

RESEARCH

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

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

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  
Correspondance

**Abstract**

**Background:** For decades, codon usage has been used as a measure of adaptation for translational efficiency of a gene’s coding sequence. These patterns of codon usage reflect both the selective and mutational environment in which the coding sequences evolved. Over this same period, gene transfer between lineages has become widely recognized as an important biological phenomenon. Nevertheless, most studies of codon usage implicitly assume that all genes within a genome evolved under the same selective and mutational environment, an assumption violated when introgression occurs.

**Results:** In order to better understand the effects of introgression on codon usage patterns and vice versa, we examine the patterns of codon usage in *Lachancea kluyveri*, a yeast which has experienced a large introgression. We quantify the effects of mutation bias and selection for translation efficiency on the codon usage pattern of the endogenous and introgressed exogenous genes using a Bayesian mixture model, ROC SEMPPR, which is built on mechanistic assumptions of protein synthesis and grounded in population genetics. We find substantial differences in codon usage between the endogenous and exogenous genes, and show that these differences can be largely attributed to a shift in mutation bias favoring A/T ending codons in the endogenous genes to C/G ending codons in the exogenous genes. Recognizing the two different signatures of mutation bias and selection improves our ability to predict protein synthesis rate by 42% and allowed us to accurately assess endogenous codon preferences. In addition, using our estimates of mutation bias and selection, we identify *Eremothecium gossypii* as the closest relative to the exogenous genes, providing an alternative hypothesis about the origin of the exogenous genes, estimate the introgression occurred  $\sim 6 \times 10^8$  generation ago, and estimate its historic and current selection against mismatched codon usage.

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

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

**Background**

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

<sup>1</sup>mutation bias is largely determined by the organism's internal or cellular environ-  
<sup>2</sup>ment, such as their DNA repair genes or UV exposure. While this mutation bias<sup>2</sup>  
<sup>3</sup>is an omnipresent evolutionary force, its impact can be obscured or amplified by<sup>3</sup>  
<sup>4</sup>selection. In contrast, the signature of selection on codon usage is largely deter-<sup>4</sup>  
<sup>5</sup>mined by an organism's cellular environment alone, such as its tRNA species, their<sup>5</sup>  
<sup>6</sup>copy number, and their post-transcriptional modifications. The strength of selec-<sup>6</sup>  
<sup>7</sup>tion on the codon usage of an individual gene is largely determined by its expression<sup>7</sup>  
<sup>8</sup>and synthesis rate which, in turn, is largely determined by the organism's external<sup>8</sup>  
<sup>9</sup>environment. In general, the strength of selection on codon usage increases with<sup>9</sup>  
<sup>10</sup>its expression level [1–3], specifically its protein synthesis rate [4]. Thus as protein<sup>10</sup>  
<sup>11</sup>synthesis increases, codon usage shifts from a process dominated by mutation to a<sup>11</sup>  
<sup>12</sup>process dominated by selection. The overall efficacy of selection on codon usage is<sup>12</sup>  
<sup>13</sup>a function of the organism's effective population size  $N_e$  which, in turn, is largely<sup>13</sup>  
<sup>14</sup>determined by its external environment. ROC SEMPPR allows us disentangle the<sup>14</sup>  
<sup>15</sup>evolutionary forces responsible for the patterns of codon usage bias (CUB) encoded<sup>15</sup>  
<sup>16</sup>in an species' genome, by explicitly modeling the combined evolutionary forces of<sup>16</sup>  
<sup>17</sup>mutation, selection, and drift [4–7]. In turn, these evolutionary forces should pro-<sup>17</sup>  
<sup>18</sup>vide biologically meaningful information about the lineage's historical cellular and<sup>18</sup>  
<sup>19</sup>external environment. 19

<sup>20</sup>Most studies implicitly assume that the CUB of a genome is shaped by a single<sup>20</sup>  
<sup>21</sup>cellular environment. As genes are horizontally transferred, introgress, or combined<sup>21</sup>  
<sup>22</sup>to form novel hybrid species, one would expect to see the influence of multiple cel-<sup>22</sup>  
<sup>23</sup>lular environments on a genomes codon usage pattern [8, 9]. Given that transferred<sup>23</sup>  
<sup>24</sup>genes are likely to be less adapted than endogenous genes to their new cellular en-<sup>24</sup>  
<sup>25</sup>vironment, we expect a greater selection against mismatched codon usage in trans-<sup>25</sup>  
<sup>26</sup>ferred genes if donor and recipient environment differ greatly in their selection bias,<sup>26</sup>  
<sup>27</sup>making such transfers less likely. More practically, if differences in codon usage of<sup>27</sup>  
<sup>28</sup>transferred genes are unaccounted for, they may distort the interpretation of codon<sup>28</sup>  
<sup>29</sup>usage patterns. Such distortion could lead to the wrong inference of codon prefer-<sup>29</sup>  
<sup>30</sup>ence for an amino acid [5, 7], underestimate the variation in protein synthesis rate,<sup>30</sup>  
<sup>31</sup>or influence mutation estimates when analyzing a genome. 31

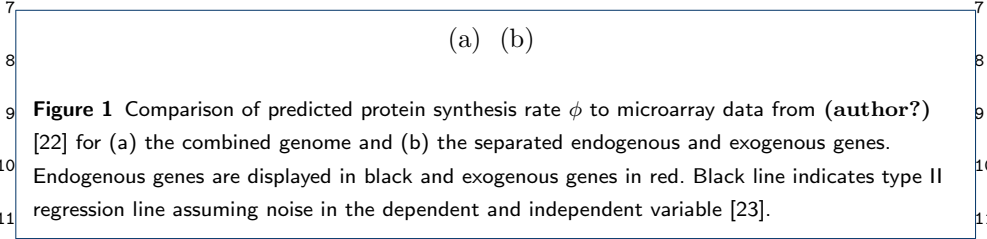
<sup>32</sup>To illustrate these ideas, we analyze the CUB of the genome of *Lachancea kluyveri*,<sup>32</sup>  
<sup>33</sup>which is sister to all other *Lachancea* species. The *Lachancea* clade diverged from the<sup>33</sup>

