

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
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}

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

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

Abstract

Background: For decades, codon usage has been used as a measure of adaptation for translational efficiency and translation accuracy of a gene's coding sequence. These patterns of codon usage reflect both the selective and mutational environment in which the coding sequences evolved. Over this same period, gene transfer between lineages has become widely recognized as an important biological phenomenon. Nevertheless, most studies of codon usage implicitly assume that all genes within a genome evolved under the same selective and mutational environment, an assumption violated when introgression occurs. In order to better understand the effects of introgression on codon usage patterns and vice versa, we examine the patterns of codon usage in *Lachancea kluyveri*, a yeast which has experienced a large introgression. We quantify the effects of mutation bias and selection for translation efficiency on the codon usage pattern of the endogenous and introgressed exogenous genes using a Bayesian mixture model, ROC SEMPPR, which is built on mechanistic assumptions about protein synthesis and grounded in population genetics.

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

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

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

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

1 **Background**

2
3 Synonymous codon usage patterns varies within a genome and between taxa, re-
4 reflecting differences in mutation bias, selection, and genetic drift. The signature of
5 mutation bias is largely determined by both the organism's internal cellular envi-
6 ronment, such as their DNA repair genes, and its external environment, such as UV
7 exposure. While mutation bias is an omnipresent evolutionary force, its impact on
8 codon usage bias (CUB) can be countered or reinforced by natural selection. In con-
9 trast to mutation bias, the signature of selection on CUB is largely determined by
10 the organism's cellular environment, such as its tRNA species, their copy number,
11 expression level, and post-transcriptional modifications. In general, the strength of
12 selection on the CUB of an individual gene increases with its expression level [1–
13 3], specifically its protein synthesis rate [4]. Thus as its protein synthesis increases,
14 CUB of a gene shifts from a process dominated by mutation to a process dominated
15 by selection. The balance of mutation and selection on codon usage is also affected
16 by genetic drift which, in turn, is an inverse function of the organism's effective
17 population size N_e . Disentangling these evolution forces is of keen interest to biolo-
18 gists; ROC SEMPPR allows us disentangle the contribution of mutation, selection,
19 and drift to the patterns of CUB encoded across an species' genome [4–7]. In turn,
20 these evolutionary parameters should provide biologically meaningful information
21 about the lineage's historical cellular and external environment.

22
23 Most studies implicitly assume that the CUB of a genome is shaped by a single
24 cellular and external environment. However, this assumption is clearly violated to
25 increasing degrees via horizontally gene transfer, large scale introgressions, and hy-
26 brid specie formation. In these scenarios, one would expect to see the signature of
27 multiple cellular environments in a genome's CUB [8, 9]. Indeed, differences in CUB
28 between linages have been proposed to have a major effect on their rates of gene
29 transfer with rates declining with differences in their CUB. On a more practical
30 level, if differences in codon usage of transferred genes are not taken into account
31 for, they may distort the interpretation of codon usage patterns. Such distortion
32 could lead to the wrong inference of codon preference for an amino acid [5, 7], un-
33 derestimate the variation in protein synthesis rate, or distort estimates of mutation
34 bias when analyzing a genome.

To illustrate these ideas, we analyze the CUB of the genome of the yeast *Lachancea kluyveri* using ROC SEMPPR, a population genetics based model of synonymous codon usage evolution that accounts for and, in turn, can estimate the contribution of mutation bias ΔM , selection bias. The mathematics of ROC SEMPPR are derived on a mechanistic description of ribosome movement along an mRNA, although the approximation of other biological mechanisms could also be consistent with the model. Broadly speaking, ROC SEMPPR allows us to quantify the cellular environment in which genes have evolved by separately estimating the effects of mutation bias and selection bias on codon usageDE between synonymous codons and protein synthesis rate ϕ to the patterns of codon usage observed within a set of genes. Briefly, the set of ΔM for an amino acid quantifies the relative differences in mutational stability or bias between the synonymous codons of the amino acid \mathbb{S} . In the absence of selection bias (or equivalently when gene expression $\phi = 0$), the equilibrium frequency of synonymous codon i is simply $\exp[-\Delta M_i] / \left(\sum_{j \in \mathbb{S}} \exp[-\Delta M_j] \right)$. Because the time units of protein production rate have no intrinsic time scale, we define the average protein production rate for a set of genes to be one, i.e. $\bar{\phi} = 1$ by definition [7]. In order to facilitate comparisons between gene sets, we express both, ΔM and $\Delta \eta$, as deviation from the mean of each synonymous codon family (see Materials and Methods for details). Nevertheless, the difference $\Delta \eta$ describes the difference in fitness between two synonymous codons relative to drift for a gene whose protein production rate ϕ is equal to the the average rate of protein production $\bar{\phi}$ across the set of genes. In other words, for a gene whose protein is expressed at the average rate, for any two given synonymous codons i and j , $\Delta \eta_i - \Delta \eta_j = N_e s$.

The Lachancea clade diverged from the Saccharomyces clade, prior to its whole genome duplication ~ 100 Mya ago [10, 11]. Since that time, *L. kluyveri*, which is sister species to all other *Lachancea spp.*, has experienced a large introgression of exogenous genes (1 Mb, 457 genes) which is found in all of its populations [12, 13], but in no other known Lachancea species [14]. The introgression replaced the left arm of the C chromosome and displays a 13% higher GC content than the endogenous *L. kluyveri* genome [12, 13]. Previous studies suggest that the source of the introgression is probably a currently unknown or potentially extinct Lachancea lineage based on gene concatenation or synteny relationships [12–15]. These char-

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

Lachancea

1 characteristics make *L. kluyveri* an ideal model to study the effects of an introgressed
2 cellular environment and the resulting mismatch in codon usage.

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

30 Results

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

32 We compared model fits of ROC SEMPPR to the entire *L. kluyveri* genome and
33 its genome partitioned into two sets of 4,864 endogenous and 497 exogenous genes

1

2

3 mikeg: cite these 'previous
4 studies'.

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 **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

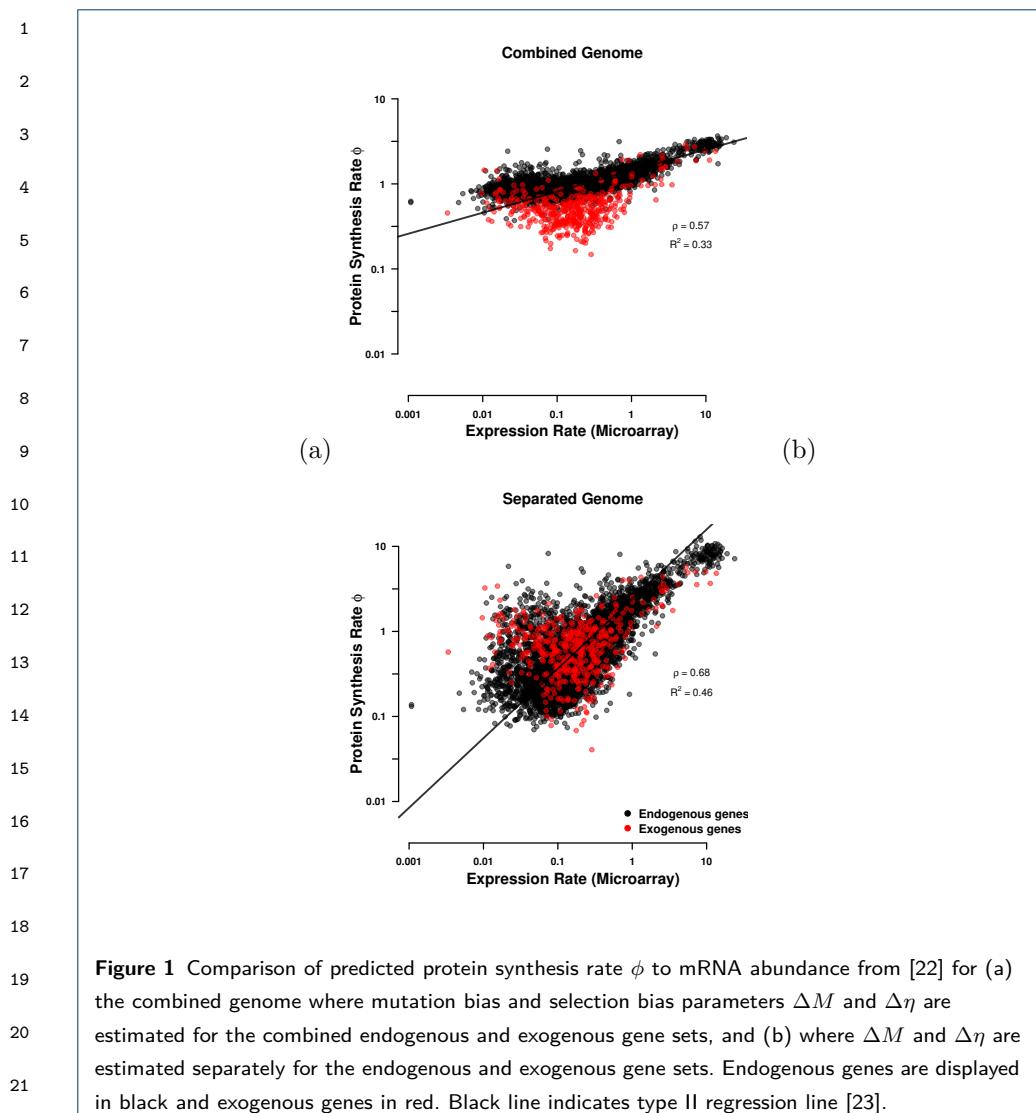
mikeg: Note that Alex has
 implemented DIC routines into
 CoDa which is better than BF
 because it is insensitive to the
 width of flat priors, which I
 naively thought didn't matter
 previously.

