixation under this simple scenario is³
⁴effectively zero. In order for the exogenous genes to have reached fixation one or⁴
⁵more exogenous loci must have provided a selective advantage not considered in⁵
⁶this study (See Discussion).⁶

⁷

⁶

⁷

⁸Discussion

⁸

⁹In order to study the evolutionary effects of the large scale introgression of the left⁹
¹⁰arm of chromosome C, we used ROC SEMPPR, a mechanistic model of ribosome¹⁰
¹¹movement along an mRNA. The usage of a mechanistic model rooted in popula-¹¹
¹²tion genetics allows us generate more nuanced quantitative parameter estimates¹²
¹³and separate the effects of mutation and selection on the evolution of codon usage.¹³
¹⁴This allowed us to calculate the selection against the introgression, and provides E .¹⁴
¹⁵*gossypii* as a potential source lineage of the introgression which was previously not¹⁵
¹⁶considered. Our parameter estimates indicate that the *L. kluyveri* genome contains¹⁶
¹⁷distinct signatures of mutation and selection bias from both an endogenous and ex-¹⁷
¹⁸ogenous cellular environment. By fitting ROC SEMPPR separately to *L. kluyveri*'s¹⁸
¹⁹endogenous and exogenous sets of genes we generate a quantitative description of¹⁹
²⁰their signatures of mutation bias and natural selection for efficient protein transla-²⁰
²¹tion.²¹

²² In contrast to other methods such as RSCU, CAI, or tAI, ROC SEMPPR does²²
²³not rely on external information such as gene expression or tRNA gene copy number²³
²⁴[5, 19]. Instead, ROC SEMPPR allows for the estimation of protein synthesis rate ϕ ²⁴
²⁵and separates the effects of mutation and selection on codon usage. In addition, [20]²⁵
²⁶showed that approaches like CAI are sensitive to amino acid composition, another²⁶
²⁷property that distinguishes the endogenous and exogenous genes [15].²⁷

²⁸ Previous work by [15] showed an increased bias towards GC rich codons in the²⁸
²⁹exogenous genes but our results provide more nuanced insights by separating the²⁹
³⁰effects of mutation bias and selection. We are able to show that the difference in GC³⁰
³¹content between endogenous and exogenous genes is mostly due to differences in³¹
³²mutation bias as 95% of exogenous codon families show a strong mutation bias to-³²
³³wards GC ending codons (Table S1). However, the exogenous genes show a selective³³

¹preference for AT ending codons for 90% of codon families (Table S2). Acknowl-¹
²edging the increased mutation bias towards GC ending codons and the difference in²
³strength of selection between endogenous and exogenous genes by separating them³
⁴also improves our estimates of protein synthesis rate ϕ by 42% relative to the full⁴
⁵genome estimate ($R^2 = 0.46, p = 0$ vs. $0.32, p = 0$, respectively).⁵

⁶ Previous studies showed that nucleotide composition can be strongly affected by⁶
⁷biased gene conversion, which, in turn would affect codon usage. Biased gene conver-⁷
⁸sion is thought to act similar to directional selection, typically favoring the fixation⁸
⁹of G/C alleles [28, 29]. Further, [30, Harrison & Charlesworth] suggested that bi-⁹
¹⁰ased gene conversion affects codon usage in *S. cerevisiae*. ROC SEMPPR, however,¹⁰
¹¹does not explicitly account for biased gene conversion. If biased gene conversion is¹¹
¹²independent of gene expression, as in the case of DNA repair, it will be absorbed¹²
¹³in our estimates of ΔM . If instead biased gene conversion forms hotspots, and¹³
¹⁴thus becomes gene specific, it will affect our estimates of protein synthesis ϕ . This¹⁴
¹⁵might be the case at recombination hotspots. Recombination, however, is very low¹⁵
¹⁶in the introgressed region (discussed below) [15, 18]. The low recombination rate¹⁶
¹⁷also indicates that the GC content had to be high before the introgression occurred.¹⁷

¹⁸ The estimates of mutation and selection bias parameters, ΔM and $\Delta \eta$, are ob-¹⁸
¹⁹tained under an equilibrium assumption. Given that the introgression is still adapt-¹⁹
²⁰ing to its new environment, this assumption is clearly violated. However, the adap-²⁰
²¹tation of the exogenous genes progresses very slowly as a quasi-static process as²¹
²²shown in this work as well as [16]. Therefore, the genome can be assumed to main-²²
²³tain an internal equilibrium at any given time. We see empirical evidence for this²³
²⁴behavior in our ability to predict gene expression and to correctly identify the low²⁴
²⁵expression genes (Figure 1b).²⁵

²⁶ Despite the violation of the equilibrium assumption, the mutation and selection²⁶
²⁷bias parameters ΔM and $\Delta \eta$ of the introgressed exogenous genes contain informa-²⁷
²⁸tion, albeit decaying, about its previous cellular environment. We selected the top²⁸
²⁹ten lineages with the highest similarity in ΔM to see if our parameters estimates²⁹
³⁰would allow us to identify a potential source lineage. The synteny relationship of³⁰
³¹these lineages with the exogenous genes was calculated as a point of comparison as³¹
³²it provides orthogonal information to our parameter estimates. Synteny with the³²
³³exogenous genes is limited to the Saccharomycetaceae clade, excluding all of the³³

¹potential source lineages identified using codon usage but *E. gossypii* (Table 2). Interestingly, this also showed that similarity in codon usage does not correlate with phylogenetic distance.³

⁴ Previous work indicated that the donor lineage of the exogenous genes has to be a, potentially unknown, Lachancea lineage [15–18]. These previous results, however, are based on species rather than gene trees, ignoring the differential adaptation rate to their novel cellular environment between genes or do not consider lineages outside of the Lachancea clade. Considering the similarity in selection bias (Figure 2b) and our calculation of selection on the exogenous genes (Figure 5b), both of which are free of any assumption about the origin of the exogenous genes, a species tree estimated from the exogenous genes will be biased towards the Lachancea clade.¹¹ Estimating individual gene trees rather than relying on a species tree provided further evidence that the exogenous genes could originate from a lineage that does not belong to the Lachancea clade. As we highlighted in this study, relatively small sets of genes with a signal of a foreign cellular environment can significantly bias the outcome of a study. The same holds true for phylogenetic inferences [31], and as we showed the signal of the original endogenous cellular environment that shaped CUB is at different stages of decay in high and low expression genes (Figure S8).¹⁸ In summary, our work does not dispute an unknown Lachancea as possible origin, but provides an alternative hypothesis based on the codon usage of the exogenous genes, phylogenetic analysis, and synteny.²¹

²² In terms of understanding the spread of the introgression, we calculated the expected 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, ~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

¹fixation probability if the current exogenous genes would introgress into *L. kluyveri*.¹

²Our estimate of the current selective cost of the mismatch of codon usage is on the²

³order of -0.0003 . The fixation probability of the current exogenous genes would³

⁴still be astronomically small $\approx 10^{-2609} \approx 0$ These results are in accordance with⁴

⁵previous work, highlighting the necessity of codon usage compatibility between en-⁵

⁶dogenous and transferred exogenous genes [32, 33]. Thus, the basic scenario of an⁶

⁷introgression between two yeast species with large N_e and where the introgression⁷

⁸solely imposes a selective cost due to codon mismatch is clearly too simplistic.⁸

⁹ One or more loci with a combined selective advantage on the order of 0.0008 ⁹

¹⁰or greater would have made the introgression change from disadvantageous to ef-¹⁰

¹¹fectively neutral or advantageous. While this scenario seems plausible, it raises¹¹

¹²the question as to why recombination events did not limit the introgression to¹²

¹³only the adaptive loci. A potential answer is the low recombination rate between¹³

¹⁴the endogenous and exogenous regions [15, 18]. Estimates of the recombination¹⁴

¹⁵rate as measured by crossovers (COs) for *L. kluyveri* are almost four times lower¹⁵

¹⁶than for *S. cerevisiae* and about half that of *Schizosaccharomyces pombe* (≈ 1.6 ¹⁶

