

Copy Correction and Concerted Evolution in the Conservation of Yeast Genes

Saumyadipta Pyne,* Steven Skiena* and Bruce Futcher^{†,1}

*Department of Computer Science and [†]Department of Microbiology and Molecular Genetics,
Stony Brook University, Stony Brook, New York 11794

Manuscript received August 4, 2004
Accepted for publication April 6, 2005

ABSTRACT

The yeast *Saccharomyces cerevisiae* and other members of the genus *Saccharomyces* are descendants of an ancient whole-genome duplication event. Although most of the duplicate genes have since been deleted, many remain, and so there are many pairs of related genes. We have found that poorly expressed genes diverge rapidly from their paralog, while highly expressed genes diverge little, if at all. This lack of divergence of highly expressed paralogous gene pairs seems to involve gene correction: one member of the pair “corrects” the sequence of its twin, and so the gene pair evolves as a unit. This correction presumably involves gene conversion and could occur via a reverse-transcribed cDNA intermediate. Such correction events may also occur in other organisms. These results support the idea that copies of poorly expressed genes are preserved when they diverge to take on new functions, while copies of highly expressed genes are preserved when they are needed to provide additional gene product for the original function.

IT is generally believed that selection for preferred codons (codon bias) increases the sequence conservation of highly expressed genes relative to poorly expressed genes (the “selection hypothesis”) (POWELL and MORIYAMA 1997). Both highly expressed and poorly expressed genes are selected for function, which means that many nonsynonymous codon changes are selected against; but in addition, for highly expressed genes, many synonymous changes are also selected against to maintain codons preferred for translational efficiency and accuracy. A corollary of this argument is that when an organism has two similar copies of a highly expressed gene, these copies should be preserved in evolution as a gene pair sharing high homology, because selection for both function and codon bias prevents the members of the pair from drifting apart.

In this study, we propose a parallel hypothesis for the conservation of duplicated highly expressed genes and show that the new hypothesis not only plays a significant role, at least in yeast, but also may be more important than selection under certain conditions.

We call our hypothesis the *correction hypothesis*. It consists of three proposals: first, that one copy of a gene can correct the sequence of a second copy; second, that correction depends on high sequence identity; and third, that the probability of correction depends on the level of gene expression. We propose two possible mechanisms of correction. In the first, correction happens through the occasional copying back of mature

mRNA into cDNA using reverse transcriptase and a subsequent recombinational interaction between the cDNA and the second copy of the chromosomal gene. In the second mechanism, correction is due to a direct recombinational interaction between the two genes of a duplicate pair.

In the following sections, we first show that correction indeed plays a significant role in the conservation of gene pairs in yeast and that correction is correlated with the level of gene expression. Evidence for correction is, first, that the conservation between the members of a highly expressed gene pair is too high to be explained by selection alone, and second, that the pattern of nucleotide substitution within and between species of *Saccharomyces* is much more compatible with the correction hypothesis than with the selection hypothesis. We examine some properties of correction and consider whether correction might play a role in other organisms.

MATERIALS AND METHODS

Saccharomyces sequences were obtained from the *Saccharomyces* Genome Database (<http://www.yeastgenome.org/>), from Washington University (<http://www.genome.wustl.edu/projects/yeast/> and http://www.genome.wustl.edu/blast/yeast_client.cgi), and from the Massachusetts Institute of Technology (<http://www-genome.wi.mit.edu/seq/Saccharomyces/>). ClustalW-based end-to-end fungal alignments of the *Saccharomyces cerevisiae* genes and their analogous sequences in the other *Saccharomyces* species were obtained from the *Saccharomyces* Genome Database (SGD; <http://www.yeastgenome.org/>) whenever available.

Genes and gene pairs from the ancient duplication were selected using the “blocks menu” page of the website of Wolfe and colleagues at <http://acer.gen.tcd.ie/cgi-bin/khwolfe/>

¹Corresponding author: Department of Microbiology and Molecular Genetics, Stony Brook University, Stony Brook, NY 11794-5222.
E-mail: bfutcher@ms.cc.sunysb.edu

blocks.pl?block=ALL. All gene pairs from the above site whose coding sequences appear at SGD under the same systematic/standard name were downloaded in an automated manner. Any gene, and thus its pair, with a name that led to confusion in its recognition, was discarded. Note especially that we have not analyzed all duplicated genes in *S. cerevisiae*. We have restricted our study to those genes thought to have been duplicated as part of a single, ancient, genome-wide duplication event. For example, many duplications near telomeres appear to be recent duplications. We have not included any of these recently duplicated, telomere-associated genes in our analysis.

We used a modified version of the Jukes-Cantor model for measuring divergence between a pair of gene sequences. The conventional Jukes-Cantor model for computing evolutionary divergence between two sequences is extended to account for "internal" indels in the pairwise alignment of sequences. (An internal indel is not a part of the continuous batch of indels that might be present at either extremity of a pairwise alignment, perhaps owing to the difference in the lengths of the two sequences.) The modified measure treats an indel as a substitution of weight c between 0 and 1, inclusive, whereas a substitution of any kind is, as in the conventional model, of weight 1. Upon making the usual approximations, the divergence is given by $-d = -((3+c)/(4+c)) \times \ln(1 - (4+c) \times p/(3+c))$, where p is the proportion of substitutions, which, in our model, is a 1: c weighted proportion of both substitutions and internal indels. Clearly, a weight of $c = 0$ reduces the new measure to the old one. Inheriting the property of the old model, the new measure is also a partial function; *i.e.*, it is not defined for certain values of the valid input p . The divergence-expression figures are plotted with the above parameter c set to 1. Varying the value of c between 0 and 1 makes little difference to the plots.

For computing the counts and ratio of synonymous and nonsynonymous substitutions, we use the Synonymous Nonsynonymous Analysis Program (SNAP) available at the HIV Sequence Database (hiv-web.lanl.gov) (NEI and GOJOBORI 1986). Complete data are available at: <http://www.cs.sunysb.edu/~compbio/Correction/>.