using the AnaCoDa software package [21]. These two gene sets where previously defined based on their striking difference in GC content [12], with very little overlap in GC content between the two sets (Figure S1a). Bayes factor strongly support the hypothesis that the *L. kluyveri* genome consists of genes with two different and distinct patterns of codon usage bias rather than a single (Bayes Factor $K = \exp[42,294]$; Table 1). We find additional support for this hypothesis when we compare our predictions of protein synthesis rate to empirically observed mRNA expression values as a proxy for protein synthesis. Specifically, we improve the variance explained by our predicted protein synthesis rates by $\sim 42\%$, from $R^2 = 0.33$ ($p \approx 0$) to 0.46 ($p \approx 0$) (Figure 1). While the implicit consideration of GC content in this analysis certainly plays a roll, it does not explain the improvement in R^2 (Figure S1b).

Comparing Differences in the Endogenous and Exogenous Codon Usage

Because ROC SEMPPR defines $\bar{\phi} = 1$, it makes the interpretation of $\Delta\eta$ as selection on codon usage of the average gene with $\phi = 1$ straightforward and gives us the ability to compare the efficacy of selection sN_e across genomes. While it may be expected for the endogenous and exogenous genes to differ in their codon usage pattern due to the large difference in GC content it is not clear how much of this difference is due to differences in the mutation bias ΔM or selection bias $\Delta\eta$ between the gene sets. To better understand the differences in the endogenous and exogenous cellular environments, we compared our parameter estimates of ΔM and $\Delta\eta$ for the two sets of genes. Our estimates of ΔM for the endogenous and exogenous genes were negatively correlated ($\rho = -0.49$, $p = 3.56 \times 10^{-5}$), indicating weak similarity with only $\sim 5\%$ of the codons share the same sign between the two mutation environments (Figure 2a). Overall, mutation bias favors codons ending in the purines A and T over codons ending in the pyrimidines G or C, respectively, as indicated by where the endogenous model fit curves intercept the left axis in Figure

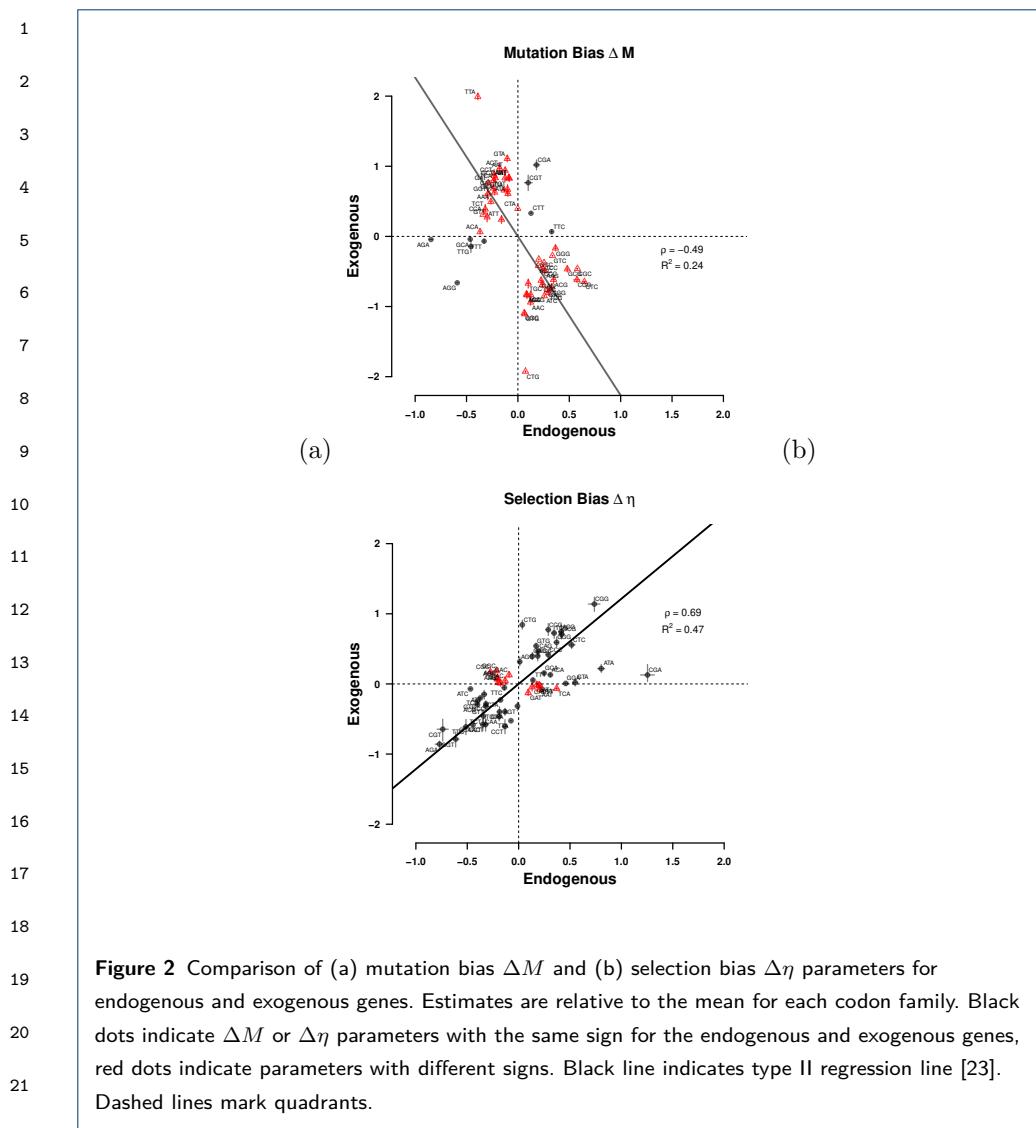
Replace $p \approx 0$ with
 $p < 10^{-15}$ or whatever.



3. One exception is Lys, L, where mutation appears to favor the codon TTG over TTA; however, this difference is very small and not statistically meaningful. The exogenous genes display behavior, favoring G over A and C over G and doing so more strongly, as indicated by shift in order and the greater distance between the exogenous model fit curves where they intercept the left axis in Figure 3).