¹⁷COs/Mb/meiosis, ≈ 6 COs/Mb/meiosis, ≈ 3 COs/Mb/meiosis) with no observed¹⁷

¹⁸crossovers in the introgressed region [18], and no observed transposable elements¹⁸

¹⁹[15]. This is presumably due to the dissimilarity in GC content and/or a lower than¹⁹

²⁰average sequence homology between the exogenous region and the one it replaced.²⁰

²¹A population bottleneck reducing the N_e of the *L. kluyveri* lineage around the time²¹

²²of the introgression could also help explain the spread of the introgression. Compati-²²

²³ble with these explanation is the possibility of several advantageous loci distributed²³

²⁴across the exogenous region drove a rapid selective sweep and/or the population²⁴

²⁵through a bottleneck speciation process.²⁵

²⁶ Assuming *E. gossypii* as potential source lineage of the exogenous region, we²⁶

²⁷illustrated how information on codon usage can be used to infer the time since²⁷

²⁸the introgression occurred using our estimates of mutation bias ΔM . The ΔM ²⁸

²⁹estimates are well suited for this task as they are free of the influence of selection²⁹

³⁰and unbiased by N_e and other scaling terms, which is in contrast to our estimates of³⁰

³¹ $\Delta\eta$ [10]. Our estimated age of the introgression of $6.2 \pm 1.2 \times 10^8$ generations is ~ 10 ³¹

³²times longer than a previous minimum estimate by [16] of 5.6×10^7 generations,³²

³³which was based on the effective population recombination rate and the population³³

¹mutation parameter [34]. Furthermore, these estimates assume that the current *E. gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the time of the introgression. Thus, if the ancestral mutation environments were more similar (dissimilar) at the time of the introgression then our result is an overestimate (underestimate).

⁶ Further, the presented work provides a template to explore the evolution of codon usage. This applies not only to species who experienced an introgression but is more generally applicable to any species.

⁹

¹⁰**Conclusion**

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

²³

²⁴**Methods**

²⁵Separating Endogenous and Exogenous Genes

²⁶A 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.

¹Model Fitting with ROC SEMPPR

²ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) [22] and R (3.4.1)²
³[35]. ROC SEMPPR was run from 10 different starting values for at least 250,000³
⁴iterations and thinned to keep every 50th iteration. After manual inspection to⁴
⁵verify that the MCMC had converged, parameter posterior means, log posterior⁵
⁶probability and log likelihood were estimated from the last 500 samples (last 10%⁶
⁷of samples).
⁷

8

⁹Model selection

¹⁰The marginal likelihood of the combined and separated model fits was calculated¹⁰
¹¹using a generalized harmonic mean estimator [36]. A variance scaling of 1.1 was¹¹
¹²used to scale the important density of the estimator. Using the estimated marginal¹²
¹³likelihoods, we calculated the Bayes factor to assess model performance. Increases¹³
¹⁴in the variance scaling increase the estimated Bayes factor, therefore we report a¹⁴
¹⁵conservative Bayes factor bases on a small variance scaling S9.
¹⁵

16

¹⁷Comparing Codon Specific Parameter Estimates and Selecting Candidate lineages

¹⁸As the choice of reference codon can reorganize codon families coding for an amino¹⁸
¹⁹acid relative to each other, all parameter estimates were interpreted relative to the¹⁹
²⁰mean for each codon family.
²⁰

21

$$\Delta M_i = \Delta M_{i,1} - \overline{\Delta M_i} \quad (1)_{22}$$

23

$$\Delta \eta_i = \Delta \eta_{i,1} - \overline{\Delta \eta_i} \quad (2)_{24}$$

²⁵Comparison of codon specific parameters (ΔM and $\Delta \eta = 2N_e q(\eta_i - \eta_j)$) was per-²⁵
²⁶formed using the function lmodel2 in the R package lmodel2 (1.7.3) [37] and R²⁶
²⁷version 3.4.1 [35]. The parameter $\Delta \eta$ can be interpreted as the difference in fitness²⁷
²⁸between codon i and j for the average gene with $\phi = 1$ scaled by the effective pop-²⁸
²⁹ulation size N_e , and the selective cost of an ATP q [4, 10]. Type II regression was²⁹
³⁰performed with re-centered parameter estimates, accounting for noise in dependent³⁰
³¹and independent variable [24].
³¹

³²Due to the greater dissimilarity of the ΔM estimates between the endogenous and³²
³³exogenous genes, and the slower decay rate of mutation bias, we decided to focus³³

¹on our estimates of mutation bias to identify potential source lineages. The top ten¹
²lineages with the highest similarity in ΔM to the exogenous genes were selected as²
³potential candidates (Figure 2).

⁴

⁵Phylogenetic Analysis

⁶Using the dataset from [21], we first identified 129 alignments for exogenous genes⁶
⁷that further contained homologous genes for *E. gossypii*, and at least one other⁷
⁸Lachancea species. We excluded all species from the alignments that do not belong⁸
⁹to the Saccharomycetaceae clade. IQTree [25] was used to identify the best fit-⁹
¹⁰ting model for each gene and to estimate the individual gene trees. Each gene tree¹⁰
¹¹was rooted using either *Saccharomyces cerevisiae*, *Saccharomyces uvarum*, *Saccha-*¹¹
¹²*romyces eubayanus* as outgroup. We calculated the most recent common ancestor¹²
¹³(MRCA) of *L. kluyveri* and *E. gossypii* as well as the MRCA of *L. kluyveri* and the¹³
¹⁴remaining Lachancea. The distance between the MRCA and the root was used to¹⁴
¹⁵asses which pairs (*L. kluyveri* and *E. gossypii*, or *L. kluyveri* and other Lachancea)¹⁵
¹⁶have a more recent common ancestor.

¹⁷

¹⁸Synteny Comparison

¹⁹We obtained complete genome sequences for all 10 candidate lineages (Table 2)¹⁹
²⁰from NCBI (on: 02-05-2017). Genomes were aligned and checked for synteny using²⁰
²¹SyMAP (4.2) with default settings [38, 39]. We assess synteny as percentage coverage²¹
²²of the exogenous gene region.

²³

²⁴Estimating Age of Introgression

²⁵We modeled the change in codon frequency over time using an exponential model²⁵
²⁶for all two codon amino acids. While our approach is equivalent to [40], we want²⁶
²⁷to explicitly state the relationship between the change in codon frequency c_1 as a²⁷
²⁸function of mutation bias ΔM as

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

²⁹
³⁰
³¹where $\mu_{i,j}$ is the rate at which codon i mutates to codon j and c_1 is the fre-³¹
³²quency of the reference codon. Initial codon frequencies $c_1(0)$ for each codon³²
³³family were taken from our mutation parameter estimates for *E. gossypii* where³³

¹ $c_1(0) = \exp[\Delta M_{\text{gos}}]/(1 + \exp[\Delta M_{\text{gos}}])$. Our estimates of ΔM_{endo} can be used to¹
²calculate the steady state of equation 3 were $\frac{dc_1}{dt} = 0$ to obtain the equality²

$$\frac{\mu_{2,1}}{\mu_{1,2} + \mu_{2,1}} = \frac{1}{1 + \exp[\Delta M_{\text{endo}}]} \quad (4)_4$$

⁵ Solving for $\mu_{1,2}$ gives us $\mu_{1,2} = \Delta M_{\text{endo}} \exp[\mu_{2,1}]$ which allows us to rewrite and⁵
⁶solve equation 3 as⁶

$$c_1(t) = \frac{1 + \exp[-X](K - 1)}{1 + \Delta M_{\text{endo}}} \quad (5)^8$$

⁹where $X = (1 + \Delta M_{\text{endo}})\mu_{2,1}t$ and $K = c_1(0)(1 + \Delta M_{\text{endo}})$.¹⁰