RESULTS

High expression, conservation, and correction: The genus *Saccharomyces* arose from an ancient whole-genome duplication event, shortly after *Saccharomyces* diverged from *Kluyveromyces* (WOLFE and SHIELDS 1997; KELLIS *et al.* 2004). Subsequent to the duplication, many individual deletion events deleted one of the copies of most of the duplicate genes. Nevertheless, the *S. cerevisiae* of today has up to 450 gene pairs (16% of the proteome) remaining from the ancient duplication (SEOIGHE and WOLFE 1999). While studying a 382-pair subset of these duplicates, we found a remarkably strong negative correlation between sequence divergence and the codon adaptation index (CAI; Figure 1). CAI (SHARP and LI 1987) is used as a surrogate for gene expression (FUTCHER *et al.* 1999). That is, the highly expressed genes have diverged less from their duplicates than the poorly expressed genes. The Pearson correlation is $r = -0.72$ ($P < 10^{-16}$ for the null hypothesis that $r = 0$). The majority of gene pairs have low codon bias (*i.e.*, low expression) and high DNA sequence divergence, while a substantial subgroup has high codon bias

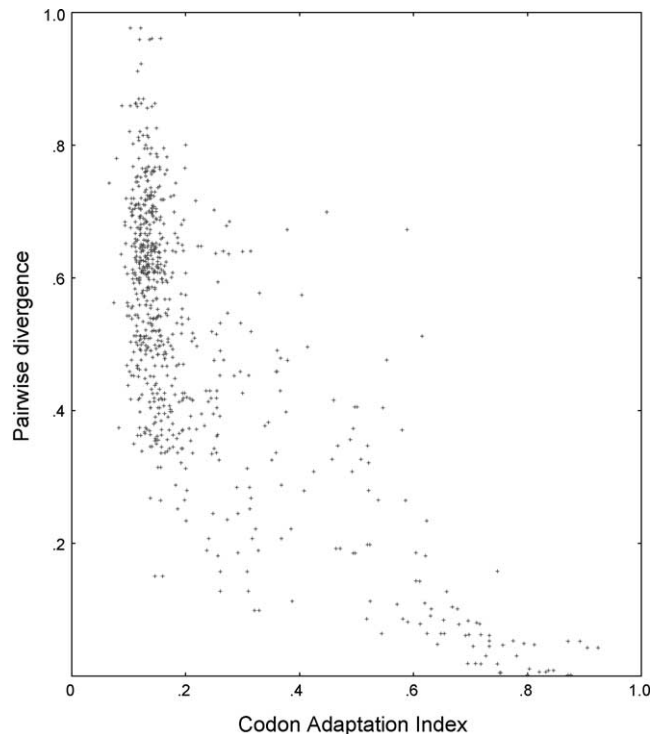


FIGURE 1.—Divergence as a function of expression in *S. cerevisiae*. For each of 764 genes (the members of qualified gene pairs from the ancient duplication; see MATERIALS AND METHODS), the closest homolog is found. The modified Jukes-Cantor divergence score for the pair, if defined (MATERIALS AND METHODS), is plotted on the y -axis, on the basis of the ClustalW pairwise alignment of the genes. The CAI of the chosen gene is plotted on the x -axis. CAI is a surrogate for the level of mRNA produced. Note that the divergence of each gene pair (g_1, g_2) is plotted twice, once against the CAI of g_1 and once against the CAI of g_2 . There is a strong negative correlation between divergence and CAI (*i.e.*, level of expression), with correlation coefficient $r = -0.72$ ($P < 10^{-16}$ for the null hypothesis that $r = 0$).

(*i.e.*, high expression) and low DNA sequence divergence.

The correction model: The correlation seen in Figure 1 is not necessarily inconsistent with the "selection" model. However, the very strong correlation and the very large differences in sequence identity were so striking that we wondered whether there might be some other explanation. In particular, some highly expressed gene pairs were >95% identical in DNA sequence despite apparently diverging many millions of years ago (WOLFE and SHIELDS 1997). This striking conservation might be explained if DNA sequence correction occurred between members of the pair. *S. cerevisiae* has a very active homologous recombination system, and recombination or gene conversion could account for pairs with very high identity.

In the most obvious model of sequence correction (Figure 2, "DNA-DNA Correction"), the two gene copies interact, and gene conversion events occur directly between the two chromosomal copies. There is some evi-

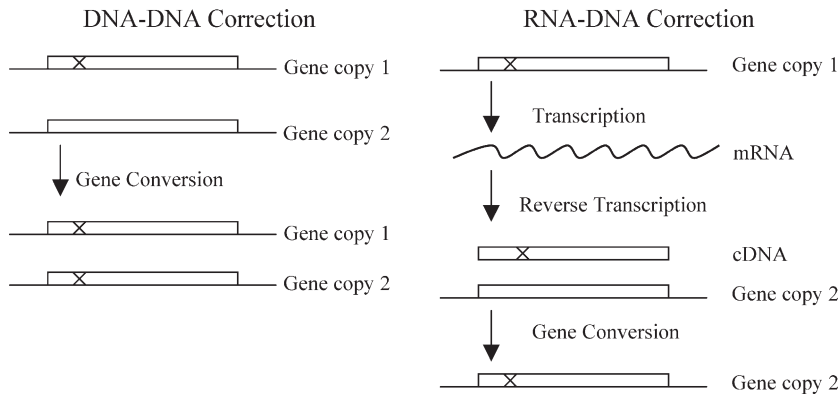


FIGURE 2.—The correction model. In DNA-DNA Correction (left), one gene interacts with a second gene and corrects it by gene conversion. In RNA-DNA Correction (right), one gene is transcribed and then copied into cDNA (or a cDNA/RNA hybrid), and this cDNA molecule interacts with a second gene and corrects it by gene conversion.

dence that very highly transcribed genes are particularly active in recombinational events, while repressed genes are relatively inactive (SAXE *et al.* 2000).