We find that the signature of selection bias $\Delta \eta$ also differs substantially between the endogenous and exogenous gene sets. In terms of their magnitude relative to mutation bias ΔM , $\Delta \eta$ for the endogenous gene set is substantially greater than for the exogenous gene set as indicated by the fact that the protein production rate ϕ where the effect of selection fails to override the effect of The difference in codon usage between endogenous and exogenous genes is striking as some amino acids have

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
mikeg: Type II is also
referred to as 'Deming
regression' or 'errors-in
variables model'. Verify that the 95% PI of
 ΔM for TTG and TTA
mikeg P revised the ΔM
results because they
weren't actually consistent
with the curves illustrated
in Figure 3. Please check
and, if changes are correct,
please revise $\Delta \eta$ results to
mirror structure of ΔM
section.



nikeg: can you revise this
text to make it more
consistent with previous
discussions of these terms?
For example, 'opposite
codon preferences' is
confusing, is it that $\Delta\eta$
values have opposite signs
between gene sets? 25
26
27
28
29
30
31
32
33

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 that the relatively small exogenous gene set ($\sim 10\%$ of genes) has a disproportionate effect on the model fit (Figure S2, S3).

1 Of the nine cases in which the endogenous and exogenous genes show differences
2 in the selectively most favored codon five cases (Asp, D; His, H; Lys, K; Asn, N;
3 and Pro, P) the endogenous genes favor the codon with the most abundant tRNA.
4 For the remaining four cases (Ile, I; Ser, S; Thr, T; and Val, V), there are no
5 tRNA genes for the wobble free cognate codon encoded in the *L. kluyveri* genome.
6 However, the codon preference of these four amino acids in the exogenous genes
7 matches the most abundant tRNA encoded in the *L. kluyveri* genome. In contrast
8 to ΔM , our estimates of selection bias $\Delta\eta$ for the endogenous and exogenous genes
9 are positively correlated ($\rho = 0.69$, $p = 9.76 \times 10^{-10}$) and show the same sign in
10 $\sim 53\%$ of the cases (Figure 2).

11

12 This striking difference in codon usage was noted previously. For example, using
13 RSCU [16], GAA (coding for Glu, E) was identified as the optimal synonymous
14 codon in the whole genome and GAG as the optimal codon in the exogenous genes
15 [12]. Our results, however, indicate that GAA is the optimal codon in both, endoge-
16 nous and exogenous genes, and that the high RSCU in the exogenous genes of GAG
17 is driven by mutation bias (Table S1 and S2). Similar effects are observed for other
18 amino acids.

19

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

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

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

3 Since the origin of the exogenous genes is currently unknown, we explored if the 2
4 information on codon usage extracted from the exogenous genes can be used to 3
5 identify a potential source lineage. We combined our estimates of mutation bias 4
6 ΔM and selection bias $\Delta\eta$ with synteny information and searched for potential 5
7 source lineages of the introgressed exogenous region. We used ΔM to identify 6
8 candidate lineages as the endogenous and exogenous genes show greater dissimilarity 7
9 in mutation bias than in selection bias. We examined 332 budding yeasts [20] and, 8
10 identified the ten lineages with the highest correlation to the exogenous ΔM 9
11 parameters as potential source lineages (Figure 4, Table 2). Two of the ten candidate 10
12 lineages utilize the alternative yeast nuclear code (NCBI codon table 12). In this 11
13 case, the codon CTG codes for Serine instead of Leucine. We therefore excluded the 12
14 Leucine codon family from our comparison of codon families; however, there was no 13
15 need to exclude Serine as CTG is not a one step neighbor of the remaining Serine 14
16 codons. A mutation between CTG and the remaining Serine codons would require 15
17 two mutations with one of them being non-synonymous, which would violate the 16
18 weak mutation assumption of ROC SEMPPR. 17

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

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

Table 2 Budding yeast lineages showing similarity in codon usage with the exogenous genes. $\rho_{\Delta M}$ and $\rho_{\Delta \eta}$ represent the Pearson correlation coefficient for exogenous 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

extinct lineage [12–14]. While it is not possible for us to dispute this hypothesis, our results provide a novel hypothesis about the origin of the exogenous genes.

To further test the plausibility of *E. gossypii* as potential source lineage, we identified 127 genes in our dataset [20] with homologous genes in *E. gossypii* and other Lachancea and used IQTree [24] to infer the phylogenetic relationship of the exogenous genes. Our results show that at least ~ 45% of exogenous genes (57/127) are more closely related to *E. gossypii* than to other Lachancea S5. Interestingly, our results also indicate that codon usage does not necessarily correlate with phylogenetic distance (Table 2).

Estimating Introgression Age

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

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

8 Estimating Selection against Codon Mismatch of the Exogenous Genes

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

24 Using our estimates of ΔM and $\Delta \eta$ from the endogenous genes and assuming the
25 current exogenous amino acid composition of genes is representative of the replaced
26 endogenous genes, we estimate the strength of selection against the exogenous genes
27 at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of
28 selection bias for the exogenous genes show that, while well correlated with the
29 endogenous genes, only nine amino acids share the same selectively preferred codon.
30 Exogenous genes are, therefore, expected to represent a significant reduction in
31 fitness for *L. kluyveri* due to mismatch in codon usage. Since $\Delta \eta$ is proportional
32 to the difference in fitness between the wild type and a mutant, we can use our
33 estimates of $\Delta \eta$ to approximate the selection against the exogenous genes Δs [7, 25].

1 We estimate that the selection against all exogenous genes due to mismatched codon
2 usage to have been $\Delta s \approx -0.0008$ at the time of the introgression and ≈ -0.0003
3 today. This reduction in Δs is primarily due to adaptive changes to the codon
4 usage of the most highly expressed, introgressed genes (Figures 5a & S8). Based
5 on the selection against the codon mismatch at the time of the introgression and
6 assuming an effective population size N_e on the order of 10^7 [26], we estimate a
7 fixation probability of $(1 - \exp[-\Delta s])/(1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$ [25] for the
8 exogenous genes. Clearly, the possibility of fixation under this simple scenario is
9 effectively zero. In order for the exogenous genes to have reached fixation one or
10 more exogenous loci must have provided a selective advantage not considered in
11 this study (See Discussion).

14 Discussion

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

28 In contrast to other methods such as RSCU, CAI, or tAI, ROC SEMPPR does
29 not rely on external information such as gene expression or tRNA gene copy number
30 [16, 18]. Instead, ROC SEMPPR allows for the estimation of protein synthesis rate ϕ
31 and separates the effects of mutation and selection on codon usage. In addition, [19]
32 showed that approaches like CAI are sensitive to amino acid composition, another
33 property that distinguishes the endogenous and exogenous genes [12].

1 Previous work by [12] showed an increased bias towards GC rich codons in the
2 exogenous genes but our results provide more nuanced insights by separating the
3 effects of mutation bias and selection. We are able to show that the difference in GC
4 content between endogenous and exogenous genes is mostly due to differences in
5 mutation bias as 95% of exogenous codon families show a strong mutation bias to-
6 wards GC ending codons (Table S1). However, the exogenous genes show a selective
7 preference for AT ending codons for 90% of codon families (Table S2). Acknowledg-
8 ing the increased mutation bias towards GC ending codons and the difference in
9 strength of selection between endogenous and exogenous genes by separating them
10 also improves our estimates of protein synthesis rate ϕ by 42% relative to the full
11 genome estimate ($R^2 = 0.46, p = 0$ vs. $0.32, p = 0$, respectively).

12 Previous studies showed that nucleotide composition can be strongly affected by
13 biased gene conversion, which, in turn would affect codon usage. Biased gene conver-
14 sion is thought to act similar to directional selection, typically favoring the fixation
15 of G/C alleles [27, 28]. Further, [29, Harrison & Charlesworth] suggested that bi-
16 ased gene conversion affects codon usage in *S. cerevisiae*. ROC SEMPPR, however,
17 does not explicitly account for biased gene conversion. If biased gene conversion is
18 independent of gene expression, as in the case of DNA repair, it will be absorbed
19 in our estimates of ΔM . If instead biased gene conversion forms hotspots, and
20 thus becomes gene specific, it will affect our estimates of protein synthesis ϕ . This
21 might be the case at recombination hotspots. Recombination, however, is very low
22 in the introgressed region (discussed below) [12, 15]. The low recombination rate
23 also indicates that the GC content had to be high before the introgression occurred.