¹¹Equation 5 was solved with a mutation rate $\mu_{2,1}$ of 3.8×10^{-10} per nucleotide per¹¹
¹²generation [41]. Current codon frequencies for each codon family where taken from¹²
¹³our estimates of ΔM from the exogenous genes. Mathematica (11.3) [42] was used¹³
¹⁴to calculate the time t_{intro} it takes for the initial codon frequencies $c_1(0)$ for each¹⁴
¹⁵codon family to equal the current exogenous codon frequencies. The same equation¹⁵
¹⁶was used to determine the time t_{decay} at which the signal of the exogenous cellular¹⁶
¹⁷environment has decayed to within 1% of the endogenous environment.¹⁷

¹⁸Estimating Selection against Codon Mismatch

¹⁹In order to estimate the selection against codon mismatch, we had to make three¹⁹
²⁰key assumptions. First, we assumed that the current exogenous amino acid sequence²⁰
²¹of a gene is representative of its ancestral state and the replaced endogenous gene²¹
²²it replaced. Second, we assume that the currently observed cellular environment of²²
²³*E. gossypii* reflects the cellular environment that the exogenous genes experienced²³
²⁴before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy²⁴
²⁵of selection between the cellular environments due to differences in either effective²⁵
²⁶population size N_e or the selective cost of an ATP q of the source lineage and *L.*²⁶
²⁷*kluyveri* can be expressed as a scaling constant and that protein synthesis rate ϕ ²⁷
²⁸has not changed between the replaced endogenous and the introgressed exogenous²⁸
²⁹genes. Using estimates for $N_e = 1.36 \times 10^7$ [27] for *Saccharomyces paradoxus* we²⁹
³⁰scale our estimates of $\Delta\eta$ which explicitly contains the effective population size N_e ³⁰
³¹[10] and define $\Delta\eta' = \frac{\Delta\eta}{N_e}$.³¹

³²All of our genome parameter estimations are scaled by lineage specific effects such³²
³³as N_e , the average, absolute gene expression level, and/or the proportionate fitness³³

¹value of an ATP. In order to account for these genome specific differences in scaling,¹
²we scale the difference in the efficacy of selection on codon usage between the donor²
³lineage and *L. kluyveri* using a linear scaling factor κ . As $\Delta\eta$ is defined as $\Delta\eta =$ ³
⁴ $2N_e q(\eta_i - \eta_j)$, we cannot distinguish if κ is a scaling on protein synthesis rate ϕ ,⁴
⁵effective population size N_e , or the selective cost of an ATP q [4, 10]. We calculated⁵
⁶the selection against each genes codon mismatch assuming additive fitness effects⁶
⁷as
⁸

$$s_g = \sum_{i=1}^{L_g} -\kappa \phi_g \Delta\eta'_i \quad (6)^9$$

10

¹¹where s_g is the overall strength of selection for translational efficiency on gene, g ¹¹
¹²in the exogenous gene set, κ is a constant, scaling the efficacy of selection between¹²
¹³the endogenous and exogenous cellular environments, L_g is length of the protein in¹³
¹⁴codons, ϕ_g is the estimated protein synthesis rate of the gene in the endogenous¹⁴
¹⁵environment, and $\Delta\eta'_i$ is the $\Delta\eta'$ for the codon at position i . As stated previously,¹⁵
¹⁶our $\Delta\eta$ are relative to the mean of the codon family. We find that the selection¹⁶
¹⁷against the introgressed genes is minimized at $\kappa \sim 5$ (Figure S7b). Thus, we expect¹⁷
¹⁸a five fold difference in the efficacy of selection between *L. kluyveri* and *E. gossypii*,¹⁸
¹⁹due to differences in either protein synthesis rate ϕ , effective population size N_e ,¹⁹
²⁰and/or the selective cost of an ATP q . Therefore, we set $\kappa = 1$ if we calculate the s_g ²⁰
²¹for the endogenous and the current exogenous genes, and $\kappa = 5$ for s_g for selection²¹
²²calculations at the time of introgression.
²²

²³ However, since we are unable to observe codon sequences of the replaced en-²³
²⁴dogenous genes and for the exogenous genes at the time of introgression, instead²⁴
²⁵of summing over the sequence, we calculate the expected codon count $E[n_{g,i}]$ for²⁵
²⁶codon i in gene g simply as the probability of observing codon i multiplied by the²⁶
²⁷number of times the corresponding amino acids is observed in gene g , yielding:
²⁷

$$E[n_{g,i}] = P(c_i | \Delta M, \Delta\eta, \phi) \times m_{a_i} \quad 29$$

$$= \frac{\exp[-\Delta M_i - \Delta\eta_i \phi_g]}{\sum_j^C \exp[-\Delta M_j - \Delta\eta_j \phi_g]} \times m_{a_i} \quad 30$$

31

³² where m_{a_i} is the number of occurrences of amino acid a that codon i codes for. Thus³²
³³ replacing the summation over the sequence length L_g in equ. (6) by a summation³³

¹over the codon set C and calculating s_g as

²

$$s_g = \sum_{i=1}^C -\kappa \phi_g \Delta \eta'_i E[n_{g,i}] \quad (7)^3$$

⁴

⁵We report the selection due to mismatched codon usage of the introgression as

⁶ $\Delta s_g = s_{\text{intro},g} - s_{\text{endo},g}$ where $s_{\text{intro},g}$ is the selection against an introgressed gene g

⁷either at the time of the introgression or presently.

⁸**Abbreviations**

⁹**AIC:** Akaike information criterion; **CAI:** Codon adaptation index; **CUB:** Codon

¹⁰usage bias; **ROC SEMPPR:** ribosome overhead costs Stochastic Evolutionary

¹¹Model of Protein Production Rate; **RSCU:** Relative synonymous codon usage;

¹²**tAI:** tRNA adaptation index

¹³

¹⁴**Declarations**

¹⁵Ethics approval and consent to participate

¹⁶Not applicable

¹⁷

¹⁸Consent to publish

¹⁹Not applicable

²⁰

²¹Availability of data and materials

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

²³Competing interests

²⁴The authors declare that they have no competing interests.

²⁵**Funding**

²⁶This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ), MCB-1546402 (A. Von Arnim and MAG), and DEB-1355033 (BCO, MAG, and RZ) with additional support from Department of Ecology &

²⁷Evolutionary Biology (EEB) at the University of Tennessee Knoxville (UTK) and the National Institute for

²⁸Mathematical and Biological Synthesis (NIMBioS), an Institute sponsored by the National Science Foundation through NSF Award DBI-1300426. CL received support as a Graduate Student Fellow from NIMBioS with

²⁹additional support from Departments of Mathematics and EEB at UTK. The funding bodies (NSF, NIMBioS, UTK) played no role in the design of the study and collection, analysis, and interpretation of the data, and the writing of

³⁰the manuscript.

³¹

³²Authors' contributions

³³CL and MAG initiated the study. CL collected and analyzed the data and wrote the manuscript. MAG and BCO

³⁴edited the manuscript. CL, MAG, BCO, and RZ contributed to the data analysis and acquiring of funding. All

³⁵Authors approved the final manuscript.

³⁶

³⁷**Acknowledgments**

³⁸The authors would like to thank Alexander Cope for helpful criticisms and suggestions for this work.

¹**Author details**

²Department of Ecology & Evolutionary Biology, University of Tennessee, 37996, Knoxville, TN, USA. ²National Institute for Mathematical and Biological Synthesis, 37996, Knoxville, TN, USA. ³Max-Planck Institute of Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, 01307, Dresden, Germany. ⁴Department of Business Analytics and Statistics, University of Tennessee, 37996, Knoxville, TN, USA.