A second model of sequence correction (Figure 2, "RNA-DNA Correction") invokes an RNA intermediate. This model depends upon the fact that *S. cerevisiae* contains a retrotransposon, Ty, which encodes a reverse transcriptase. Occasionally, this reverse transcriptase makes cDNA copies of normal, cellular genes (XU and BOEKE 1990), and these cDNA copies can recombine with their chromosomal homologs (DERR *et al.* 1991; DERR and STRATHERN 1993). It has been proposed that reverse transcription followed by homologous recombination explains why so few genes in *S. cerevisiae* have introns (BALTIMORE 1985; FINK 1987). In these proposals, a gene with an intron produces a transcript; the transcript is spliced; the spliced transcript is converted to cDNA by Ty reverse transcriptase; and then the cDNA interacts with the chromosomal gene and removes the intron by gene conversion. Similarly, in our model of correction via an RNA intermediate (Figure 2), an mRNA produced by gene copy 1 is converted to cDNA by Ty reverse transcriptase. This cDNA then interacts with gene copy 2, and by gene conversion corrects gene copy 2 into an exact duplicate of gene copy 1. Over succeeding generations, gene copies 1 and 2 may again drift apart, but then will undergo another round of correction, again making the two genes identical. An attractive feature of this model is that the probability of correction is obviously directly proportional to the level of gene expression: the more mRNA that is made, the higher the probability that some of it will be converted to cDNA. This could explain the strong correlation between pairwise homology and expression level. However, a weakness of this model is that RNA-mediated gene conversion is much rarer than DNA-DNA events (DERR and STRATHERN 1993).

A feature of both of these correction models is that gene conversion will occur only between sequences that have a very high percentage of identity; even a small number of mismatches drastically reduces the frequency of gene conversion (MODRICH and LAHUE 1996; DATTA *et al.* 1997; CHEN and JINKS-ROBERTSON 1998). Once

members of a gene pair have drifted sufficiently far apart, they would no longer be able to correct each other.

In either model, genes expressed at a high level correct each other frequently because highly transcribed genes are recombinationally active (SAXE *et al.* 2000), as in the DNA-DNA model, or because highly transcribed genes make more RNA, as in the RNA-DNA model, and so do not drift apart. Because they do not drift apart, they remain eligible for future correction events. Genes expressed at a low level correct each other infrequently, and so sometimes drift far apart between correction events, greatly reducing the probability of future correction. Thus, there would be two groups of gene pairs—a highly expressed, highly conserved group and a poorly expressed, poorly conserved group.

Correction vs. selection: We wished to distinguish the selection model from the correction models and took advantage of the fact that five other species of *Saccharomyces* (*castelli*, *kluyveri*, *mikatae*, *paradoxus*, and *byanus*) have recently been sequenced (CLIFTEN *et al.* 2003; KELLIS *et al.* 2003). These species diverged from *S. cerevisiae* at various times (*paradoxus*, 10 MYA; *byanus*, 20 MYA; *mikatae*, 20 MYA; *castelli*, 50 MYA; *kluyveri*, very roughly 75 MYA), whereas the duplication of the *Saccharomyces* genome preceded most of these speciation events (the probable exception being the speciation of *S. kluyveri*). We reasoned as follows: if the sequence of a gene is maintained solely by selection, and not by correction, then the rate at which the similarity of gene 1 copy 1 and gene 1 copy 2 drift apart within *S. cerevisiae* will be roughly the same as that of gene 1 copy 1 of *S. cerevisiae* drifting from their orthologs in each of the other species. In other words, if there is no correction, then each gene in each species will diverge independently and at roughly the same rate, regardless of whether there is a duplicate gene in the same cell. If anything, divergence will be faster when there is a duplicate copy in the same cell, since the duplicate can provide important functions lost by its mutating partner. Alternatively, if correction is a significant force, then gene 1 copy 1 and gene 1 copy 2 in *S. cerevisiae* will drift apart more slowly (if at all) than gene 1 copy 1 of *S. cerevisiae* and its nearest

ortholog in each of the other species, since of course there will not be any correction between species. (Note that for purposes of this argument, it does not matter whether the genes are unique or duplicate pairs in the other species.) Thus, for each *S. cerevisiae* gene in a list of 382 *S. cerevisiae* gene pairs, we found the closest ortholog in each of the other five species and compared the divergence of these five orthologs and of the *S.*

cerevisiae paralog. Divergence was examined as a correlate of the codon adaptation index (SHARP and LI 1987). Results are shown in Figure 3.

There are several noteworthy points. First, for genes expressed at low and medium levels, the divergence is roughly proportional to the time since divergence. This is true both between species and within *S. cerevisiae*. Second, for genes expressed at high levels, the divergence is decreased; *i.e.*, highly expressed genes tend to be more conserved. This conservation is consistent with the idea that selection is important in preserving highly expressed genes; presumably some of the effect is due to selection for preferred codons. Nevertheless, the divergence between the highly expressed *S. cerevisiae* genes and their *castelli* or *kluyveri* homologs is still consider-