<sup>1</sup>Saccharomyces clade, prior to its whole genome duplication  $\sim 100$  Mya ago [10, 11].<sup>1</sup>  
<sup>2</sup>Since that time, *L. kluyveri* has experienced a large introgression of exogenous genes<sup>2</sup>  
<sup>3</sup>which is found in all of its populations [12], but in no other known Lachancea species<sup>3</sup>  
<sup>4</sup>[13]. The introgression replaced the left arm of the C chromosome and displays a<sup>4</sup>  
<sup>5</sup>13% higher GC content than the endogenous *L. kluyveri* genome [12, 14]. Previous<sup>5</sup>  
<sup>6</sup>studies suggest that the source of the introgression is likely a currently unknown<sup>6</sup>  
<sup>7</sup>or potentially extinct Lachancea lineage based on gene concatenation or synteny<sup>7</sup>  
<sup>8</sup>relationships [12–15]. These characteristics make *L. kluyveri* an ideal model to study<sup>8</sup>  
<sup>9</sup>the effects of an introgressed cellular environment and the resulting mismatch in<sup>9</sup>  
<sup>10</sup>codon usage.<sup>10</sup>

<sup>11</sup> Using ROC SEMPPR, a Bayesian population genetics model based on a mecha-<sup>11</sup>  
<sup>12</sup>nistic description of ribosome movement along an mRNA, allows us to quantify the<sup>12</sup>  
<sup>13</sup>cellular environment in which genes have evolved by separately estimating the ef-<sup>13</sup>  
<sup>14</sup>fects of mutation bias and selection bias on codon usage. ROC SEMPPR’s resulting<sup>14</sup>  
<sup>15</sup>predictions of protein synthesis rates have been shown to be on par with laboratory<sup>15</sup>  
<sup>16</sup>measurements [5, 7]. In contrast to often used heuristic approaches to study codon<sup>16</sup>  
<sup>17</sup>usage [16–18], ROC SEMPPR explicitly incorporates and distinguishes between<sup>17</sup>  
<sup>18</sup>mutation and selection effects on codon usage and properly weights by amino acid<sup>18</sup>  
<sup>19</sup>usage [19]. We use ROC SEMPPR to independently describe two cellular environ-<sup>19</sup>  
<sup>20</sup>ments reflected in the *L. kluyveri* genome; the signature of the current environment<sup>20</sup>  
<sup>21</sup>in the endogenous genes and the decaying signature of the exogenous environment<sup>21</sup>  
<sup>22</sup>in the introgressed genes. Our results indicate that the difference in GC content<sup>22</sup>  
<sup>23</sup>between endogenous and exogenous genes is mostly due to the differences in muta-<sup>23</sup>  
<sup>24</sup>tion bias of their ancestral environments. Accounting for these different signatures<sup>24</sup>  
<sup>25</sup>of mutation bias and selection bias of the endogenous and exogenous sets of genes<sup>25</sup>  
<sup>26</sup>substantially improves our ability to predict present day protein synthesis rates.<sup>26</sup>  
<sup>27</sup>These endogenous and exogenous gene set specific estimates of mutation bias and<sup>27</sup>  
<sup>28</sup>selection bias, in turn, allow us to address more refined questions of biological im-<sup>28</sup>  
<sup>29</sup>portance. For example, they allow us to provide an alternative hypothesis about the<sup>29</sup>  
<sup>30</sup>origin of the introgression and identify *E. gossypii* as the nearest sampled relative<sup>30</sup>  
<sup>31</sup>of the source of the introgressed genes out of the 332 budding yeast lineages with<sup>31</sup>  
<sup>32</sup>sequenced genomes [20]. While this hypothesis is in contrast previous work [12–15],<sup>32</sup>  
<sup>33</sup>we find support for it in gene trees and synteny. We also estimate the age of the<sup>33</sup>

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

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



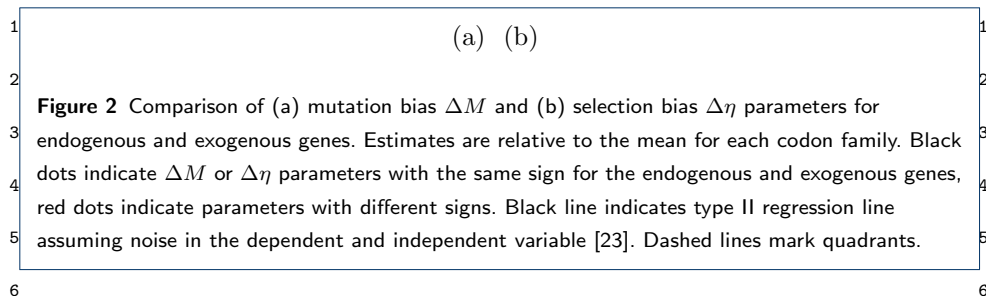
**Figure 1** Comparison of predicted protein synthesis rate  $\phi$  to microarray data from (author?) [22] for (a) the combined genome and (b) the separated endogenous and exogenous genes. Endogenous genes are displayed in black and exogenous genes in red. Black line indicates type II regression line assuming noise in the dependent and independent variable [23].

introgression to be on the order of 0.2 - 1.7 Mya, estimate the selection against these genes, both at the time of introgression and now, and predict a detectable signature of CUB to persist in the introgressed genes for another 0.3 - 2.8 Mya, highlighting the sensitivity of our approach.

## Results

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

We used our software package AnaCoDa [21] to compare model fits of ROC SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two sets of 4,864 endogenous and 497 exogenous genes. ROC SEMPPR is a statistical model that relates the effects of mutation bias  $\Delta M$  and selection bias  $\Delta \eta$  between synonymous codons, and protein synthesis rate  $\phi$  to explain the observed codon usage patterns. 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 ( $K = \exp(42, 294)$ ; 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 proxy for protein synthesis. Specifically, the explanatory power between our predictions and observed values improved by  $\sim 42\%$ , from  $R^2 = 0.33$  to 0.46 (Figure 1).



## Comparing Differences in the Endogenous and Exogenous Codon Usage

To better understand the differences in the endogenous and exogenous cellular environments, we compared our parameter estimates of mutation bias  $\Delta M$  and selection  $\Delta \eta$  for the two sets of genes. Our estimates of  $\Delta M$  for the endogenous and exogenous genes were negatively correlated ( $\rho = -0.49$ ), indicating weak similarity with only  $\sim 5\%$  of the codons share the same sign between the two mutation environments (Figure 2a). Overall, the endogenous genes only show a selection preference for C and G ending codons in  $\sim 58\%$  of the codon families. In contrast, the exogenous genes display a strong preference for A and T ending codons in  $\sim 89\%$  of the codon families.