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

32 Despite the violation of the equilibrium assumption, the mutation and selection
33 bias parameters ΔM and $\Delta \eta$ of the introgressed exogenous genes contain informa-

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

10 Previous work indicated that the donor lineage of the exogenous genes has to be
11 a, potentially unknown, Lachancea lineage [12–15]. These previous results, however,
12 are based on species rather than gene trees, ignoring the differential adaptation rate
13 to their novel cellular environment between genes or do not consider lineages outside
14 of the Lachancea clade. Considering the similarity in selection bias (Figure 2b) and
15 our calculation of selection on the exogenous genes (Figure 5b), both of which
16 are free of any assumption about the origin of the exogenous genes, a species tree
17 estimated from the exogenous genes may be biased towards the Lachancea clade.
18 Estimating individual gene trees rather than relying on a species tree provided
19 further evidence that the exogenous genes could originate from a lineage that does
20 not belong to the Lachancea clade. As we highlighted in this study, relatively small
21 sets of genes with a signal of a foreign cellular environment can significantly bias
22 the outcome of a study. The same holds true for phylogenetic inferences [30], and as
23 we showed the signal of the original endogenous cellular environment that shaped
24 CUB is at different stages of decay in high and low expression genes (Figure S8).
25 In summary, our work does not dispute an unknown Lachancea as possible origin,
26 but provides an alternative hypothesis based on the codon usage of the exogenous
27 genes, phylogenetic analysis, and synteny.

28 In terms of understanding the spread of the introgression, we calculated the ex-
29 pected selective cost of codon mismatch between the *L. kluyveri* and *E. gossypii*
30 lineages. Under our working hypothesis, the majority of the introgressed would have
31 imposed a selective cost due to codon mismatch. Nevertheless, $\sim 30\%$ of low expres-
32 sion exogenous genes ($\phi < 1$) appeared to be exapted at the time of the introgres-
33 sion. This exaptation is due to the mutation bias in the endogenous genes matching

1 the selection bias in the exogenous genes for GC ending codons. Our estimate of
2 the selective cost of codon mismatch on the order of -0.0008 . While this selective
3 cost may not seem very large, assuming *L. kluyveri* had a large N_e , the fixation
4 probability of the introgression is the astronomically small value of $\approx 10^{-6952} \approx 0$.
5 While this estimate heavily depends on the working hypothesis that the exogenous
6 genes originated from the *E. gossypii* lineage, we can also calculate the hypothetical
7 fixation probability if the current exogenous genes would introgress into *L. kluyveri*.
8 Our estimate of the current selective cost of the mismatch of codon usage is on the
9 order of -0.0003 . The fixation probability of the current exogenous genes would still
10 be astronomically small $\approx 10^{-2609} \approx 0$. Thus, the basic scenario of an introgression
11 between two yeast species with large N_e and where the introgression solely imposes
12 a selective cost due to codon mismatch is clearly too simplistic.

13 One or more loci with a combined selective advantage on the order of 0.0008
14 or greater would have made the introgression change from disadvantageous to ef-
15 fectively neutral or advantageous. While this scenario seems plausible, it raises
16 the question as to why recombination events did not limit the introgression to
17 only the adaptive loci. A potential answer is the low recombination rate between
18 the endogenous and exogenous regions [12, 15]. Estimates of the recombination
19 rate as measured by crossovers (COs) for *L. kluyveri* are almost four times lower
20 than for *S. cerevisiae* and about half that of *Schizosaccharomyces pombe* (≈ 1.6
21 COs/Mb/meiosis, ≈ 6 COs/Mb/meiosis, ≈ 3 COs/Mb/meiosis) with no observed
22 crossovers in the introgressed region [15], and no observed transposable elements
23 [12]. This is presumably due to the dissimilarity in GC content and/or a lower than
24 average sequence homology between the exogenous region and the one it replaced.
25 A population bottleneck reducing the N_e of the *L. kluyveri* lineage around the time
26 of the introgression could also help explain the spread of the introgression. Compati-
27 ble with these explanation is the possibility of several advantageous loci distributed
28 across the exogenous region drove a rapid selective sweep and/or the population
29 through a bottleneck speciation process.

30 Assuming *E. gossypii* as potential source lineage of the exogenous region, we
31 illustrated how information on codon usage can be used to infer the time since
32 the introgression occurred using our estimates of mutation bias ΔM . The ΔM
33 estimates are well suited for this task as they are free of the influence of selection

1 and unbiased by N_e and other scaling terms, which is in contrast to our estimates of
2 $\Delta\eta$ [7]. Our estimated age of the introgression of $6.2 \pm 1.2 \times 10^8$ generations is ~ 10
3 times longer than a previous minimum estimate by [13] of 5.6×10^7 generations,
4 which was based on the effective population recombination rate and the population
5 mutation parameter [31]. Furthermore, these estimates assume that the current *E.*
6 *gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the
7 time of the introgression. Thus, if the ancestral mutation environments were more
8 similar (dissimilar) at the time of the introgression then our result is an overestimate
9 (underestimate).

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

13 Conclusion

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

27 Materials and Methods

28 Separating Endogenous and Exogenous Genes

29 A GC-rich region was identified by [12] in the *L. kluyveri* genome extending
30 from position 1 to 989,693 of chromosome C. This region was later identified as
31 an introgression by [13]. We obtained the *L. kluyveri* genome from SGD Project
32 <http://www.yeastgenome.org/download-data/> (on 09-27-2014) and the annotation
33 for *L. kluyveri* NRRL Y-12651 (assembly ASM14922v1) from NCBI (on 12-09-

1 2014). We assigned 457 genes located on chromosome C with a location within the
2 ~ 1 Mb window to the exogenous gene set. All other 4864 genes of the *L. kluyveri*
3 genome were assigned to the exogenous genes.

4

5 Model Fitting with ROC SEMPPR

6 ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) [21] and R (3.4.1)
7 [32]. ROC SEMPPR was run from 10 different starting values for at least 250,000
8 iterations and thinned to every 50th iteration. After manual inspection to verify that
9 the MCMC had converged, parameter posterior means, log posterior probability and
10 log likelihood were estimated from the last 500 samples (last 10% of samples).

11

12 Model selection

13 The marginal likelihood of the combined and separated model fits was calculated
14 using a generalized harmonic mean estimator [33]. A variance scaling of 1.1 was
15 used to scale the important density of the estimator. Using the estimated marginal
16 likelihoods, we calculated the Bayes factor to assess model performance. Increases
17 in the variance scaling increase the estimated Bayes factor, therefore we report a
18 conservative Bayes factor bases on a small variance scaling S9.

19

20 Comparing Codon Specific Parameter Estimates and Selecting Candidate lineages

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

$$24 \Delta M_i = \Delta M_{i,1} - \overline{\Delta M_i} \quad (1) \quad 24$$

$$25$$
$$26 \Delta \eta_i = \Delta \eta_{i,1} - \overline{\Delta \eta_i} \quad (2) \quad 26$$

27 Comparison of codon specific parameters (ΔM and $\Delta \eta = 2N_e q(\eta_i - \eta_j)$) was per-
28 formed using the function lmodel2 in the R package lmodel2 (1.7.3) [34] and R
29 version 3.4.1 [32]. The parameter $\Delta \eta$ can be interpreted as the difference in fitness
30 between codon i and j for the average gene with $\phi = 1$ scaled by the effective pop-
31 ulation size N_e , and the selective cost of an ATP q [4, 7]. Type II regression was
32 performed with re-centered parameter estimates, accounting for noise in dependent
33 and independent variable [23].

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

6 Phylogenetic Analysis

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

18 Synteny Comparison

19 We obtained complete genome sequences for all 10 candidate lineages (Table 2)
20 from NCBI (on: 02-05-2017). Genomes were aligned and checked for synteny using
21 SyMAP (4.2) with default settings [35, 36]. We assess synteny as percentage coverage
22 of the exogenous gene region.