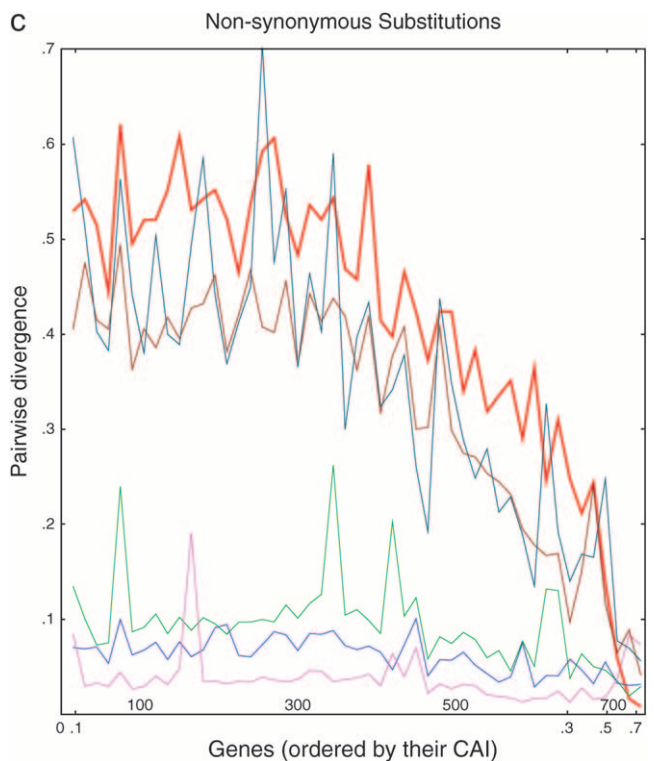
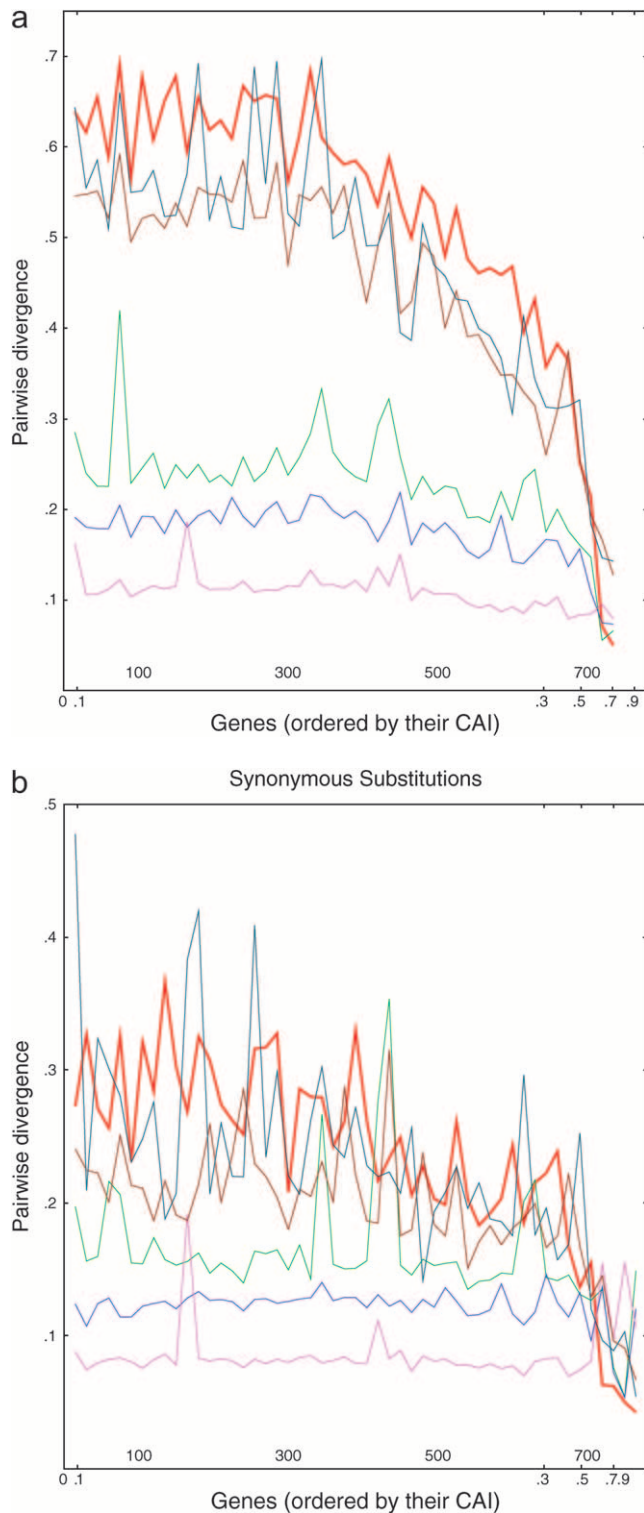


FIGURE 3.—Divergence as a function of expression in six yeasts. (a) Each of 764 ancient duplicated genes in *S. cerevisiae* is compared to its closest homolog in *S. cerevisiae* (intraspecies comparison) or to its closest homolog in another species of *Saccharomyces* (interspecies comparisons). From top to bottom, on the left side of the graph, the lines are *S. kluyveri* (thin line, blue-green), *S. cerevisiae* (thick line, red), *S. castellii* (brown), *S. bayanus* (green), *S. mikatae* (blue), and *S. paradoxus* (purple). Divergence scores were calculated (MATERIALS AND METHODS) for each pairwise comparison and plotted on the y-axis. Along the x-axis, the genes are sorted from left to right in the increasing order of the CAI of the chosen *cerevisiae* gene. Divergence scores were averaged over consecutive disjoint sets of 15 genes each to smooth the curves. The number of genes plotted in each curve varies from 752 to 758, depending on the existence/availability of interspecies homologs, whether the Jukes-Cantor score is defined, etc. b and c are the same as a, but are confined to synonymous and nonsynonymous substitutions, respectively.

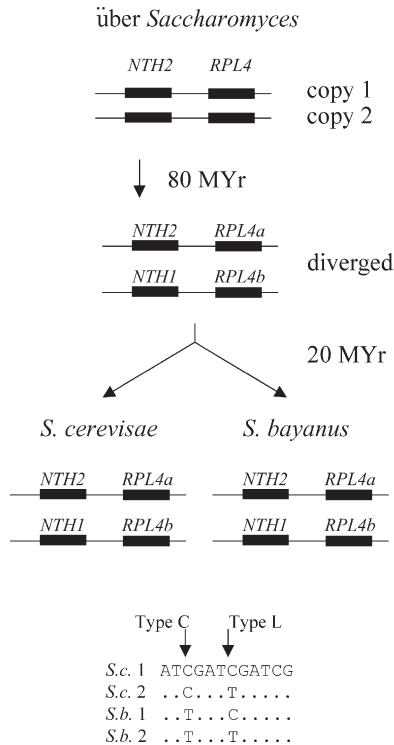


FIGURE 4.—Gene lineages and patterns of nucleotide substitution. The evolutionary lineages of the *NTH1/2* and *RPL4a/b* genes are shown. At a type C (correction) position, the intraspecies duplicates share a distinguishing nucleotide, while at a type L (lineage) position, the interspecies homologs, which are most closely related by descent, share a distinguishing nucleotide.