For example, the endogenous genes show a mutational bias for A and T ending codons in  $\sim 95\%$  of the codon families (the exception being Phe, F). The exogenous genes display an equally consistent mutational bias towards C and G ending codons (Table S1). In contrast to  $\Delta M$ , our estimates of  $\Delta \eta$  for the endogenous and exogenous genes were positively correlated ( $\rho = 0.69$ ) and showing the same sign in  $\sim 53\%$  of codons between the two selection environments (Figure 2). ROC SEMPPR constraints  $E[\phi] = 1$ , allowing us to interpret  $\Delta \eta$  as selection on codon usage of the average gene with  $\phi = 1$  and gives us the ability to compare the efficacy of selection  $sN_e$  across genomes.

We find that the efficacy of selection within each codon family differs between sets of genes. The difference in codon usage between endogenous and exogenous genes is striking as some amino acids have opposite codon preferences. As a result, our estimates of the optimal codon differ in nine cases between endogenous and exogenous genes (Figure 3, Table S2). For example, the usage of the Asparagine (Asn, N) codon AAC is increased in highly expressed endogenous genes but the same codon is depleted in highly expressed exogenous genes. For Aspartic acid (Asp, D), the combined genome shows the same codon preference in highly expressed genes

**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.

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 disproportional effect on the model fit (Figure S1, S2).

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

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

### Determining Source of Exogenous Genes

We combined our estimates of mutation bias  $\Delta M$  and selection bias  $\Delta\eta$  with synteny information and searched for potential source lineages of the introgressed exogenous region. We examined 332 budding yeasts [20] and, identified the ten lineages with

**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  $\Delta M$  and  $\Delta \eta$ , 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 zsoitii</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>Middelhovenomyces 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

the highest correlation for the  $\Delta M$  parameters as potential source lineages (Figure 13, Table 2). We used  $\Delta M$  to identify candidate lineages as the endogenous and exogenous genes show greater dissimilarity in mutation bias than in selection bias. 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 in our comparison of codon families, however, there was no need to exclude Serine as well as CTG is not a one step neighbor of the remaining Serine codons. The endogenous *L. kluyveri* genome exhibits codon usage very similar to most (77 %) yeast lineages examined, indicating that most of the examined yeasts share a similar codon usage (Figure S3). Only ~ 17% of all examined yeast show a positive correlation in both,  $\Delta M$  and  $\Delta \eta$  with the exogenous genes, whereas the vast majority of lineages (~ 83%) show a negative correlation for  $\Delta M$ , only 21 % show a negative correlation for  $\Delta \eta$ .

Comparing synteny between the exogenous genes, which are restricted to the left arm of chromosome C, and the determined candidate yeast species we find that *E. gossypii* is the only species that displays high synteny (Table 2). Furthermore, the synteny relationship between the exogenous region and other yeasts appears to be limited to Saccharomycetaceae clade. Given these results, we conclude that of the 332 examined yeast lineages the *E. gossypii* lineage is the most likely source of the introgressed exogenous genes. This result is in contrast to previous studies which studied the exogenous genes and chromosome recombination in the Lachancea clade



**Figure 4** Correlation coefficients of  $\Delta M$  and  $\Delta \eta$  of the exogenous genes with 332 examined budding yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta \eta$  of the lineages with the exogenous parameter estimates. Blue triangles indicate the *Lachancea* and red diamonds indicate *Eremothecium* species. All regressions were performed using a type II regression assuming noise in the dependent and independent variable [23].

and concluded that the exogenous region originated from within the *Lachancea* clade [12–14]. To validate our results, we identified 121 genes in our dataset [20] with homologous gene in *E. gossypii* and *L. thermotolerance* and used IQTree [24] to infer the phylogenetic relationship of the exogenous genes. Our results show that ~ 60% of exogenous genes (73/121) are more closely related to *E. gossypii* than to other *Lachancea*. Interestingly, our results also indicate that codon usage does not necessarily correlate with phylogenetic distance (Table 2).

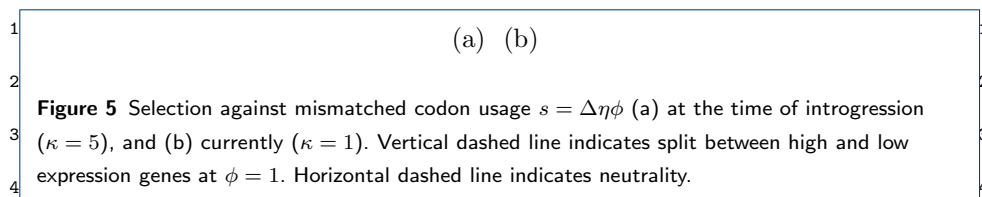
#### Estimating Introgression Age

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 (author?) [12].

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

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

We define the selection against inefficient codon usage as the difference between the fitness on the log scale of an expected, replaced endogenous gene and the exogenous gene,  $s \propto \phi \Delta \eta$  due to the mismatch in codon usage parameters (See Methods for details). As the introgression occurred before the diversification of *L. kluyveri* and



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

Using our estimates of  $\Delta M$  and  $\Delta\eta$  from the endogenous genes and assuming the current exogenous amino acid composition of genes is representative of the replaced endogenous genes, we estimate the selection against the exogenous genes at the time of introgression (Figure 5a) and currently (Figure 5b). Estimates of selection bias for the exogenous genes show that, while well correlated with the endogenous genes, only nine amino acids share the same selectively preferred codon. Exogenous genes are, therefore, expected to represent a significant reduction in fitness for *L. kluyveri* due to mismatch in codon usage. We estimate that the selection against the exogenous genes due to mismatched codon usage to have been  $\Delta s \approx -0.0008$  at the time of the introgression and  $\approx -0.0003$  today. This reduction in  $\Delta s$  is primarily due to adaptive changes to the codon usage of the most highly expressed, introgressed genes (Figures 5a & S6). Based on the selection against the codon mismatch at the time of the introgression and assuming an effective population size  $N_e$  on the order of  $10^7$  [25], we approximate a fixation probability of  $(1 - \exp[-\Delta s]) / (1 - \exp[-2\Delta s N_e]) \approx 10^{-6952}$  [26] for the exogenous genes. Clearly, the possibility of fixation under this simple scenario is effectively zero (See Discussion).