⁴**References**

1. Gouy, M., Gautier, C.: Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Research* **10**, 2
6 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, 10
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? **11**
12 *Biochem Soc Trans.* **21**(4), 835–841 (1993)
8. Shah, P., Gilchrist, M.A.: Explaining complex codon usage patterns with selection for translational efficiency, **13**
mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences U.S.A* **108**(25), **14**
10231–10236 (2011)
9. Wallace, E.W., Airoldi, E.M., Drummond, D.A.: Estimating selection on synonymous codon usage from noisy **15**
experimental data. *Molecular Biology and Evolution* **30**, 1438–1453 (2013)
10. Gilchrist, M.A., Chen, W.C., Shah, P., Landerer, C.L., Zaretzki, R.: Estimating gene expression and **16**
codon-specific translational efficiencies, mutation biases, and selection coefficients from genomic data alone. **17**
Genome Biology and Evolution **7**, 1559–1579 (2015)
11. Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., Danchin, A.: Evidence for horizontal gene transfer in **18**
Escherichia coli speciation. *Journal of Molecular Biology* **222**(4), 851–856 (1991)
12. Lawrence, J.G., Ochman, H.: Amelioration of bacterial genomes: Rates of change and exchange. *Journal of **19**
Molecular Biology* **44**, 383–397 (1997)
13. Marcket-Houben, M., Gabaldón, T.: Beyond the whole-genome duplication: Phylogenetic evidence for an ancient **21**
interspecies hybridization in the baker's yeast lineage. *PLoS Biology* **13**(8), 1002220 (2015)
14. Beimforde, C., Feldberg, K., Nylander, S., Rikkinen, J., Tuovila, H., Dörfelt, H., Gube, M., Jackson, D.J., **22**
Reitner, J., Seyfullah, L.J., Schmidt, A.R.: Estimating the phanerozoic history of the ascomycota lineages: **23**
combining fossil and molecular data. *Mol. Phylogenet. Evol.* **78**, 386–398 (2014)
15. Payen, C., Fischer, G., Marck, C., Proux, C., Sherman, D.J., Coppée, J.-Y., Johnston, M., Dujon, B., **24**
Neuvéglise, C.: Unusual composition of a yeast chromosome arm is associated with its delayed replication. **25**
Genome Research **19**(10), 1710–1721 (2009)
16. Friedrich, A., Reiser, C., Fischer, G., Schacherer, J.: Population genomics reveals chromosome-scale **26**
heterogeneous evolution in a protoploid yeast. *Molecular Biology and Evolution* **32**(1), 184–192 (2015)
17. Vakirlis, N., Sarilar, V., Drillon, G., Fleiss, A., Agier, N., Meyniel, J.-P., Blanpain, L., Carbone, A., Devillers, H., **27**
Dubois, K., Gillet-Markowska, A., Graziani, S., Huu-Vang, N., Poirel, M., Reisser, C., Schott, J., Schacherer, **28**
J., Lafontaine, I., Llorente, B., Neuvéglise, C., Fischer, G.: Reconstruction of ancestral chromosome **29**
architecture and gene repertoire reveals principles of genome evolution in a model yeast genus. *Genome **30**
research* **26**(7), 918–932 (2016)
18. Brion, C., Legrand, S., Peter, J., Caradec, C., Pflieger, D., Hou, J., Friedrich, A., Llorente, B., Schacherer, J., **31**
Variation of the meiotic recombination landscape and properties over a broad evolutionary distance in yeasts. **32**
PLoS Genetics **13**(8), 1006917 (2017)
19. dos Reis, M., Savva, R., Wernisch, L.: Solving the riddle of codon usage preferences: a test for translational **33**
selection. *Nucleic Acids Research* **32**(17), 5036–5044 (2004)
20. Cope, A.L., Hettich, R.L., Gilchrist, M.A.: Quantifying codon usage in signal peptides: Gene expression and **33**

- 1 amino acid usage explain apparent selection for inefficient codons. *Biochimica et Biophysica Acta (BBA) - Biomembranes* **1860**(12), 2479–2485 (2018) 1
- 2 21. Shen, X.X., Opulente, D.A., Kominek, J., Zhou, X., Steenwyk, J.L., Buh, K.V., Haase, M.A.B., Wisecaver, 2
J.H., Wang, M., Doering, D.T., Boudouris, J.T., Schneider, R.M., Langdon, Q.K., Ohkuma, M., Endoh, R., 3
Takashima, M., Manabe, R., Čadež, N., Libkind, D., Rosa, C., DeVirgilio, J., Hulfachor, A.B., Groenewald, M., 4
Kurtzman, C., Hittinger, C.T., Rokas, A.: Tempo and mode of genome evolution in the budding yeast 4
subphylum. *Cell* **175**(6), 1533–154520 (2018) 5
- 6 22. Landerer, C., Cope, A., Zaretzki, R., Gilchrist, M.A.: AnaCoDa: analyzing codon data with bayesian mixture 6
models. *Bioinformatics* **34**(14), 2496–2498 (2018)
- 7 23. Tsankov, A.M., Thompson, D.A., Socha, A., Regev, A., Rando, O.J.: The role of nucleosome positioning in the 7
evolution of gene regulation. *PLoS Biol* **8**(7), 1000414 (2010)
- 8 24. Sokal, R.R., Rohlf, F.J.: *Biometry - The principles and practice of statistics in biological*, pp. 547–555. W. H. 8
Freeman, New York, NY (1981) 9
- 9 25. Nguyen, L.T., Schmidt, H.A., von Haeseler, A., Minh, B.Q.: Iq-tree: A fast and effective stochastic algorithm 9
for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **32**(1), 268–274 (2015) 10
- 10 26. Sella, G., Hirsh, A.E.: The application of statistical physics to evolutionary biology. *Proceedings of the National 11
Academy of Sciences of the United States of America* **102**, 9541–9546 (2005) 11
- 11 27. Wagner, A.: Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* **22**, 12
1365–1374 (2005)
- 12 28. Nagylaki, T.: Evolution of a finite population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**, 13
6278–6281 (1983) 14
- 13 29. Nagylaki, T.: Evolution of a large population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**, 15
5941–5945 (1983) 15
- 14 30. Harrison, R.J., Charlesworth, B.: Biased gene conversion affects patterns of codon usage and amino acid usage 16
in the *Saccharomyces* sensu stricto group of yeasts. *Molecular Biology and Evolution* **28**(1), 117–129 (2011) 16
- 15 31. Salichos, L., Rokas, A.: Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature* 17
497, 327–331 (2013) 17
- 16 32. Medrano-Soto, A., Moreno-Hagelsieb, G., Vinuesa, P., Christen, J.A., Collado-Vides, J.: Successful lateral 18
transfer requires codon usage compatibility between foreign genes and recipient genomes. *Molecular Biology 19
and Evolution* **21**(10), 1884–1894 (2004) 19
- 20 33. Tuller, T., Girshovich, Y., Sella, Y., Kreimer, A., Freilich, S., Kupiec, M., Gophna, U., Ruppin, E.: Association 20
between translation efficiency and horizontal gene transfer within microbial communities. *Nucleic Acids 21
Research* **39**(11), 4743–4755 (2011). doi:10.1093/nar/gkr054 21
- 21 34. Ruderfer, D.M., Pratt, S.C., Seidl, H.S., Kruglyak, L.: Population genomic analysis of outcrossing and 22
recombination in yeast. *Nature Genetics* **38**(9), 1077–1081 (2006) 22
- 22 35. R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical 23
Computing, Vienna, Austria (2013). R Foundation for Statistical Computing. <http://www.R-project.org/> 24
- 23 36. Gronau, Q.F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D.S., Forster, J.J., 25
Wagenmakers, E.J., Steingrover, H.: A tutorial on bridge sampling. *Journal of Mathematical Psychology* **81**, 25
80–97 (2017) 26
- 24 37. Legendre, P.: Lmodel2: Model II Regression. (2018). R package version 1.7-3. 27
- 25 38. Soderlund, C., Nelson, W., Shoemaker, A., Paterson, A.: Symap A system for discovering and viewing syntenic 28
regions of fpc maps. *Genome Research* **16**, 1159–1168 (2006) 28
- 26 39. Soderlund, C., Bomhoff, M., Nelson, W.: Symap v3.4: a turnkey synteny system with application to plant 29
genomes. *Nucleic Acids Research* **39**(10), 68 (2011) 30
- 27 40. Marais, G., Charlesworth, B., Wright, S.I.: Recombination and base composition: the case of the highly 31
self-fertilizing plant *Arabidopsis thaliana*. *Genome Biology* **5**, 45 (2004) 31
- 28 41. Lang, G.I., Murray, A.W.: Estimating the per-base-pair mutation rate in the yeast *Saccharomyces cerevisiae*. 32
Genetics **178**(1), 67–82 (2008) 32
- 29 42. Wolfram Research Inc.: Mathematica 11. (2017). <http://www.wolfram.com> 33