23 Estimating Age of Introgression

24 We modeled the change in codon frequency over time using an exponential model
25 for all two codon amino acids. While our approach is equivalent to [37], we want
26 to explicitly state the relationship between the change in codon frequency c_1 as a
27 function of mutation bias ΔM as

$$30 \quad \frac{dc_1}{dt} = -\mu_{1,2}c_1 - \mu_{2,1}(1 - c_1) \quad (3)$$

31 where $\mu_{i,j}$ is the rate at which codon i mutates to codon j and c_1 is the fre-
32 quency of the reference codon. Initial codon frequencies $c_1(0)$ for each codon

1 family were taken from our mutation parameter estimates for *E. gossypii* where
 2 $c_1(0) = \exp[\Delta M_{\text{gos}}]/(1 + \exp[\Delta M_{\text{gos}}])$. Our estimates of ΔM_{endo} can be used to
 3 calculate the steady state of equation 3 were $\frac{dc_1}{dt} = 0$ to obtain the equality

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

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

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

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

12 Equation 5 was solved with a mutation rate $\mu_{2,1}$ of 3.8×10^{-10} per nucleotide per
 13 generation [38]. Current codon frequencies for each codon family were taken from
 14 our estimates of ΔM from the exogenous genes. Mathematica (11.3) [39] was used
 15 to calculate the time t_{intro} it takes for the initial codon frequencies $c_1(0)$ for each
 16 codon family to equal the current exogenous codon frequencies. The same equation
 17 was used to determine the time t_{decay} at which the signal of the exogenous cellular
 18 environment has decayed to within 1% of the endogenous environment.

20 Estimating Selection against Codon Mismatch

21 In order to estimate the selection against codon mismatch, we had to make three
 22 key assumptions. First, we assumed that the current exogenous amino acid sequence
 23 of a gene is representative of its ancestral state and the replaced endogenous gene
 24 it replaced. Second, we assume that the currently observed cellular environment of
 25 *E. gossypii* reflects the cellular environment that the exogenous genes experienced
 26 before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy
 27 of selection between the cellular environments due to differences in either effective
 28 population size N_e or the selective cost of an ATP q of the source lineage and *L.*
 29 *kluyveri* can be expressed as a scaling constant and that protein synthesis rate ϕ
 30 has not changed between the replaced endogenous and the introgressed exogenous
 31 genes. Using estimates for $N_e = 1.36 \times 10^7$ [26] for *Saccharomyces paradoxus* we
 32 scale our estimates of $\Delta\eta$ which explicitly contains the effective population size N_e
 33 [7] and define $\Delta\eta' = \frac{\Delta\eta}{N_e}$.

1 All of our genome parameter estimations are scaled by lineage specific effects
 2 such as N_e , the average, absolute gene expression level, and/or the proportionate
 3 fitness value of an ATP. In order to account for these genome specific differences in
 4 scaling, we scale the difference in the efficacy of selection on codon usage between
 5 the donor lineage and *L. kluyveri* using a linear scaling factor κ . As $\Delta\eta$ is defined as
 6 $\Delta\eta = 2N_e q(\eta_i - \eta_j)$, we cannot distinguish if κ is a scaling on protein synthesis rate
 7 ϕ , effective population size N_e , or the selective cost of an ATP q [4, 7]. We calculated
 8 the selection against each genes codon mismatch assuming additive fitness effects
 9 as

$$10 \\
 11 s_g = \sum_{i=1}^{L_g} -\kappa \phi_g \Delta\eta'_i \quad (6) \\
 12$$

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

25 However, since we are unable to observe codon sequences of the replaced en-
 26 dogenous genes and for the exogenous genes at the time of introgression, instead
 27 of summing over the sequence, we calculate the expected codon count $E[n_{g,i}]$ for
 28 codon i in gene g simply as the probability of observing codon i multiplied by the
 29 number of times the corresponding amino acids is observed in gene g , yielding:
 30

$$31 E[n_{g,i}] = P(c_i | \Delta M, \Delta\eta, \phi) \times m_{a_i} \\
 32 = \frac{\exp[-\Delta M_i - \Delta\eta_i \phi_g]}{\sum_j^C \exp[-\Delta M_j - \Delta\eta_j \phi_g]} \times m_{a_i} \\
 33$$

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

$$4 \\ 5 s_g = \sum_{i=1}^C -\kappa \phi_g \Delta \eta'_i E[n_{g,i}] \quad (7) \\ 6$$

7 We report the selection due to mismatched codon usage of the introgression as
8 $\Delta s_g = s_{\text{intro},g} - s_{\text{endo},g}$ where $s_{\text{intro},g}$ is the selection against an introgressed gene g
9 either at the time of the introgression or presently.

10 Acknowledgments

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

12 Availability of data and materials

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

15 Authors' contributions

16 CL and MAG initiated the study. CL collected and analyzed the data and wrote the manuscript. MAG and BCO
17 edited the manuscript. CL, MAG, BCO, and RZ contributed to the data analysis and acquiring of funding. All
18 Authors approved the final manuscript.

19 Funding

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

26 Ethics approval and consent to participate

27 Not applicable

28 Consent for publication

29 Not applicable

30 Competing interests

31 The authors declare that they have no competing interests.

32 Author details

33 ¹Department of Ecology & Evolutionary Biology, University of Tennessee, 37996, Knoxville, TN, USA. ²National
34 Institute for Mathematical and Biological Synthesis, 37996, Knoxville, TN, USA. ³Max-Planck Institute of
35 Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, 01307, Dresden, Germany. ⁴Department of Business
36 Analytics and Statistics, University of Tennessee, 37996, Knoxville, TN, USA.

37 References

- 38 1. Gouy, M., Gautier, C.: Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Research* **10**,
39 7055–7074 (1982)
- 40 2. Ikemura, T.: Codon usage and tRNA content in unicellular and multicellular organisms. *Molecular Biology and
41 Evolution* **2**, 13–34 (1985)
- 42 3. Bulmer, M.: The selection-mutation-drift theory of synonymous codon usage. *Genetics* **129**, 897–907 (1990)