## Discussion

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

Previous work by [14] showed an increased preference for GC rich codons in the exogenous genes but our results provide more nuanced insights by separating the effects of mutation bias and selection. We are able to show that the difference in GC content between endogenous and exogenous genes is mostly due to differences in mutation bias as 95% of exogenous codon families show a strong mutation bias towards GC ending codons (Table S1). However, the exogenous genes show a selective preference for AT ending codons for 90% of codon families (Table S2). Acknowledging the increased mutation bias towards GC ending codons and the difference in strength of selection between endogenous and exogenous genes by separating them also improves our estimates of protein synthesis rate  $\phi$  by 42% relative to the full genome estimate ( $R^2 = 0.46$  vs.  $0.32$ , respectively).

The mutation and selection bias parameters  $\Delta M$  and  $\Delta\eta$  of the introgressed exogenous genes contain information, albeit decaying, about its previous cellular environment. We selected the top ten lineages with the highest similarity in  $\Delta M$  to see if our parameters estimates would allow us to identify a potential source lineage. The synteny relationship of these lineages with the exogenous genes was calculated as a point of comparison as it provides orthogonal information to our parameter estimates. Synteny with the exogenous genes is limited to the Saccharomycetaceae clade, excluding all of the potential source lineages identified using codon usage but

<sup>1</sup>*E. gossypii* (Table 2). Interestingly, this also showed that similarity in codon usage<sup>1</sup>  
<sup>2</sup>does not correlate with phylogenetic distance.<sup>2</sup>

<sup>3</sup> Previous work indicated that the donor lineage of the exogenous genes has to be a,<sup>3</sup>  
<sup>4</sup>potentially unknown, Lachancea lineage [12–15]. These previous results, however,<sup>4</sup>  
<sup>5</sup>are based on species rather than genes trees ignoring the differential adaptation<sup>5</sup>  
<sup>6</sup>rate to their novel cellular environment between genes or due not consider lineages<sup>6</sup>  
<sup>7</sup>outside of the Lachancea clade. Considering the similarity in selection bias (Figure<sup>7</sup>  
<sup>8</sup>2b) and our calculation of selection on the exogenous genes (Figure 5b), both of<sup>8</sup>  
<sup>9</sup>which are free of any assumption about the origin of the exogenous genes, a species<sup>9</sup>  
<sup>10</sup>tree estimated from the exogenous genes may be biased towards the Lachancea<sup>10</sup>  
<sup>11</sup>clade. Estimating individual gene trees rather than relying on a species tree provided<sup>11</sup>  
<sup>12</sup>further evidence that the exogenous genes could originate from a lineage that does<sup>12</sup>  
<sup>13</sup>not belong to the Lachancea clade. As we highlighted in this study, relatively small<sup>13</sup>  
<sup>14</sup>sets of genes with a signal of a foreign cellular environment can significantly bias<sup>14</sup>  
<sup>15</sup>the outcome of a study. The same holds true for phylogenetic inferences [27], and as<sup>15</sup>  
<sup>16</sup>we showed the signal of the original endogenous cellular environment that shaped<sup>16</sup>  
<sup>17</sup>CUB is at different stages of decay in high and low expression genes (Figure S6).<sup>17</sup>  
<sup>18</sup>In summary, our work does not dispute an unknown Lachancea as possible origin,<sup>18</sup>  
<sup>19</sup>but provides an alternative hypothesis based on the codon usage of the exogenous<sup>19</sup>  
<sup>20</sup>genes, phylogenetic analysis, and synteny.<sup>20</sup>

<sup>21</sup> In terms of understanding the spread of the introgression, we calculated the ex-<sup>21</sup>  
<sup>22</sup>pected selective cost of codon mismatch between the *L. kluyveri* and *E. gossypii*<sup>22</sup>  
<sup>23</sup>lineages. Under our working hypothesis, the majority of the introgressed would have<sup>23</sup>  
<sup>24</sup>imposed a selective cost due to codon mismatch. Nevertheless,  $\sim 30\%$  of low expres-<sup>24</sup>  
<sup>25</sup>sion exogenous genes ( $\phi < 1$ ) appeared to be exapted at the time of the introgres-<sup>25</sup>  
<sup>26</sup>sion. This exaptation is due to the mutation bias in the endogenous genes matching<sup>26</sup>  
<sup>27</sup>the selection bias in the exogenous genes for GC ending codons. Our estimate of<sup>27</sup>  
<sup>28</sup>the selective cost of codon mismatch on the order of  $-0.0008$ . While this selective<sup>28</sup>  
<sup>29</sup>cost may not seem very large, assuming *L. kluyveri* had a large  $N_e$ , the fixation<sup>29</sup>  
<sup>30</sup>probability of the introgression is the astronomically small value of  $\approx 10^{-6952} \approx 0$ .<sup>30</sup>  
<sup>31</sup>Thus, the basic scenario of an introgression between two yeast species with large  $N_e$ <sup>31</sup>  
<sup>32</sup>and where the introgression solely imposes a selective cost due to codon mismatch<sup>32</sup>  
<sup>33</sup>is clearly too simplistic.<sup>33</sup>

For example, one or more loci with a combined selective advantage on the order of 0.0008 or greater would have made the introgression change from disadvantageous to effectively neutral or advantageous. While this scenario seems plausible, it raises the question as to why recombination events did not limit the introgression to only the adaptive loci. A potential answer is the low recombination rate between the endogenous and exogenous regions [14, 15]. This is presumably due to the dissimilarity in GC content and/or a lower than average sequence homology between the exogenous region and the one it replaced. A population bottleneck reducing the  $N_e$  of the *L. kluyveri* lineage around the time of the introgression could also help explain the spread of the introgression. Compatible with these explanation is the possibility of several advantageous loci distributed across the exogenous region drove a rapid selective sweep and/or the population through a bottleneck speciation process.

Assuming *E. gossypii* as potential source lineage of the exogenous region, we illustrated how information on codon usage can be used to infer the time since the introgression occurred using our estimates of mutation bias  $\Delta M$ . The  $\Delta M$  estimates are well suited for this task as they are free of the influence of selection and unbiased by  $N_e$  and other scaling terms, which is in contrast to our estimates of  $\Delta\eta$  [7]. Our estimated age of the introgression of  $6.2 \pm 1.2 \times 10^8$  generations is  $\sim 10$  times longer than a previous minimum estimate by [12] of  $5.6 \times 10^7$  generations, which was based on the effective population recombination rate and the population mutation parameter [28]. 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).