- ¹43. Zhang, Y., Ponty, Y., Blanchette, M., Lécuyer, E., Waldspühl, J.: Sparcs: a web server to analyze
2 (un)structured regions in coding rna sequences. Nucleic Acids Research **41**, 480–485 (2013)
- 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**Supplementary Material**

²Supporting 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

⁴most likely codon in low expression genes for each amino acid in: *E. gossypii*, in the endogenous and

⁵exogenous genes of *L. kluyveri*, and in the combined *L. kluyveri* genome without accounting for the

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

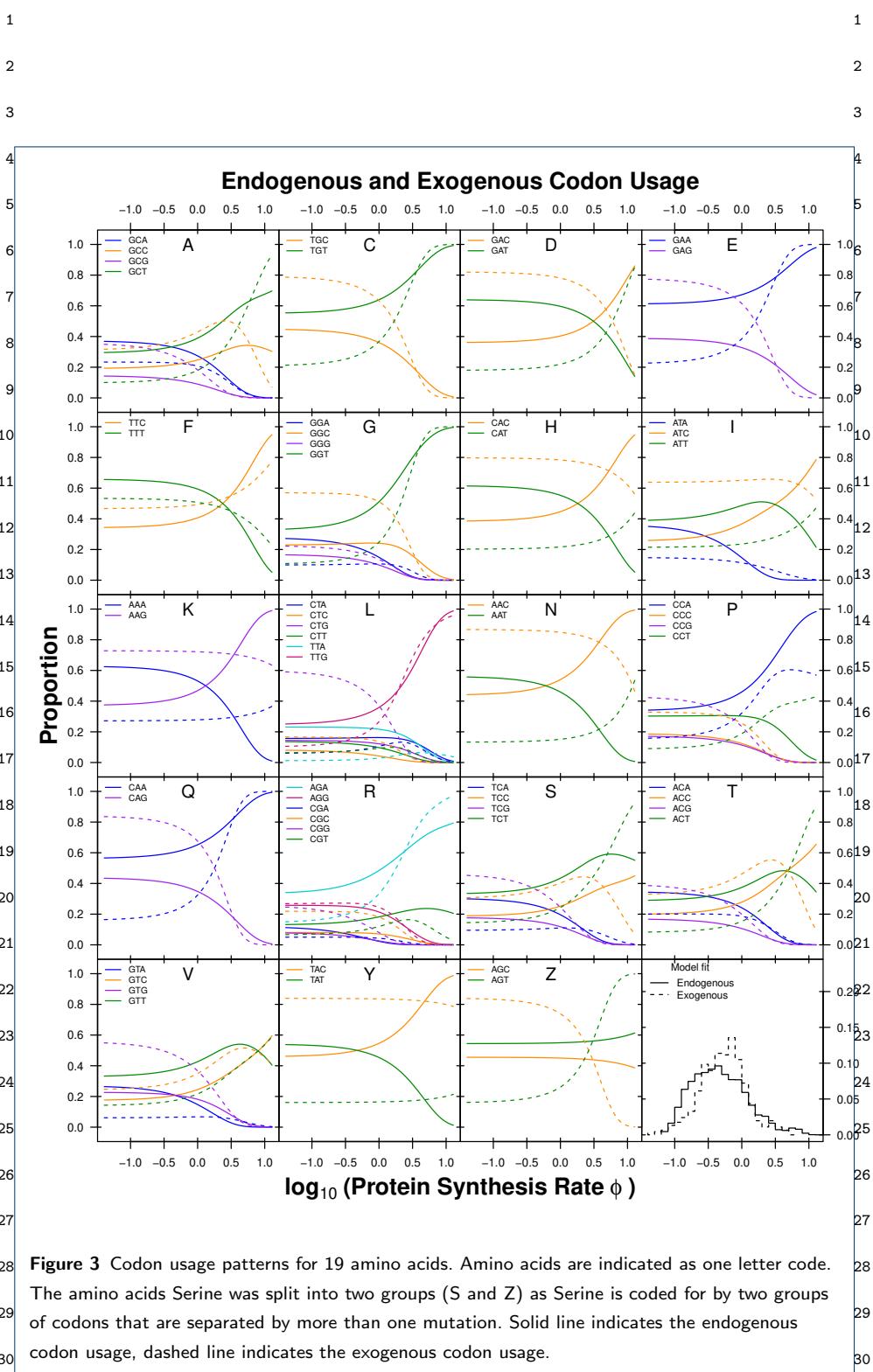


Figure 3 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dashed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

31

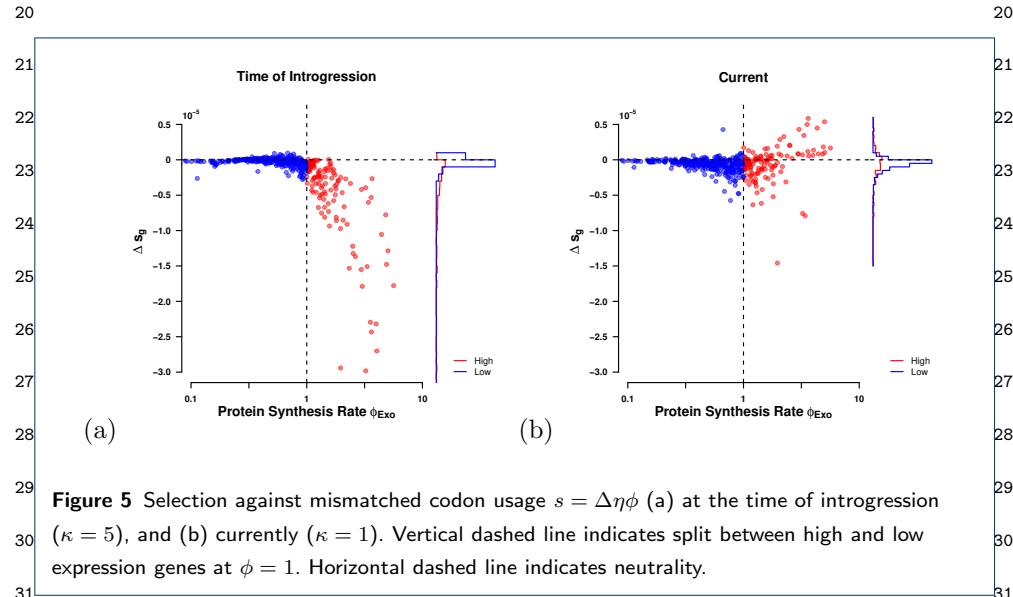
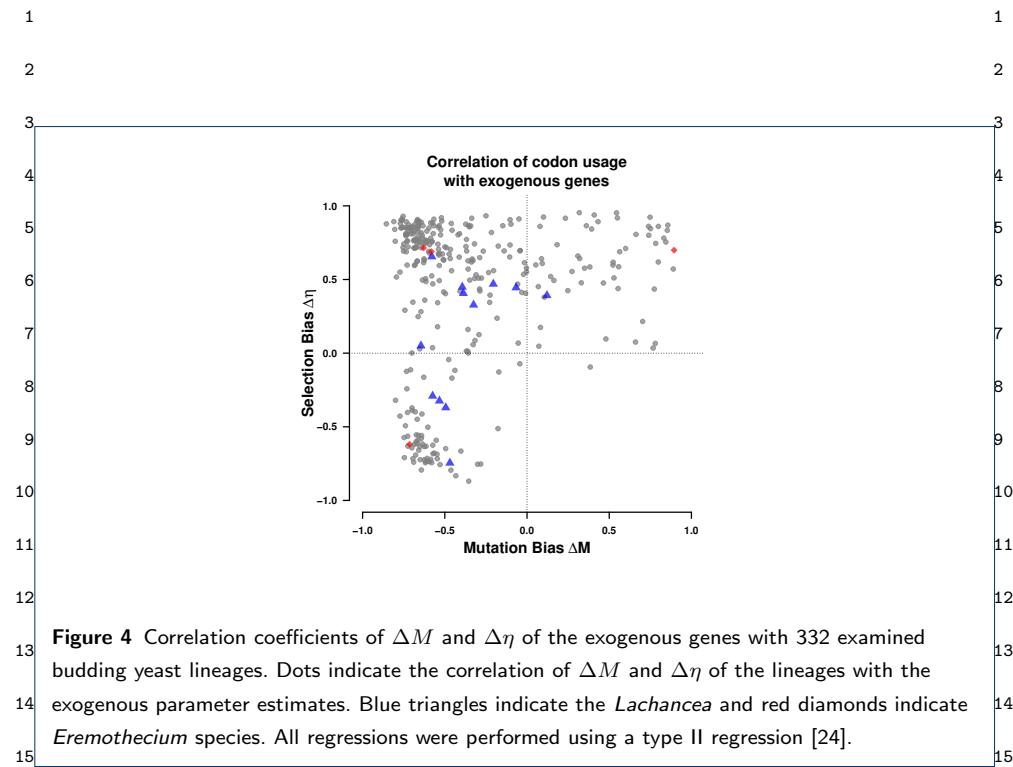
31