able, showing that there is still sequence space into which these genes can diverge while still maintaining function. Third and most striking, for highly expressed genes, there is very little divergence between the two *S. cerevisiae* copies. The red line in Figure 3 for the intraspecies *cerevisiae-cerevisiae* comparisons initially (*i.e.*, at lowest CAI) shows very high divergence scores, but then falls at higher CAIs, crossing through all the interspecies comparisons, until finally at the highest CAI the *cerevisiae-cerevisiae* comparisons have the lowest divergence. That is, two highly expressed *cerevisiae* copies may have only a few mismatches and >95% DNA sequence identity, despite the fact that the two genes diverged long ago, and despite the fact that many more mismatches are present between the same *S. cerevisiae* genes and their closest orthologs in all the other species, which diverged more recently. In summary, in intraspecies comparisons, we see a striking lack of divergence between pairs of very highly expressed genes, even though it is clear from interspecies comparisons that such genes can diverge. Since correction can occur within a species but not between species, we take this as evidence for correction.

The pattern of nucleotide substitution: Figure 4 shows the evolution of the *RPL4a/RPL4b* genes (encoding ri-

bosomal proteins) and the *NTH1/NTH2* genes (encoding neutral trehalase). In *S. cerevisiae*, *RPL4a* and *NTH2* are linked on chromosome 2, while their duplicates, *RPL4b* and *NTH1*, are linked on chromosome 4. These gene pairs are part of the “block 3” syntenic duplicated region defined by Wolfe and colleagues (<http://acer.gen.tcd.ie/~khwolfe/yeast/nova/>), and they are derived from the genome-wide duplication of 100 MYA. These gene duplicates have also survived in *S. bayanus*. Thus, we can now align four gene copies for each type of gene (*S.c.* *RPL4a*, *S.c.* *RPL4b*, *S.b.* *RPL4a*, and *S.b.* *RPL4b*, or *S.c.* *NTH1*, *S.c.* *NTH2*, *S.b.* *NTH1*, and *S.b.* *NTH2*) and ask about the patterns of nucleotide substitution.

Any pair of duplicated genes from the ancient duplication event has had ~80 MY in which to diverge before the separation of *S. cerevisiae* from *S. bayanus* ~20 MYA. Thus, in the absence of correction, one would expect *cerevisiae* copy 1 and its ortholog *bayanus* copy 1 to share certain nucleotide changes, while *cerevisiae* copy 2 and its ortholog *bayanus* copy 2 to share other changes, because copy 1 and copy 2 have had ~80 MY to diverge, while the two orthologs of copy 1 in the two species have had only 20 MY. We call this pattern of nucleotide substitution (where the orthologs in different species share a distinguishing nucleotide) the “L” pattern, for “lineage.” In contrast, if correction occurs, then *cerevisiae* copy 1 and *cerevisiae* copy 2 (*i.e.*, the paralogs) will share certain nucleotide changes (because the change has been copied from 1 to 2 or vice versa), while *bayanus* copy 1 and *bayanus* copy 2 will share other changes. We call this pattern of substitution (where paralogs within the species share a distinguishing nucleotide) the “C” pattern, for “correction.” The bottom of Figure 4 shows these two different patterns of nucleotide substitution. Finally, mutational noise will sometimes generate a situation in which *cerevisiae* copy 1 will share a distinguishing nucleotide with copy 2 (*i.e.*, the nonorthologous gene) in *bayanus*. We call this a type N pattern, for “noise”; its frequency is important for estimating the number of C and L patterns that might be due to noise.

Figure 5 shows a sample four-way alignment of part of the *RPL4a,b* and *NTH1,2* genes, and Table 1 shows results for the full-length four-way alignments. It is clear and striking that the highly expressed *RPL4* genes show exclusively the type C pattern of nucleotide substitution, arguing that they have undergone correction, while the tightly linked *NTH1,2* genes show mainly type L substitution (305 positions). Although 31 type C substitutions are seen in the *NTH1 vs. NTH2* comparison, there are also 35 type N substitutions, arguing that the type C substitutions in these genes are simply mutational noise, and not correction. Thus, as predicted, this pair of highly expressed genes shows primarily (in this case, exclusively) the correction pattern of substitution, while the poorly expressed but linked genes show primarily or (after allowing for noise) exclusively the lineage pat-

S.c. RPL4a ATCTCCACTTCTGATGTCACCAGAATTATCAACTCTTCTGAAATCCAATCTGCTATCAGA
S.c. RPL4b ATCTCCACTTCTGATGTCACCAGAATTATCAACTCTTCTGAAATCCAATCTGCTATCAGA
S.b. c664 ATCTCCACTTCTGATGTCACCAGAATCATCAACTCTTCTGAAATCCAATCTGCTGTTAGA
S.b. c21426 ATCTCCACTTCTGATGTCACCAGAATCATCAACTCTTCTGAAATCCAATCTGCTGTTAGA

Type C: 5. Type L: 0. Type N: 0.

S.c. NTH1 ATGAGTCAAGTTAATACAAAGCCAAGGACCGGTAGCCCAAGGCCGTCAAAGAAGATTATCA
S.c. NTH2 ATGGTAGATTTTTACCAAAAGTAAACGGAATAAATCCTCCATCCGAAGGTAATGATGGT
S.b. NTH1 ATGAGCAAAGTTAATGCAAAGCCAAGGACCTGTGGTTCAGGGCCGTCAAAGAAGATTATCA
S.b. NTH2 ATGGTAGAATTTTACCAAAAGTACAGAAATAAACCCGCCCGGATGTTGATATGATG

Type C: 0. Type L: 29. Type N: 1.

tern of substitution. We consider this very strong evidence for the correction model.

We extended this analysis to most of the syntenic duplicated blocks defined by Wolfe and co-workers (<http://acer.gen.tcd.ie/~khwolfe/yeast/nova/>). We examined all blocks (a) that contained at least one gene pair with a high codon bias and (b) where we could find two different *bayanus* orthologs of the high-codon-bias gene(s). There were ~28 eligible blocks. For all of these genes (both high and low bias) where two *bayanus* orthologs existed (168 genes total), we did the four-way alignments, as shown in Figure 5, and noted the number of C (correction), L (lineage), and N (noise) nucleotide substitutions. These results are shown in Table 2. The main results are as follows.