## Conclusion

Overall, 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

<sup>1</sup>used to study codon usage like CAI [16] or tAI [18], ROC SEMPPR incorporates<sup>1</sup>  
<sup>2</sup>the effects of mutation bias and amino acid composition explicitly [19]. We highlight<sup>2</sup>  
<sup>3</sup>potential issues when estimating codon preferences, as estimates can be biased by<sup>3</sup>  
<sup>4</sup>the signature of a second, historical cellular environment. In addition, we show<sup>4</sup>  
<sup>5</sup>how quantitative estimates of mutation bias and selection relative to drift can be<sup>5</sup>  
<sup>6</sup>obtained from codon data and used to infer the fitness cost of an introgression as<sup>6</sup>  
<sup>7</sup>well as its history and potential future. 7

8

8

## <sup>9</sup>**Materials and Methods** 9

### <sup>10</sup>Separating Endogenous and Exogenous Genes 10

<sup>11</sup>A GC-rich region was identified by [14] in the *L. kluyveri* genome extending from<sup>11</sup>  
<sup>12</sup>position 1 to 989,693 of chromosome C. This region was later identified as an<sup>12</sup>  
<sup>13</sup>introgression by [12]. We obtained the *L. kluyveri* genome from SGD Project<sup>13</sup>  
<sup>14</sup><http://www.yeastgenome.org/download-data/> (on 09-27-2014) and the annota-<sup>14</sup>  
<sup>15</sup>tion for *L. kluyveri* NRRL Y-12651 (assembly ASM14922v1) from NCBI (on 12-09-<sup>15</sup>  
<sup>16</sup>2014). We assigned 457 genes located on chromosome C with a location within the<sup>16</sup>  
<sup>17</sup> $\sim 1$  Mb window to the exogenous gene set. All other 4864 genes of the *L. kluyveri*<sup>17</sup>  
<sup>18</sup>genome were assigned to the exogenous genes. 18

19

19

### <sup>20</sup>Model Fitting with ROC SEMPPR 20

<sup>21</sup>ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) [21] and R (3.4.1)<sup>21</sup>  
<sup>22</sup>[29]. ROC SEMPPR was run from 10 different starting values for at least 250,000<sup>22</sup>  
<sup>23</sup>iterations and thinned to every 50th iteration. After manual inspection to verify that<sup>23</sup>  
<sup>24</sup>the MCMC had converged, parameter posterior means, log posterior probability and<sup>24</sup>  
<sup>25</sup>log likelihood were estimated from the last 500 samples (last 10% of samples). 25

26

26

### <sup>27</sup>Model selection 27

<sup>28</sup>The marginal likelihood of the combined and separated model fits was calculated<sup>28</sup>  
<sup>29</sup>using a generalized harmonic mean estimator [30]. A variance scaling of 1.1 was<sup>29</sup>  
<sup>30</sup>used to scale the important density of the estimator. Using the estimated marginal<sup>30</sup>  
<sup>31</sup>likelihoods, we calculated the Bayes factor to assess model performance. Increases<sup>31</sup>  
<sup>32</sup>in the variance scaling increase the estimated Bayes factor, therefore we report a<sup>32</sup>  
<sup>33</sup>conservative Bayes factor bases on a small variance scaling S7. 33

# <sup>1</sup>Comparing Codon Specific Parameter Estimates and Selecting Candidate lineages<sup>1</sup>

<sup>2</sup>As the choice of reference codon can reorganize codon families coding for an amino<sup>2</sup>  
<sup>3</sup>acid relative to each other, all parameter estimates were interpreted relative to the<sup>3</sup>  
<sup>4</sup>mean for each codon family.<sup>4</sup>

$$\Delta M_i = \Delta M_{i,1} - \overline{\Delta M_i} \quad (1)^6$$

$$\Delta \eta_i = \Delta \eta_{i,1} - \overline{\Delta \eta_i} \quad (2)^8$$

<sup>9</sup>Comparison of codon specific parameters ( $\Delta M$  and  $\Delta \eta = 2N_e q(\eta_i - \eta_j)$ ) was per-<sup>10</sup>  
<sup>11</sup>formed using the function `lmodel2` in the R package `lmodel2` (1.7.3) [31] and R<sup>11</sup>  
<sup>12</sup>version 3.4.1 [29]. The parameter  $\Delta \eta$  can be interpreted as the difference in fitness<sup>12</sup>  
<sup>13</sup>between codon  $i$  and  $j$  for the average gene with  $\phi = 1$  scaled by the effective pop-<sup>13</sup>  
<sup>14</sup>ulation size  $N_e$ , and the selective cost of an ATP  $q$  [4, 7]. Type II regression was<sup>14</sup>  
<sup>15</sup>performed with re-centered parameter estimates, accounting for noise in dependent<sup>15</sup>  
<sup>16</sup>and independent variable [23].<sup>16</sup>

<sup>17</sup>Due to the greater dissimilarity of the  $\Delta M$  estimates between the endogenous and<sup>17</sup>  
<sup>18</sup>exogenous genes, and the slower decay rate of mutation bias, we decided to focus<sup>18</sup>  
<sup>19</sup>on our estimates of mutation bias to identify potential source lineages. The top ten<sup>19</sup>  
<sup>20</sup>lineages with the highest similarity in  $\Delta M$  to the exogenous genes were selected as<sup>20</sup>  
<sup>21</sup>potential candidates (Figure 2).<sup>21</sup>

## <sup>22</sup>Phylogenetic Analysis<sup>22</sup>

<sup>23</sup>Using the dataset from [20], we first identified 121 alignments for exogenous genes<sup>23</sup>  
<sup>24</sup>and further contained homologous genes for *E. gossypii*, and *L. thermotolerance*.<sup>24</sup>  
<sup>25</sup>We excluded all species from the alignments that do not belong to the Saccharomyc-<sup>25</sup>  
<sup>26</sup>etaceae clade. IQTree [24] was used to identify the best fitting model for each gene<sup>26</sup>  
<sup>27</sup>and to estimate the individual gene trees. The distance between *L. kluyveri*, *E.*<sup>27</sup>  
<sup>28</sup>*gossypii*, and *L. thermotolerance* was calculated for each tree to identify genes for<sup>28</sup>  
<sup>29</sup>which exogenous genes are more closely related to *E. gossypii* or *L. thermotolerance*.<sup>29</sup>  
<sup>30</sup>

## <sup>31</sup>Synteny Comparison<sup>31</sup>

<sup>32</sup>We obtained complete genome sequences for all 10 candidate lineages (Table 2)<sup>32</sup>  
<sup>33</sup>from NCBI (on: 02-05-2017). Genomes were aligned and checked for synteny using<sup>33</sup>

<sup>1</sup>SyMAP (4.2) with default settings [32, 33]. We assess synteny as percentage coverage<sup>1</sup>  
<sup>2</sup>of the exogenous gene region.<sup>2</sup>