- 1 4. Gilchrist, M.A.: Combining models of protein translation and population genetics to predict protein production
2 rates from codon usage patterns. *Molecular Biology and Evolution* **24**(11), 2362–2372 (2007) 1
3 5. Shah, P., Gilchrist, M.A.: Explaining complex codon usage patterns with selection for translational efficiency,
4 mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences U.S.A* **108**(25),
5 10231–10236 (2011) 2
6 6. Wallace, E.W., Airoldi, E.M., Drummond, D.A.: Estimating selection on synonymous codon usage from noisy
7 experimental data. *Molecular Biology and Evolution* **30**, 1438–1453 (2013) 3
8 7. Gilchrist, M.A., Chen, W.C., Shah, P., Landerer, C.L., Zaretzki, R.: Estimating gene expression and
9 codon-specific translational efficiencies, mutation biases, and selection coefficients from genomic data alone.
10 *Genome Biology and Evolution* **7**, 1559–1579 (2015) 4
11 8. Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., Danchin, A.: Evidence for horizontal gene transfer in
12 *Escherichia coli* speciation. *Journal of Molecular Biology* **222**(4), 851–856 (1991) 5
13 9. Lawrence, J.G., Ochman, H.: Amelioration of bacterial genomes: Rates of change and exchange. *Journal of
14 Molecular Biology* **44**, 383–397 (1997) 6
15 10. Marcket-Houben, M., Gabaldón, T.: Beyond the whole-genome duplication: Phylogenetic evidence for an ancient
16 interspecies hybridization in the baker's yeast lineage. *PLoS Biology* **13**(8), 1002220 (2015) 7
17 11. Beimforde, C., Feldberg, K., Nylander, S., Rikkinen, J., Tuovila, H., Dörfelt, H., Gube, M., Jackson, D.J.,
18 Reitner, J., Seyfullah, L.J., Schmidt, A.R.: Estimating the phanerozoic history of the ascomycota lineages:
19 combining fossil and molecular data. *Mol. Phylogenet. Evol.* **78**, 386–398 (2014) 8
20 12. Payen, C., Fischer, G., Marck, C., Proux, C., Sherman, D.J., Coppée, J.-Y., Johnston, M., Dujon, B.,
21 Neuvéglise, C.: Unusual composition of a yeast chromosome arm is associated with its delayed replication.
22 *Genome Research* **19**(10), 1710–1721 (2009) 9
23 13. Friedrich, A., Reiser, C., Fischer, G., Schacherer, J.: Population genomics reveals chromosome-scale
24 heterogeneous evolution in a protoplid yeast. *Molecular Biology and Evolution* **32**(1), 184–192 (2015) 10
25 14. Vakirlis, N., Sarilar, V., Drillon, G., Fleiss, A., Agier, N., Meyniel, J.-P., Blanpain, L., Carbone, A., Devillers, H.,
26 Dubois, K., Gillet-Markowska, A., Graziani, S., Huu-Vang, N., Poirel, M., Reisser, C., Schott, J., Schacherer,
27 J., Lafontaine, I., Llorente, B., Neuvéglise, C., Fischer, G.: Reconstruction of ancestral chromosome
28 architecture and gene repertoire reveals principles of genome evolution in a model yeast genus. *Genome
29 research* **26**(7), 918–932 (2016) 11
30 15. 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) 12
33 16. Sharp, P.M., Li, W.H.: The codon adaptation index - a measure of directional synonymous codon usage bias,
34 and its potential applications. *Nucleic Acids Research* **15**, 1281–1295 (1987) 13
35 17. Wright, F.: The 'effective number of codons' used in a gene. *Genet* **87**, 23–29 (1990) 14
36 18. dos Reis, M., Savva, R., Wernisch, L.: Solving the riddle of codon usage preferences: a test for translational
37 selection. *Nucleic Acids Research* **32**(17), 5036–5044 (2004) 15
38 19. Cope, A.L., Hettich, R.L., Gilchrist, M.A.: Quantifying codon usage in signal peptides: Gene expression and
39 amino acid usage explain apparent selection for inefficient codons. *Biochimica et Biophysica Acta (BBA) -
40 Biomembranes* **1860**(12), 2479–2485 (2018) 16
41 20. Shen, X.X., Opulente, D.A., Kominek, J., Zhou, X., Steenwyk, J.L., Buh, K.V., Haase, M.A.B., Wisecaver,
42 J.H., Wang, M., Doering, D.T., Boudouris, J.T., Schneider, R.M., Langdon, Q.K., Ohkuma, M., Endoh, R.,
43 Takashima, M., Manabe, R., Čadež, N., Libkind, D., Rosa, C., DeVirgilio, J., Hulfachor, A.B., Groenewald, M.,
44 Kurtzman, C., Hittinger, C.T., Rokas, A.: Tempo and mode of genome evolution in the budding yeast
45 subphylum. *Cell* **175**(6), 1533–154520 (2018) 17
46 21. Landerer, C., Cope, A., Zaretzki, R., Gilchrist, M.A.: AnaCoDa: analyzing codon data with bayesian mixture
47 models. *Bioinformatics* **34**(14), 2496–2498 (2018) 18
48 22. Tsankov, A.M., Thompson, D.A., Socha, A., Regev, A., Rando, O.J.: The role of nucleosome positioning in the
49 evolution of gene regulation. *PLoS Biol* **8**(7), 1000414 (2010) 19
50 23. Sokal, R.R., Rohlf, F.J.: *Biometry - The principles and practice of statistics in biological*, pp. 547–555. W. H.
51 Freeman, New York, NY (1981) 20

- 1 24. Nguyen, L.T., Schmidt, H.A., von Haeseler, A., Minh, B.Q.: Iq-tree: A fast and effective stochastic algorithm
2 for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **32**(1), 268–274 (2015) 1
- 3 25. Sella, G., Hirsh, A.E.: The application of statistical physics to evolutionary biology. *Proceedings of the National
Academy of Sciences of the United States of America* **102**, 9541–9546 (2005) 2
- 4 26. Wagner, A.: Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* **22**,
1365–1374 (2005) 3
- 5 27. Nagylaki, T.: Evolution of a finite population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**,
6278–6281 (1983) 4
- 6 28. Nagylaki, T.: Evolution of a large population under gene conversion. *Proc. Natl. Acad. Sci. U. S. A.* **80**,
5941–5945 (1983) 5
- 7 29. Harrison, R.J., Charlesworth, B.: Biased gene conversion affects patterns of codon usage and amino acid usage
8 in the *Saccharomyces sensu stricto* group of yeasts. *Molecular Biology and Evolution* **28**(1), 117–129 (2011) 6
- 9 30. Salichos, L., Rokas, A.: Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature*
497, 327–331 (2013) 7
- 10 31. 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) 8
- 11 32. 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/> 9
- 12 33. Gronau, Q.F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D.S., Forster, J.J.,
Wagenmakers, E.J., Steingroever, H.: A tutorial on bridge sampling. *Journal of Mathematical Psychology* **81**,
80–97 (2017) 10
- 13 34. Legendre, P.: Lmodel2: Model II Regression. (2018). R package version 1.7-3.
<https://CRAN.R-project.org/package=lmodel2> 11
- 14 35. Soderlund, C., Nelson, W., Shoemaker, A., Paterson, A.: Symap A system for discovering and viewing syntenic
regions of fpc maps. *Genome Research* **16**, 1159–1168 (2006) 12
- 15 36. Soderlund, C., Bomhoff, M., Nelson, W.: Symap v3.4: a turnkey synteny system with application to plant
genomes. *Nucleic Acids Research* **39**(10), 68 (2011) 13
- 16 37. 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) 14
- 17 38. Lang, G.I., Murray, A.W.: Estimating the per-base-pair mutation rate in the yeast *Saccharomyces cerevisiae*.
Genetics **178**(1), 67–82 (2008) 15
- 18 39. Wolfram Research Inc.: Mathematica 11. (2017). <http://www.wolfram.com> 16
- 19 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. kluveri*, and in the combined *L. kluveri* genome without accounting for the
 two cellular environments.

	Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	Combined	
7	Ala A	GCG	GCA	GCG	GCG	6
8	Cys C	TGC	TGT	TGC	TGC	7
9	Asp D	GAC	GAT	GAC	GAC	8
10	Glu E	GAG	GAA	GAG	GAG	9
11	Phe F	TTC	TTT	TTT	TTT	10
12	Gly G	GGC	GGT	GGC	GGC	11
13	His H	CAC	CAT	CAC	CAC	12
14	Ile I	ATC	ATT	ATC	ATA	13
15	Lys K	AAG	AAA	AAG	AAA	14
16	Leu L	CTG	TTG	CTG	CTG	15
17	Asn N	AAC	AAT	AAC	AAT	16
18	Pro P	CCG	CCA	CCG	CCG	17
19	Gln Q	CAG	CAA	CAG	CAG	18
20	Arg R	CGC	AGA	AGG	CGG	19
21	Ser ₄ S	TCG	TCT	TCG	TCG	20
22	Thr T	ACG	ACA	ACG	ACG	21
23	Val V	GTG	GTT	GTG	GTG	22
24	Tyr Y	TAC	TAT	TAC	TAC	23
25	Ser ₂ Z	AGC	AGT	AGC	AGC	24
26						25
27						26
28						27
29						28
30						29
31						30
32						31
33						32

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	10				
11	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