First, we compared the low-bias genes to the high-bias genes. The high-bias genes tend to have a high C/L ratio (weighted mean C/L ratio = 2.1), while the low-bias genes tend to have a low C/L ratio (weighted mean

FIGURE 5.—Alignment of the *RPL4a,b* and *NTH1,2* genes. A representative 60 nucleotides of the *RPL4a,b* and *NTH1,2* genes from *S. cerevisiae* (*S.c.*) and *S. bayanus* (*S.b.*) are aligned. Type C (correction) substitutions are underlined and in boldface type; type L (lineage) substitutions are in boldface type, and type N (noise) substitutions are underlined.

C/L ratio = 0.078), and these ratios differ significantly between the two groups ($P < 10^{-15}$ by a chi-square test).

Second, of 128 low-bias gene pairs, 126 have a low C/L ratio, as predicted. The two exceptions are *LYS20-LYS21* and *PPH21-PPH22*, both members of block 12. These two genes in block 12 are close together on chromosome 4 and in the same orientation. That is, *LYS20* is close to *LYS21*, and *PPH21* is close to *PPH22*, all on chromosome 4. This situation is perhaps favorable for DNA-DNA conversion (*e.g.*, during meiotic mispairing of these tightly linked tandem syntenic regions), and this may explain the high C/L ratio for these two gene pairs. That is, DNA-DNA-based correction may have occurred in the recent past, and DNA-DNA correction may not be as dependent on expression level as the RNA-DNA correction that we suggest for the majority of genes.

Third, of 33 high-bias genes, 20 have a C/L ratio >1 , as predicted, but 13 have a C/L ratio <1 ; for instance, in the most extreme case, the *RPS1A, RPS1B* gene pair has 5 correction substitutions, but 30 lineage substitutions. Thus a substantial minority of the high-bias genes do not seem to be undergoing correction. For the 20 high-bias genes with C/L >1 , the total number of events is 400 C, 57 L, and 29 N, for a C/L ratio of 7.0; and for the 13 high-bias genes with C/L <1 , the total number of events is 40 C, 154 L, and 17 N, for a C/L ratio of 0.26. The distribution of the normalized number of C substitutions for the 33 high-bias genes appears to be bimodal, and statistical tests reject the hypothesis of a single normally distributed population ($P < 10^{-7}$ using a Shapiro-Wilk normality test). Thus, there seem to be two kinds of high-bias genes: one kind that undergoes concerted evolution via correction and a second kind that does not. Possible reasons for these two populations are considered in the DISCUSSION.

Fourth, the 20 high-bias gene pairs with a C/L ratio >1 had a mean of 3.8% base-pair mismatches, while the 13 high-bias genes with a C/L ratio <1 had a mean of 8.2% base-pair mismatches, confirming that correction is associated with high sequence identity, while lack of correction is associated with divergence.

Fifth, we judged six gene pairs to be of medium codon bias. All six pairs showed the low C/L ratio typical of low-bias genes.

TABLE 1

Correction and lineage substitutions in the *RPL4* and *NTH1/2* genes

	Length	Mismatch positions	Type C	Type L	Type N
<i>RPL4a,b</i>	1089	69	40	0	0
<i>NTH1,2</i>	2355	970	31	305	35

“Length” is the length of the nucleotide alignment, and in these cases is the length of the open reading frame of the gene. “Mismatch positions” is the number of positions where all four nucleotides in the alignment are not identical, *e.g.*, an alignment of T:T:T:C or T:T:C:C, etc. “Type C” is a mismatch position of the C type, characteristic of correction, where the two *cerevisiae* genes share the same nucleotide and the two *bayanus* genes share a different nucleotide (*e.g.*, T:T:C:C). “Type L” is a mismatch position of the L type, characteristic of the gene lineage, where *cerevisiae* gene 1 and its *bayanus* homolog share the same nucleotide, while *cerevisiae* gene 2 and its *bayanus* homolog share a different nucleotide (*e.g.*, T:C:T:C). “Type N” is a mismatch position characteristic of mutational noise, where *cerevisiae* gene 1 and its *bayanus* nonhomolog share the same nucleotide, while *cerevisiae* gene 2 and its *bayanus* nonhomolog share a different nucleotide (*e.g.*, T:C:C:T). Positions of the type T:T:T:C (*i.e.*, one nucleotide at odds with the other three) are mismatch positions, but are not C or L or N.