32

32

33

33



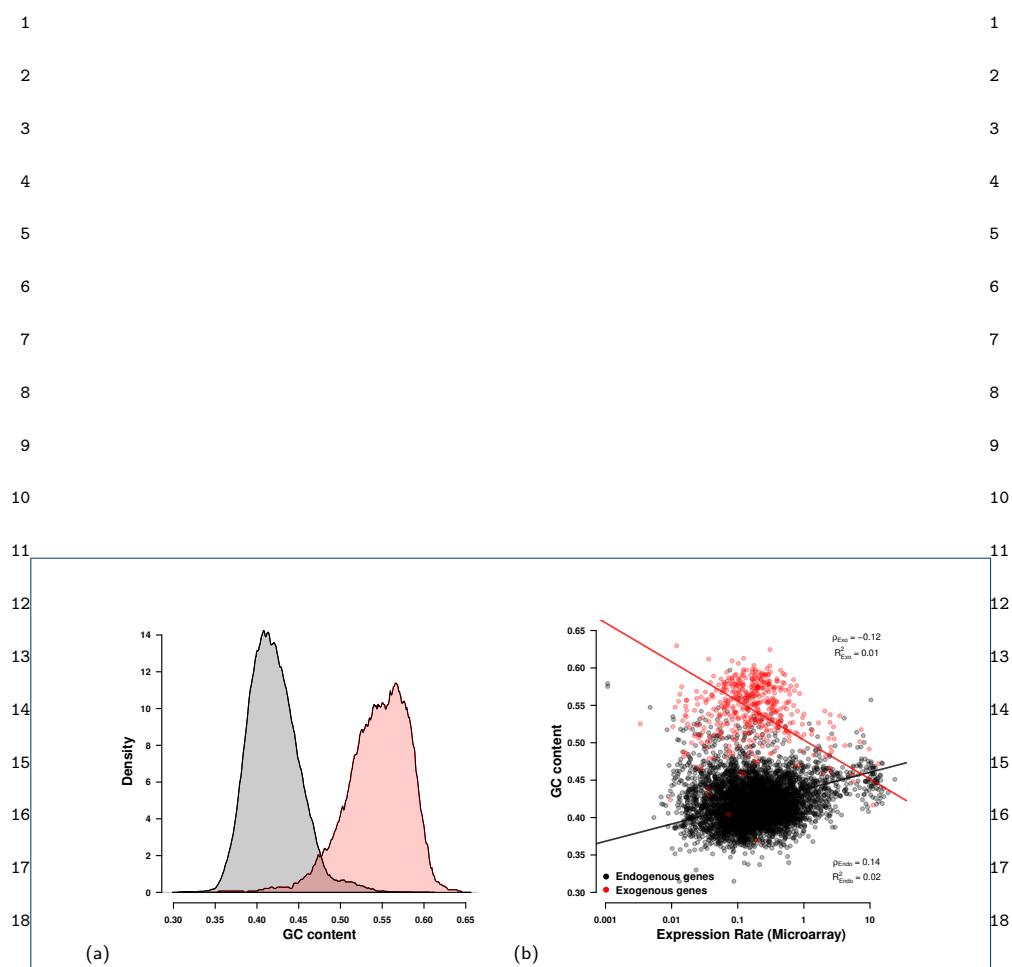


Figure S1 Endogenous and exogenous genes have distinct GC content. (a) Distribution of GC content content in the endogenous and exogenous genes. (b) Correlation of endogenous and exogenous GC content with measured gene expression. While the endogenous GC content shows a slight positive correlation with gene expression ($\rho = 0.14, p = 1.2 \times 10^{-21}$), the exogenous GC content is negatively correlated with gene expression ($\rho = -0.12, p = 0.014$).

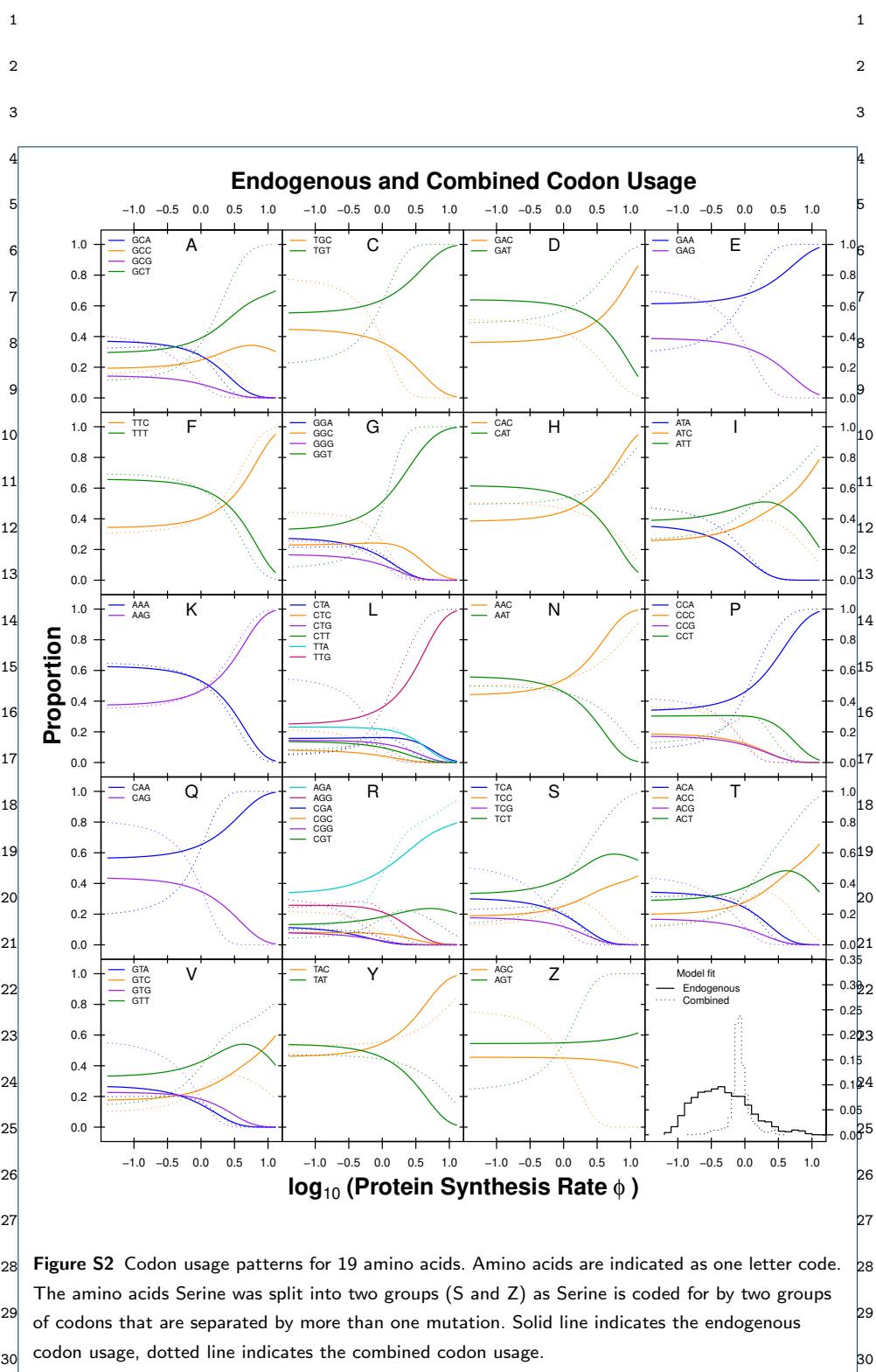


Figure S2 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dotted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