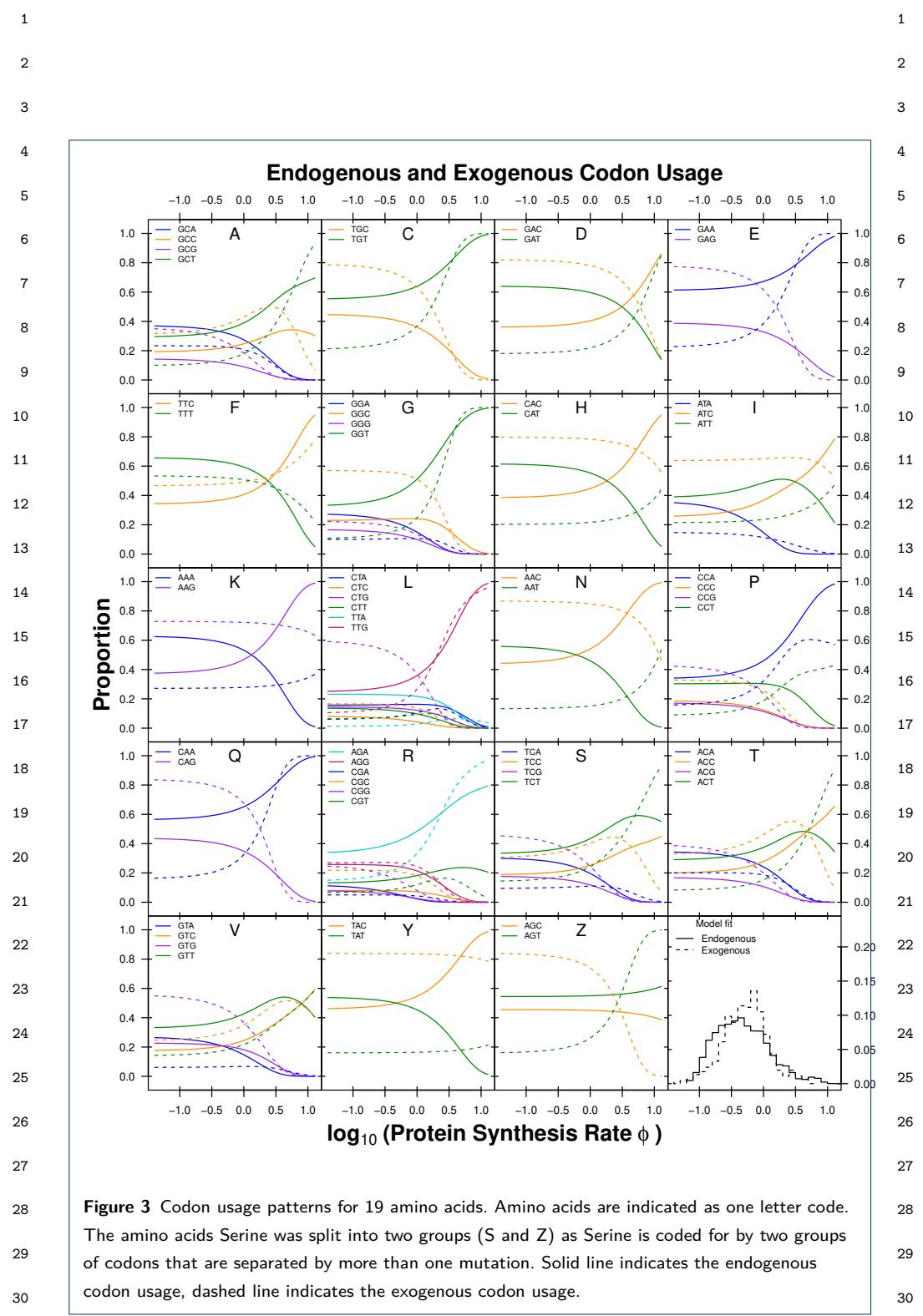
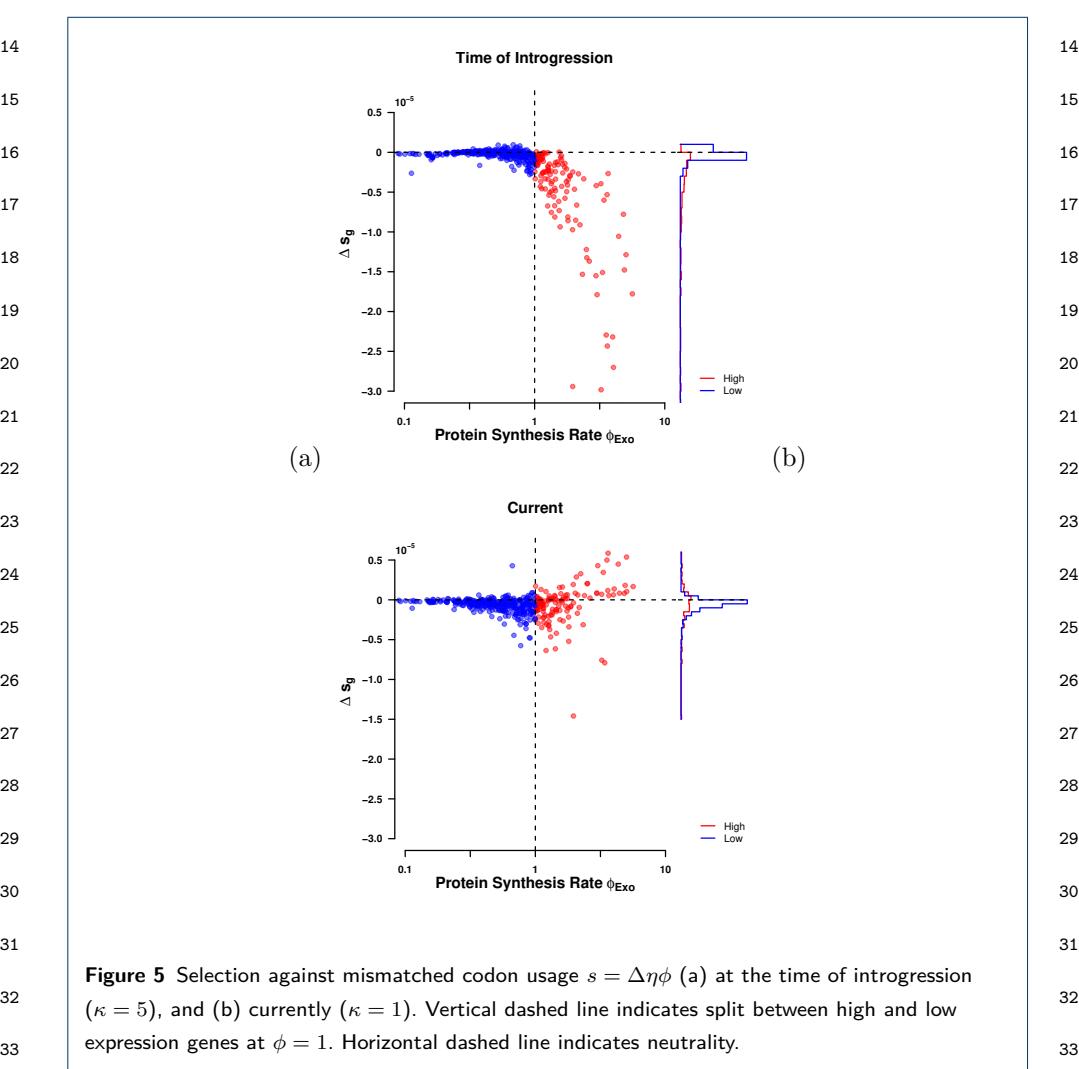
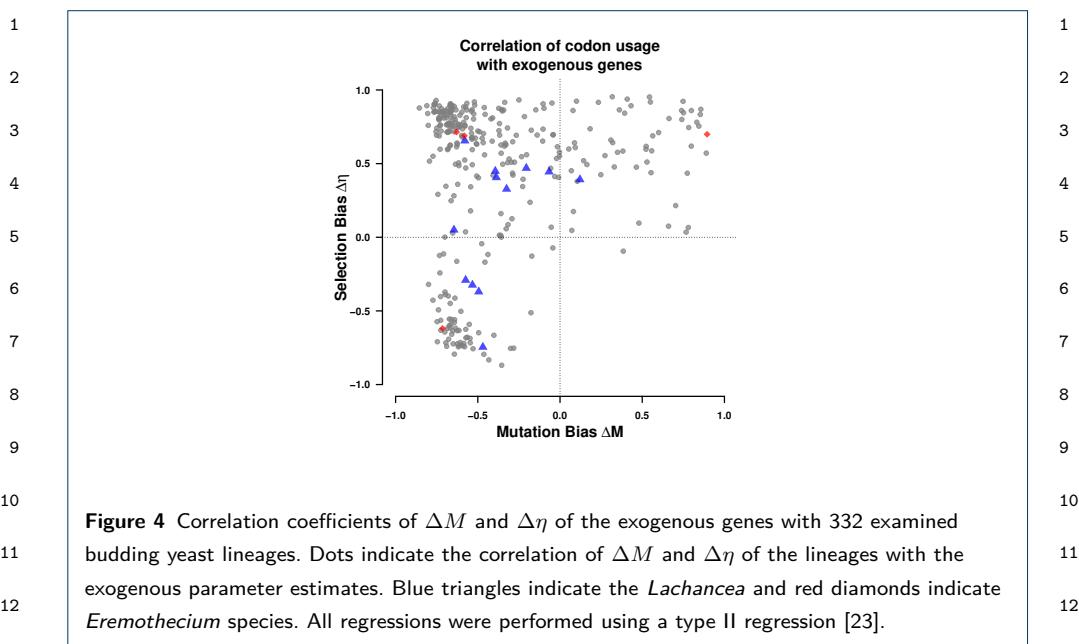
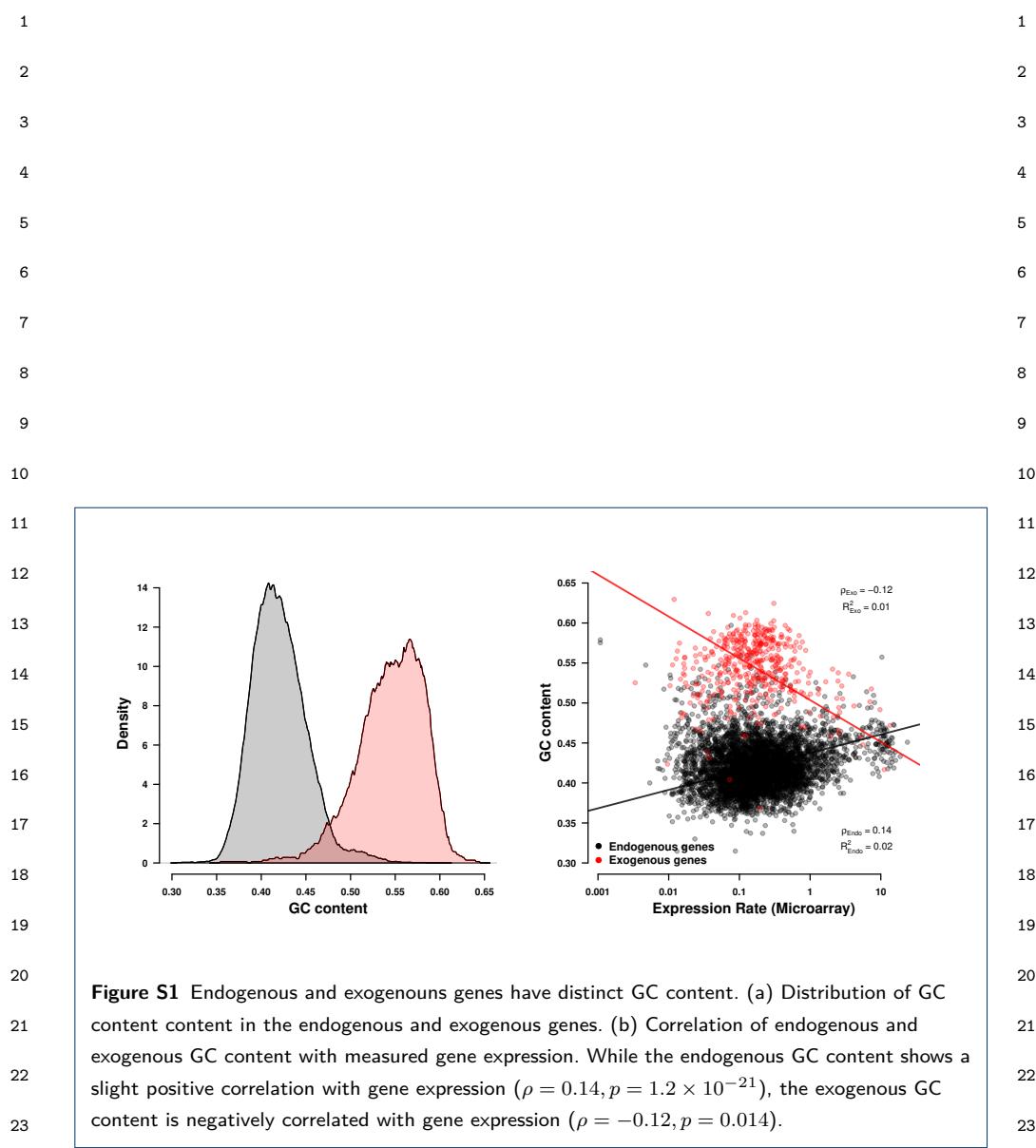