TABLE 2
Ratios of correction to lineage substitutions for duplicated genes

Block no.	Gene pair	CAI	C	L	N	C/L
7	<i>NHP6B-NHP6A</i>	Low	3	26	5	
7	<i>YMC2-YMC1</i>	Low	9	132	10	
7	<i>TKL2-TKL1</i>	Low	26	320	18	
7	<i>TEF2-TEF1</i>	High	44	0	0	High
8	<i>SSE2-SSE1</i>	Medium	26	252	25	
8	<i>SMY2-YPL105C</i>	Low	28	504	32	
8	<i>YBR177C-YPL095C</i>	Low	19	243	27	
8	<i>RPS6B-RPS6A</i>	High	12	0	0	High
8	<i>SMP1-RLM1</i>	Low	6	74	4	
8	<i>YBR183W-YPL087W</i>	Low	18	203	9	
8	<i>RPS9B-RPS9A</i>	High	17	2	2	High
8	<i>RPL21A-RPL21B</i>	High	4	14	1	Low
8	<i>YBR197C-YPL077C</i>	Low	7	163	10	
8	<i>KTR4-KTR6</i>	Low	15	426	9	
8	<i>KTR3-KTR6</i>	Low	10	343	11	
10	<i>RPS14A-RPS14B</i>	High	8	0	0	High
12	<i>ARF2-ARF1</i>	Low	9	19	2	
12	<i>RPL35B-RPL35A</i>	High	13	1	0	High
12	<i>PPH21-PPH22</i>	Low	87	20	6	High
12	<i>LYS21-LYS20</i>	Low	116	37	3	High
15	<i>BDF2-BDF1</i>	Low	20	482	19	
15	<i>RPS29B-RPS29A</i>	High	2	8	1	Low
19	<i>PPZ2-PPZ1</i>	Low	29	335	28	
19	<i>YDR438W-YML018C</i>	Low	15	251	8	
19	<i>YDR450W-RPS18B</i>	High	1	4	2	Low
19	<i>YDR451C-YOX1</i>	Low	19	240	9	
26	<i>RPS26B-RPS26A</i>	High	4	12	0	Low
26	<i>PMD1-MDS3</i>	Low	62	935	64	
29	<i>YGR221C-YHR149C</i>	Low	18	431	20	
29	<i>YGR230W-SPO12</i>	Low	6	77	9	
29	<i>KEL2-KEL1</i>	Low	36	499	34	
29	<i>YAP1802-YAP1801</i>	Low	27	387	20	
29	<i>YGR243W-YHR162W</i>	Low	4	64	5	
29	<i>SOL4-SOL3</i>	Low	16	147	13	
29	<i>ENO1-ENO2</i>	High	48	14	7	High
29	<i>GND2-GND1</i>	Medium	20	100	14	
30	<i>YGR004W-YLR324W</i>	Low	14	365	14	
30	<i>STF2-YLR327C</i>	Low	2	37	4	
30	<i>YGR010W-YLR328W</i>	Low	17	175	11	
30	<i>RPS25A-RPS25B</i>	High	2	3	0	Low
30	<i>ORM1-YLR350W</i>	Low	8	109	5	
30	<i>BUD9-BUD8</i>	Low	15	210	24	
30	<i>YGR043C-TAL1</i>	Medium	13	164	8	
30	<i>SCM4-YLR356W</i>	Low	6	158	4	
30	<i>RSC1-RSC2</i>	Low	26	553	33	
30	<i>ROM1-ROM2</i>	Low	54	673	46	
30	<i>YGR071C-YLR373C</i>	Low	35	671	31	
32	<i>YGL139W-YPL221W</i>	Low	33	371	26	
32	<i>RPL1B-RPL1A</i>	High	29	1	0	High
32	<i>PCL10-PCL8</i>	Low	21	324	14	
32	<i>YGL133W-YPL216W</i>	Low	50	724	50	
33	<i>YGL084C-YPL189W</i>	Low	35	273	33	
33	<i>YGL082W-YPL191C</i>	Low	18	217	13	
33	<i>RPL7A-RPL7B</i>	High	37	5	1	High
33	<i>AFT1-YPL202C</i>	Low	9	388	10	
33	<i>PUS2-PUS1</i>	Low	19	266	14	
34	<i>RPL11B-RPL11A</i>	High	30	7	0	High
34	<i>DBF2-DBF20</i>	Low	23	215	18	
34	<i>ASK10-YPR115W</i>	Low	46	628	44	

(continued)

TABLE 2
(Continued)

Block no.	Gene pair	CAI	C	L	N	C/L
34	<i>CLB1-CLB2</i>	Low	13	100	17	
34	<i>CLB6-CLB5</i>	Low	14	220	22	
34	<i>RPS23A-RPS23B</i>	High	5	0	3	High
34	<i>MEP1-MEP3</i>	Low	16	188	21	
34	<i>ASN2-ASN1</i>	Low	45	112	20	
34	<i>YGR131W-NCE102</i>	Low	7	95	8	
34	<i>YGR136W-YPR154W</i>	Low	12	88	8	
34	<i>YGR141W-YPR157W</i>	Low	26	230	21	
34	<i>SKN1-KRE6</i>	Low	43	290	31	
35	<i>YHL017W-PTM1</i>	Low	6	176	7	
35	<i>YHL012W-UGP1</i>	Low	17	365	8	
35	<i>LAG1-YKL008C</i>	Low	12	126	15	
35	<i>RPL14B-RPL14A</i>	High	13	5	0	High
35	<i>YHR001W-YKR003W</i>	Low	18	180	24	
37	<i>YHR115C-YNL116W</i>	Low	25	228	18	
37	<i>TOM72-TOM70</i>	Low	25	385	24	
37	<i>EPT1-CPT1</i>	Low	15	215	13	
37	<i>YHR131C-YNL144C</i>	Low	32	477	36	
37	<i>YHR133C-YNL156C</i>	Low	6	220	9	
37	<i>YCK1-YCK2</i>	Low	25	243	22	
37	<i>SPS100-YGP1</i>	Medium	12	212	10	
37	<i>RPL42B-RPL42A</i>	High	3	2	0	High
38	<i>UBP7-UBP11</i>	Low	29	426	28	
38	<i>YIL151C-YKR096W</i>	Low	57	637	52	
38	<i>RPL40A-RPL40B</i>	High	0	8	3	Low
39	<i>TPM2-TPM1</i>	Low	5	86	4	
39	<i>RPL16A-RPL16B</i>	High	5	14	3	Low
39	<i>FKH1-FKH2</i>	Low	23	291	17	
39	<i>SIM1-SUN4</i>	Low	5	11	4	
39	<i>YIL121W-YNL065W</i>	Low	21	390	17	
39	<i>YIL120W-YNL065W</i>	Low	21	406	19	
39	<i>POR2-POR1</i>	Low	17	153	9	
39	<i>YIL113W-MSG5</i>	Low	11	100	10	
39	<i>COX5B-COX5A</i>	Low	9	51	3	
39	<i>SEC24-YNL049C</i>	Low	41	449	29	
39	<i>YIL105C-YNL047C</i>	Low	29	373	28	
39	<i>PRK1-ARK1</i>	Low	24	336	15	
40	<i>RPL17B-RPL17A</i>	High	5	14	1	Low
40	<i>HAL5-KKQ8</i>	Low	31	459	32	
40	<i>TPK1-TPK3</i>	Low	22	162	13	
40	<i>CIS3-PIR3</i>	High	5	99	9	
41	<i>YUR1-KTR2</i>	Low	25	212	13	
41	<i>TIF2-TIF1</i>	High	58	1	0	High
41	<i>GLG2-GLG1</i>	Low	22	261	20	
41	<i>RPS21B-RPS21A</i>	High	5	2	5	High
41	<i>LCB3-LBP2</i>	Low	14	133	9	
41	<i>MRS3-MRS4</i>	Low	7	140	9	
41	<i>TRK1-TRK2</i>	Low	42	456	35	
41	<i>NCA3-UTH1</i>	Low	22	147	9	
41	<i>YJL112W-CAF4</i>	Low	27	423	18	
41	<i>GZF3-DAL80</i>	Low	7	169	7	
41	<i>YJL105W-YKR029C</i>	Low	15	459	13	
41	<i>CHS6-YKR027W</i>	Low	22	408	23	
41	<i>SAP185-SAP190</i>	Low	46	594	36	
41	<i>YJL084C-YKR021W</i>	Low	39	623	30	
41	<i>YJL083W-IRS4</i>	Low	0	2	0	
41	<i>YJL082W-YKR018C</i>	Low	32	354	35	
44	<i>CNA1-CNA2</i>	Low	25	258	22	
44	<i>RPS1A-RPS1B</i>	High	5	30	3	Low
44	<i>SIR3-ORC1</i>	Low	30	743	30	