<sup>3</sup><sup>3</sup>

<sup>4</sup>Estimating Age of Introgression<sup>4</sup>

<sup>5</sup>We modeled the change in codon frequency over time using an exponential model<sup>5</sup>  
<sup>6</sup>for all two codon amino acids, and describing the change in codon  $c_1$  as<sup>6</sup>

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

<sup>9</sup>where  $\mu_{i,j}$  is the rate at which codon  $i$  mutates to codon  $j$  and  $c_1$  is the fre-<sup>9</sup>  
<sup>10</sup>quency of the reference codon. Initial codon frequencies  $c_1(0)$  for each codon fam-<sup>10</sup>  
<sup>11</sup>ily where taken from our mutation parameter estimates for *E. gossypii* where<sup>11</sup>  
<sup>12</sup> $c_1(0) = \exp[\Delta M_{\text{gos}}]/(1 + \exp[\Delta M_{\text{gos}}])$ . Our estimates of  $\Delta M_{\text{endo}}$  can be used to<sup>12</sup>  
<sup>13</sup>calculate the steady state of equation 3 were  $\frac{dc_1}{dt} = 0$  to obtain the equality<sup>13</sup>

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

<sup>17</sup>Solving for  $\mu_{1,2}$  gives us  $\mu_{1,2} = \Delta M_{\text{endo}} \exp[\mu_{2,1}]$  which allows us to rewrite and<sup>17</sup>  
<sup>18</sup>solve equation 3 as<sup>18</sup>

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

<sup>21</sup>where  $X = (1 + \Delta M_{\text{endo}})\mu_{2,1}t$  and  $K = c_1(0)(1 + \Delta M_{\text{endo}})$ .<sup>21</sup>

<sup>22</sup>Equation 5 was solved with a mutation rate  $\mu_{2,1}$  of  $3.8 \times 10^{-10}$  per nucleotide per<sup>22</sup>  
<sup>23</sup>generation [34]. Current codon frequencies for each codon family where taken from<sup>23</sup>  
<sup>24</sup>our estimates of  $\Delta M$  from the exogenous genes. Mathematica (11.3) [35] was used<sup>24</sup>  
<sup>25</sup>to calculate the time  $t_{\text{intro}}$  it takes for the initial codon frequencies  $c_1(0)$  for each<sup>25</sup>  
<sup>26</sup>codon family to equal the current exogenous codon frequencies. The same equation<sup>26</sup>  
<sup>27</sup>was used to determine the time  $t_{\text{decay}}$  at which the signal of the exogenous cellular<sup>27</sup>  
<sup>28</sup>environment has decayed to within 1% of the endogenous environment.<sup>28</sup>

<sup>30</sup>Estimating Selection against Codon Mismatch<sup>30</sup>

<sup>31</sup>In order to estimate the selection against codon mismatch, we had to make three<sup>31</sup>  
<sup>32</sup>key assumptions. First, we assumed that the current exogenous amino acid sequence<sup>32</sup>  
<sup>33</sup>of a gene is representative of its ancestral state and the replaced endogenous gene<sup>33</sup>



it replaced. Second, we assume that the currently observed cellular environment of *E. gossypii* reflects the cellular environment that the exogenous genes experienced before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy of selection between the cellular environments due to differences in either effective population size  $N_e$  or the selective cost of an ATP  $q$  of the source lineage and *L. kluyveri* can be expressed as a scaling constant and that protein synthesis rate  $\phi$  has not changed between the replaced endogenous and the introgressed exogenous genes. Using estimates for  $N_e = 1.36 \times 10^7$  [25] for *Saccharomyces paradoxus* we scale our estimates of  $\Delta\eta$  which explicitly contains the effective population size  $N_e$  [7] and define  $\Delta\eta' = \frac{\Delta\eta}{N_e}$ .

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

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

where  $s_g$  is the overall strength of selection for translational efficiency on gene,  $g$  in the exogenous gene set,  $\kappa$  is a constant, scaling the efficacy of selection between the endogenous and exogenous cellular environments,  $L_g$  is length of the protein in codons,  $\phi_g$  is the estimated protein synthesis rate of the gene in the endogenous environment, and  $\Delta\eta'_i$  is the  $\Delta\eta'$  for the codon at position  $i$ . As stated previously, our  $\Delta\eta$  are relative to the mean of the codon family. We find that the selection against the introgressed genes is minimized at  $\kappa \sim 5$  (Figure S5b). Thus, we expect a five fold difference in the efficacy of selection between *L. kluyveri* and *E. gossypii*, due to differences in either protein synthesis rate  $\phi$ , effective population size  $N_e$ , and/or the selective cost of an ATP  $q$ . Therefore, we set  $\kappa = 1$  if we calculate the  $s_g$

<sup>1</sup>for the endogenous and the current exogenous genes, and  $\kappa = 5$  for  $s_g$  for selection<sup>1</sup>  
<sup>2</sup>calculations at the time of introgression.<sup>2</sup>

<sup>3</sup> However, since we are unable to observe codon sequences of the replaced en-<sup>3</sup>  
<sup>4</sup>dogamous genes and for the exogenous genes at the time of introgression, instead<sup>4</sup>  
<sup>5</sup>of summing over the sequence, we calculate the expected codon count  $E[n_{g,i}]$  for<sup>5</sup>  
<sup>6</sup>codon  $i$  in gene  $g$  simply as the probability of observing codon  $i$  multiplied by the<sup>6</sup>  
<sup>7</sup>number of times the corresponding amino acids is observed in gene  $g$ , yielding:<sup>7</sup>

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

<sup>12</sup>where  $m_{a_i}$  is the number of occurrences of amino acid  $a$  that codon  $i$  codes for. Thus,<sup>12</sup>  
<sup>13</sup>replacing the summation over the sequence length  $L_g$  in equ. (6) by a summation<sup>13</sup>  
<sup>14</sup>over the codon set  $C$  and calculating  $s_g$  as<sup>14</sup>

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

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

## <sup>21</sup>Acknowledgments<sup>21</sup>

<sup>22</sup>The authors would like to thank Alexander Cope for helpful criticisms and suggestions for this work.<sup>22</sup>

## <sup>23</sup>Availability of data and materials<sup>23</sup>

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

## <sup>25</sup>Authors' contributions<sup>25</sup>

<sup>26</sup>CL and MAG have conceptualized the study and written the manuscript. CL performed the experiments and<sup>25</sup>  
<sup>26</sup>analyzed the data. BCO and RZ helped design experiments and provided expertise for analyzing the data. All<sup>26</sup>  
<sup>27</sup>Authors approved the final manuscript.<sup>27</sup>