Figure 3 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dashed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





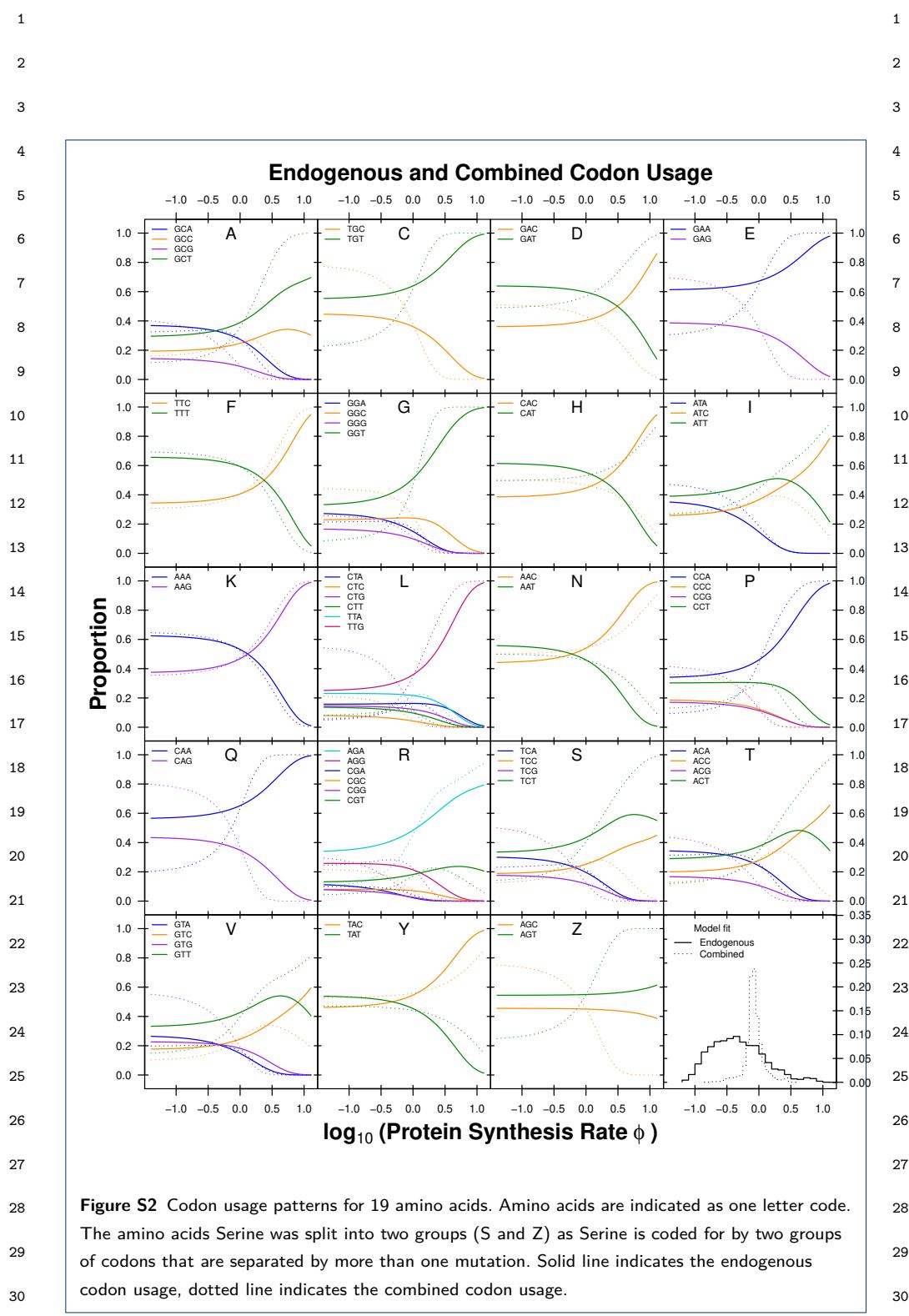


Figure S2 Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dotted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