31

31

32

32

33

33

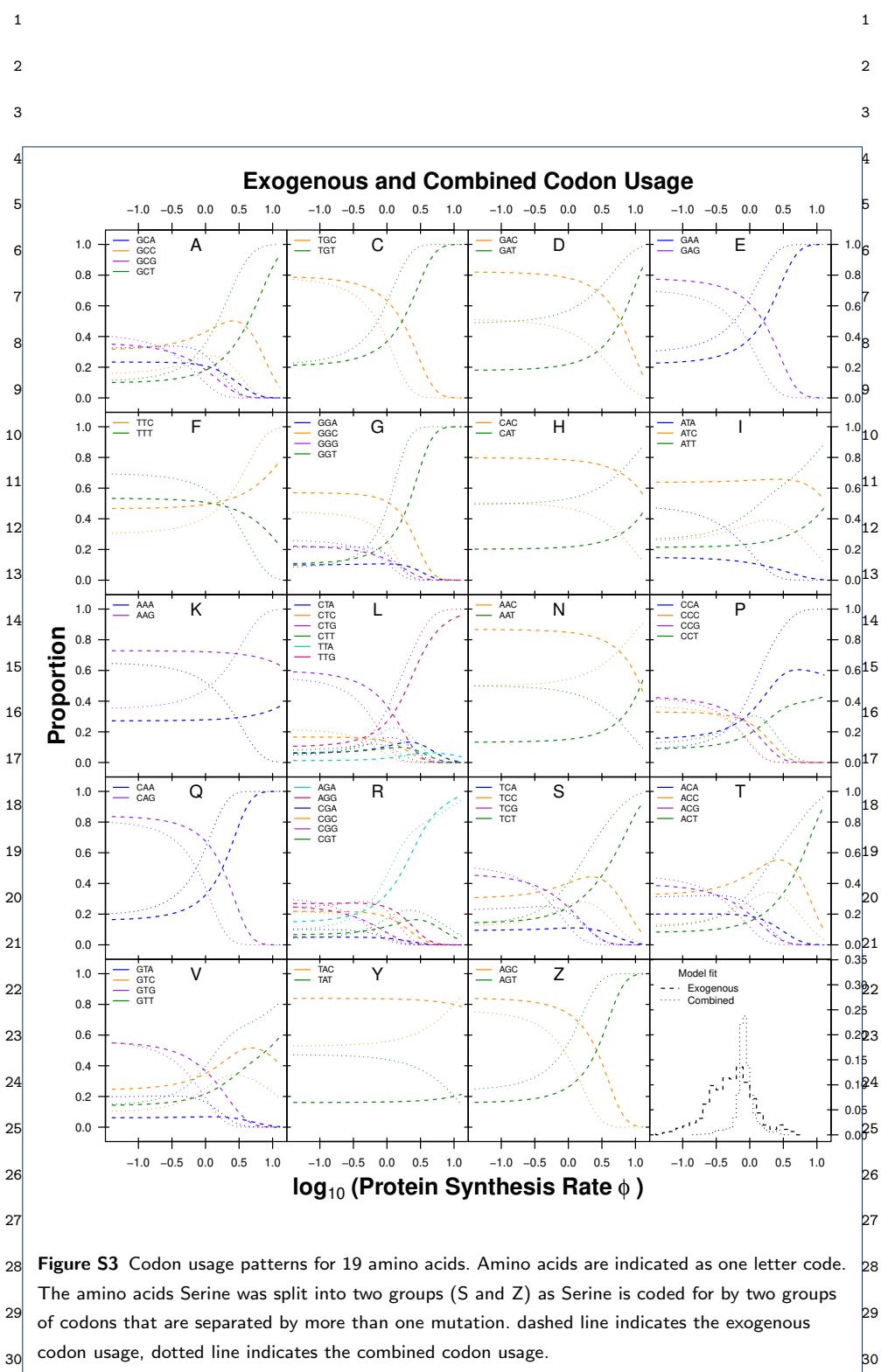


Figure S3 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. dashed line indicates the exogenous codon usage, dotted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

31

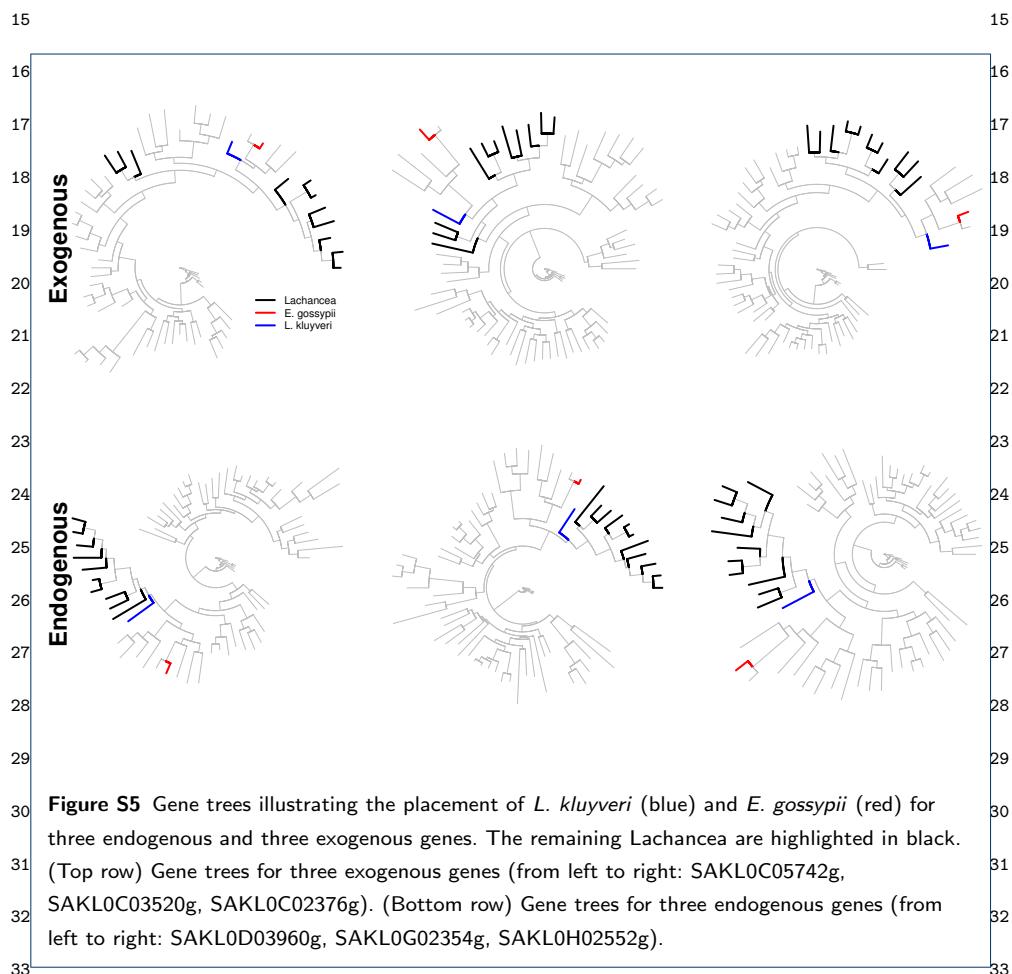
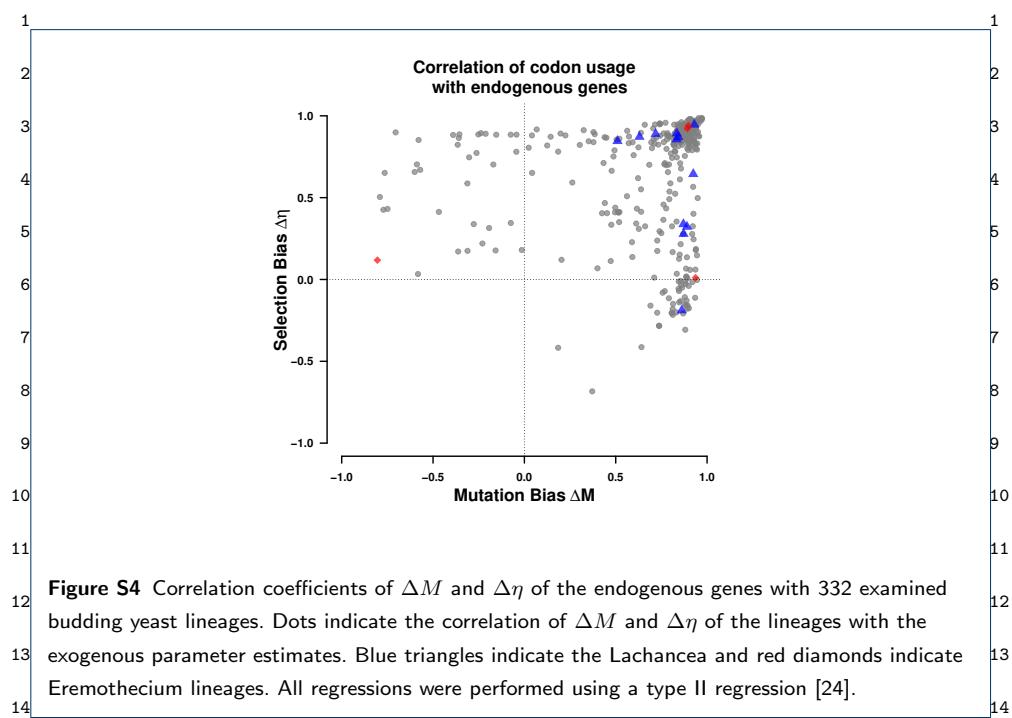
31