## <sup>28</sup>Funding<sup>28</sup>

<sup>29</sup>This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ) and DEB-1355033 (BCO, MAG,<sup>28</sup>  
<sup>29</sup>and RZ) with additional support from The University of Tennessee Knoxville. CL received support as a Graduate<sup>29</sup>  
<sup>30</sup>Student Fellow at the National Institute for Mathematical and Biological Synthesis, an Institute sponsored by the<sup>30</sup>  
<sup>30</sup>National Science Foundation through NSF Award DBI-1300426, with additional support from UTK.<sup>30</sup>

## <sup>31</sup>Ethics approval and consent to participate<sup>31</sup>

<sup>32</sup>Not applicable<sup>32</sup>

## <sup>33</sup>Consent for publication<sup>33</sup>

<sup>33</sup>Not applicable<sup>33</sup>

# <sup>1</sup>Competing interests

<sup>2</sup>The authors declare that they have no competing interests.

## <sup>3</sup>Author details

<sup>1</sup>Department of Ecology & Evolutionary Biology, University of Tennessee, 37996, Knoxville, TN, USA. <sup>2</sup>National  
<sup>4</sup>Institute for Mathematical and Biological Synthesis, 37996, Knoxville, TN, USA. <sup>3</sup>Max-Planck Institute of  
<sup>5</sup>Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, 01307, Dresden, Germany. <sup>4</sup>Department of Business  
<sup>6</sup>Analytics and Statistics, University of Tennessee, 37996, Knoxville, TN, USA.

## <sup>7</sup>References

- <sup>1</sup> Gouy, M., Gautier, C.: Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Research* **10**,  
<sup>8</sup> 7055–7074 (1982) 8
- <sup>9</sup> 2. Ikemura, T.: Codon usage and tRNA content in unicellular and multicellular organisms. *Molecular Biology and*  
<sup>10</sup> *Evolution* **2**, 13–34 (1985) 9
- <sup>11</sup> 3. Bulmer, M.: The selection-mutation-drift theory of synonymous codon usage. *Genetics* **129**, 897–907 (1990) 10
- <sup>12</sup> 4. Gilchrist, M.A.: Combining models of protein translation and population genetics to predict protein production  
<sup>13</sup> rates from codon usage patterns. *Molecular Biology and Evolution* **24**(11), 2362–2372 (2007) 11
- <sup>14</sup> 5. Shah, P., Gilchrist, M.A.: Explaining complex codon usage patterns with selection for translational efficiency,  
<sup>15</sup> mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences U.S.A* **108**(25),  
<sup>16</sup> 10231–10236 (2011) 13
- <sup>17</sup> 6. Wallace, E.W., Airoidi, E.M., Drummond, D.A.: Estimating selection on synonymous codon usage from noisy  
<sup>18</sup> experimental data. *Molecular Biology and Evolution* **30**, 1438–1453 (2013) 14
- <sup>19</sup> 7. Gilchrist, M.A., Chen, W.C., Shah, P., Landerer, C.L., Zaretzki, R.: Estimating gene expression and  
<sup>20</sup> codon-specific translational efficiencies, mutation biases, and selection coefficients from genomic data alone.  
<sup>21</sup> *Genome Biology and Evolution* **7**, 1559–1579 (2015) 16
- <sup>22</sup> 8. Médigue, C., Rouxel, T., Vigier, P., Hénaut, A., Danchin, A.: Evidence for horizontal gene transfer in  
<sup>23</sup> *Escherichia coli* speciation. *Journal of Molecular Biology* **222**(4), 851–856 (1991) 17
- <sup>24</sup> 9. Lawrence, J.G., Ochman, H.: Amelioration of bacterial genomes: Rates of change and exchange. *Journal of*  
<sup>25</sup> *Molecular Biology* **44**, 383–397 (1997) 18
- <sup>26</sup> 10. Marcet-Houben, M., Gabaldón, T.: Beyond the whole-genome duplication: Phylogenetic evidence for an ancient  
<sup>27</sup> interspecies hybridization in the baker's yeast lineage. *PLoS Biology* **13**(8), 1002220 (2015) 19
- <sup>28</sup> 11. Beimforde, C., Feldberg, K., Nylinder, S., Rikkinen, J., Tuovila, H., Dörfelt, H., Gube, M., Jackson, D.J.,  
<sup>29</sup> Reitner, J., Seyfullah, L.J., Schmidt, A.R.: Estimating the phanerozoic history of the ascomycota lineages:  
<sup>30</sup> combining fossil and molecular data. *Mol. Phylogenet. Evol.* **78**, 386–398 (2014) 21
- <sup>31</sup> 12. Friedrich, A., Reiser, C., Fischer, G., Schacherer, J.: Population genomics reveals chromosome-scale  
<sup>32</sup> heterogeneous evolution in a protoploid yeast. *Molecular Biology and Evolution* **32**(1), 184–192 (2015) 22
- <sup>33</sup> 13. Vakirlis, N., Sarilar, V., Drillon, G., Fleiss, A., Agier, N., Meyniel, J.-P., Blanpain, L., Carbone, A., Devillers, H.,  
<sup>34</sup> Dubois, K., Gillet-Markowska, A., Graziani, S., Huu-Vang, N., Poiriel, M., Reisser, C., Schott, J., Schacherer,  
<sup>35</sup> J., Lafontaine, I., Llorente, B., Neuvéglise, C., Fischer, G.: Reconstruction of ancestral chromosome  
<sup>36</sup> architecture and gene repertoire reveals principles of genome evolution in a model yeast genus. *Genome*  
<sup>37</sup> *research* **26**(7), 918–32 (2016) 23
- <sup>38</sup> 14. Payen, C., Fischer, G., Marck, C., Proux, C., Sherman, D.J., Coppée, J.-Y., Johnston, M., Dujon, B.,  
<sup>39</sup> Neuvéglise, C.: Unusual composition of a yeast chromosome arm is associated with its delayed replication.  
<sup>40</sup> *Genome Research* **19**(10), 1710–1721 (2009) 25
- <sup>41</sup> 15. Brion, C., Legrand, S., Peter, J., Caradec, C., Pflieger, D., Hou, J., Friedrich, A., Llorente, B., Schacherer, J.:  
<sup>42</sup> Variation of the meiotic recombination landscape and properties over a broad evolutionary distance in yeasts.  
<sup>43</sup> *PLoS Genetics* **13**(8), 1006917 (2017) 26
- <sup>44</sup> 16. Sharp, P.M., Li, W.H.: The codon adaptation index - a measure of directional synonymous codon usage bias,  
<sup>45</sup> and its potential applications. *Nucleic Acids Research* **15**, 1281–1295 (1987) 27
- <sup>46</sup> 17. Wright, F.: The 'effective number of codons' used in a gene. *Genet* **87**, 23–29 (1990) 28
- <sup>47</sup> 18. dos Reis, M., Savva, R., Wernisch, L.: Solving the riddle of codon usage preferences: a test for translational  
<sup>48</sup> selection. *Nucleic Acids Research* **32**(17), 5036–5044 (2004) 29