32

32

33

33



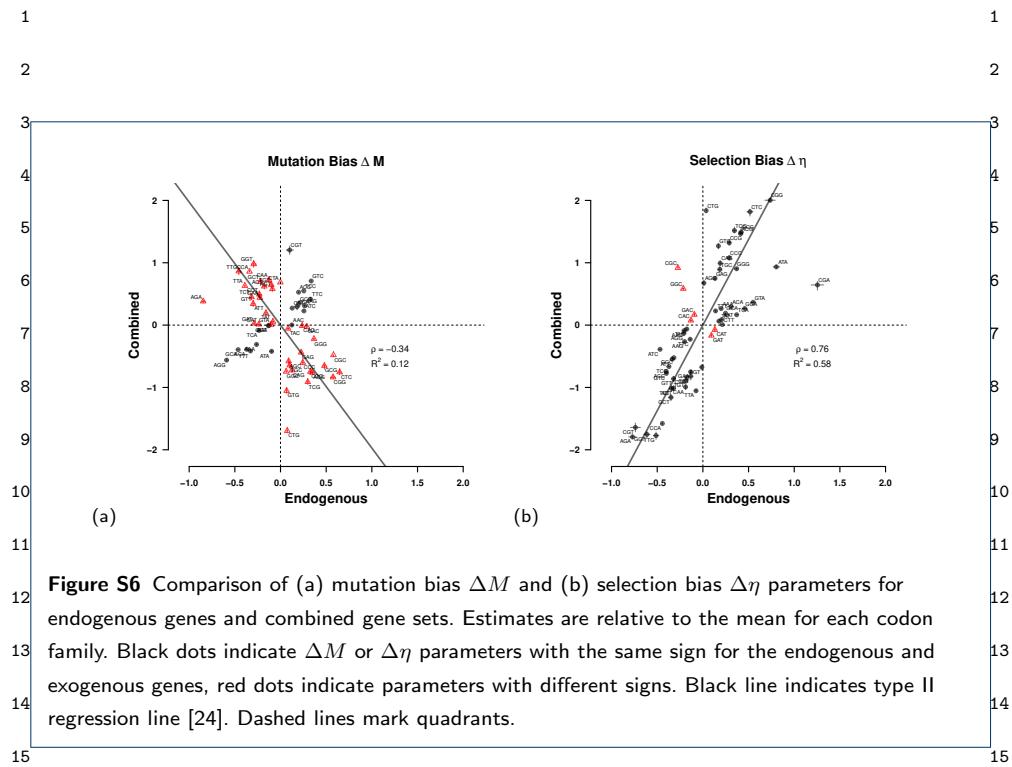


Figure S6 Comparison of (a) mutation bias ΔM and (b) selection bias $\Delta \eta$ parameters for endogenous genes and combined gene sets. Estimates are relative to the mean for each codon family. Black dots indicate ΔM or $\Delta \eta$ parameters with the same sign for the endogenous and exogenous genes,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

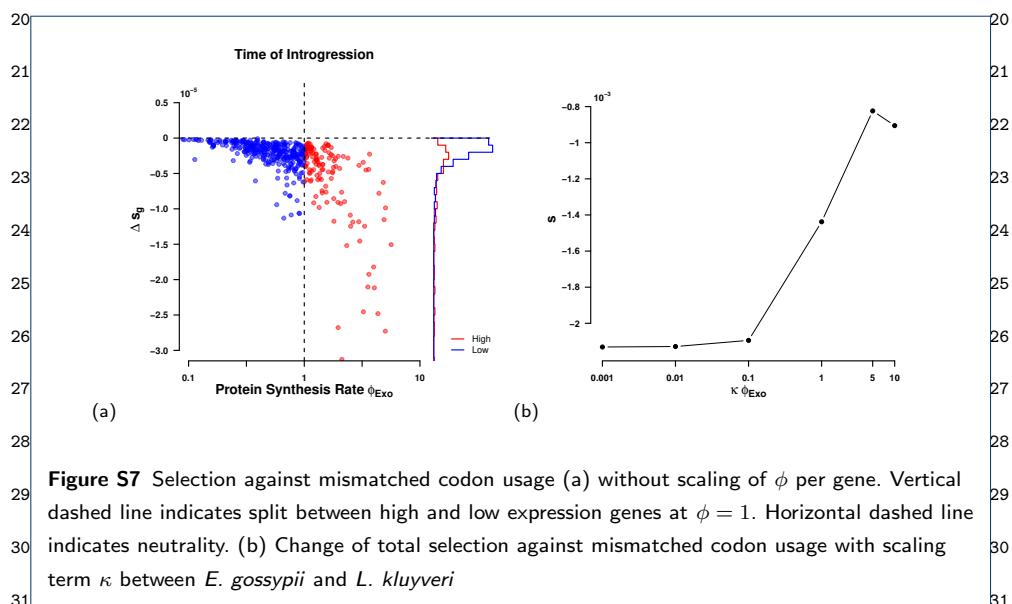


Figure S7 Selection against mismatched codon usage (a) without scaling of ϕ per gene. Vertical dashed line indicates split between high and low expression genes at $\phi = 1$. Horizontal dashed line indicates neutrality. (b) Change of total selection against mismatched codon usage with scaling term κ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