19. Cope, A.L., Hettich, R.L., Gilchrist, M.A.: Quantifying codon usage in signal peptides: Gene expression and amino acid usage explain apparent selection for inefficient codons. *Biochimica et Biophysica Acta (BBA) - Biomembranes* **1860**(12), 2479–2485 (2018)
20. Shen, X.X., Opulente, D.A., Kominek, J., Zhou, X., Steenwyk, J.L., Buh, K.V., Haase, M.A.B., Wisecaver, J.H., Wang, M., Doering, D.T., Boudouris, J.T., Schneider, R.M., Langdon, Q.K., Ohkuma, M., Endoh, R., Takashima, M., Manabe, R., Čadež, N., Libkind, D., Rosa, C., DeVirgilio, J., Hulfachor, A.B., Groenewald, M., Kurtzman, C., Hittinger, C.T., Rokas, A.: Tempo and mode of genome evolution in the budding yeast subphylum. *Cell* **175**(6), 1533–1545 (2018)
21. Landerer, C., Cope, A., Zaretski, R., Gilchrist, M.A.: AnaCoDa: analyzing codon data with bayesian mixture models. *Bioinformatics* **34**(14), 2496–2498 (2018)
22. Tsankov, A.M., Thompson, D.A., Socha, A., Regev, A., Rando, O.J.: The role of nucleosome positioning in the evolution of gene regulation. *PLoS Biol* **8**(7), 1000414 (2010)
23. Sokal, R.R., Rohlf, F.J.: *Biometry - The principles and practice of statistics in biological*, pp. 547–555. W. H. Freeman, ??? (1981)
24. Nguyen, L.T., Schmidt, H.A., von Haeseler, A., Minh, B.Q.: Iq-tree: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **32**(1), 268–274 (2015)
25. Wagner, A.: Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* **22**, 1365–1374 (2005)
26. Sella, G., Hirsh, A.E.: The application of statistical physics to evolutionary biology. *Proceedings of the National Academy of Sciences of the United States of America* **102**, 9541–9546 (2005)
27. Salichos, L., Rokas, A.: Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature* **497**, 327–331 (2013)
28. 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)
29. R Core Team: *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria (2013). <http://www.R-project.org/>
30. Gronau, Q.F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D.S., Forster, J.J., Wagenmakers, E.J., Steingrover, H.: A tutorial on bridge sampling. *Journal of Mathematical Psychology* **81**, 80–97 (2017)
31. Legendre, P.: *Lmodel2: Model II Regression*. (2018). R package version 1.7-3. <https://CRAN.R-project.org/package=lmodel2>
32. 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)
33. Soderlund, C., Bomhoff, M., Nelson, W.: Symap v3.4: a turnkey syntenic system with application to plant genomes. *Nucleic Acids Research* **39**(10), 68 (2011)
34. Lang, G.I., Murray, A.W.: Estimating the per-base-pair mutation rate in the yeast *Saccharomyces cerevisiae*. *Genetics* **178**(1), 67–82 (2008)
35. Wolfram Research Inc.: *Mathematica 11*. (2017). <http://www.wolfram.com>

# Supplementary Material

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

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

Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	Combined
Ala A	GCG	GCA	GCG	GCG
Cys C	TGC	TGT	TGC	TGC
Asp D	GAC	GAT	GAC	GAC
Glu E	GAG	GAA	GAG	GAG
Phe F	TTC	TTT	TTT	TTT
Gly G	GGC	GGT	GGC	GGC
His H	CAC	CAT	CAC	CAC
Ile I	ATC	ATT	ATC	ATA
Lys K	AAG	AAA	AAG	AAA
Leu L	CTG	TTG	CTG	CTG
Asn N	AAC	AAT	AAC	AAT
Pro P	CCG	CCA	CCG	CCG
Gln Q	CAG	CAA	CAG	CAG
Arg R	CGC	AGA	AGG	CGG
Ser <sub>4</sub> S	TCG	TCT	TCG	TCG
Thr T	ACG	ACA	ACG	ACG
Val V	GTG	GTT	GTG	GTG
Tyr Y	TAC	TAT	TAC	TAC
Ser <sub>2</sub> Z	AGC	AGT	AGC	AGC

1		1
2		2
3		3
4		4
5		5
6		6
7		7
8		8
9		9
10	<b>Table S2</b> Synonymous selection codon preference based on our estimates of $\Delta\eta$ . Shown are the most likely codon in high expression genes for each amino acid in: <i>E. gossypii</i> , in the endogenous and exogenous genes of <i>L. kluyveri</i> , and in the combined <i>L. kluyveri</i> genome without accounting for the two cellular environments.	
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

Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	Combined
Ala A	GCT	GCT	GCT	GCT
Cys C	TGT	TGT	TGT	TGT
Asp D	GAT	GAC	GAT	GAT
Glu E	GAA	GAA	GAA	GAA
Phe F	TTT	TTC	TTC	TTC
Gly G	GGA	GGT	GGT	GGT
His H	CAT	CAC	CAT	CAT
Ile I	ATA	ATC	ATT	ATT
Lys K	AAA	AAG	AAA	AAG
Leu L	TTA	TTG	TTG	TTG
Asn N	AAT	AAC	AAT	AAC
Pro P	CCA	CCA	CCT	CCA
Gln Q	CAA	CAA	CAA	CAA
Arg R	AGA	AGA	AGA	AGA
Ser <sub>4</sub> S	TCA	TCC	TCT	TCT
Thr T	ACT	ACC	ACT	ACT
Val V	GTT	GTC	GTT	GTT
Tyr Y	TAT	TAC	TAT	TAC
Ser <sub>2</sub> Z	AGT	AGT	AGT	AGT