(continued)

TABLE 2
(Continued)

Block no.	Gene pair	CAI	C	L	N	C/L
44	<i>RPL6B-RPL6A</i>	High	6	14	3	Low
44	<i>FPR4-FPR3</i>	Low	10	188	15	
44	<i>HMG2-HMG1</i>	Low	47	575	41	
45	<i>YLR266C-YRR1</i>	Low	38	505	38	
45	<i>YLR270W-YOR173W</i>	Low	12	179	10	
45	<i>BRR5-YOR179C</i>	Low	9	90	13	
45	<i>RPS30A-RPS30B</i>	High	7	0	0	High
45	<i>GSP1-GSP2</i>	Medium	4	47	4	
45	<i>EXG1-SPR1</i>	Low	17	197	10	
47	<i>YMR222C-YOR280C</i>	Low	8	151	13	
47	<i>RPS10B-RPS10A</i>	High	3	10	0	Low
47	<i>YMR233W-YOR295W</i>	Low	11	168	14	
47	<i>YMR237W-BUD7</i>	Low	32	342	22	
47	<i>RPL20A-RPL20B</i>	High	12	10	0	High
47	<i>ZRC1-COT1</i>	Low	14	223	18	
47	<i>FAA4-FAA1</i>	Low	51	307	35	
48	<i>MMT1-MMT2</i>	Low	16	245	18	
48	<i>YMR180C-CET1</i>	Low	14	229	15	
48	<i>YMR181C-YPL229W</i>	Low	6	120	5	
48	<i>RGM1-YPL230W</i>	Low	14	117	9	
48	<i>SSO2-SSO1</i>	Low	10	119	12	
48	<i>YMR192W-YPL249C</i>	Low	26	500	24	
48	<i>RPL36A-RPL36B</i>	High	3	11	1	Low
48	<i>YMR195W-YPL250C</i>	Low	6	85	2	
48	<i>CIK1-VIK1</i>	Low	21	484	26	
48	<i>CLN1-CLN2</i>	Low	12	284	14	
49	<i>MCK1-YOL128C</i>	Low	15	305	13	
49	<i>RPS19B-RPS19A</i>	High	8	4	0	High
49	<i>TRF5-TRF4</i>	Low	21	331	22	
49	<i>CLA4-SKM1</i>	Low	30	369	25	
49	<i>MSB3-MSB4</i>	Low	21	294	25	
49	<i>RFC3-RFC4</i>	Low	12	172	7	
51	<i>DED1-DBP1</i>	Low	20	291	18	
51	<i>YOR222W-YPL134C</i>	Low	12	163	4	
51	<i>YOR226C-YPL135W</i>	Low	5	65	5	
51	<i>YOR227W-YPL137C</i>	Low	60	707	42	
51	<i>YOR229W-UME1</i>	Low	12	322	12	
51	<i>WTM1-UME1</i>	Low	15	308	11	
51	<i>MKK1-MKK2</i>	Low	28	260	24	
51	<i>KIN4-YPL141C</i>	Low	27	482	35	
51	<i>RPL33B-RPL33A</i>	High	2	7	2	Low
51	<i>HES1-KES1</i>	Low	18	197	12	
IV:VIII	<i>STP1-STP2</i>	Low	24	319	17	
IV:VIII	<i>RPL27B-RPL27A</i>	High	4	19	0	Low
VII:VII	<i>TIF4631-TIF4632</i>	Low	35	492	29	
VII:VII	<i>RPL24B-RPL24A</i>	High	15	2	4	High
VII:X	<i>RNR4-RNR2</i>	Medium	25	215	3	
VII:X	<i>BUB1-MAD3</i>	Low	20	351	24	
VII:X	<i>TDH3-TDH2</i>	High	10	0	7	High
VIII:X	<i>RPS4B-RPS4A</i>	High	26	1	0	High

“Block no.” is from Wolfe and colleagues (<http://acer.gen.tcd.ie/~khwolfe/yeast/nova/>). Blocks were analyzed only if they contained at least one gene pair of high CAI (see below), and gene pairs in such blocks were analyzed only if two different homologs could be found in *S. bayanus* (i.e., analysis was carried out only when it was possible to make the four-way alignment). CAI was considered “high” if both of the *cerevisiae* genes had a CAI >0.80; CAI was considered “medium” if at least one of the *cerevisiae* genes had a CAI >0.45 but at least one was <0.80; otherwise, CAI was considered “low.” C is the number of correction substitutions in the four-way alignment; L, the number of lineage substitutions in the four-way alignment; N, the number of noise substitutions in the four-way alignment. For all gene pairs with high CAI, the C/L ratio is noted as “high” (>1) or “low” (<1). With two exceptions, all gene pairs with a low CAI had low C/L ratios, and so are not noted. The two exceptions are noted by italics. Correction, lineage, and noise-type substitutions are defined in Figures 4 and 5 and Table 1.

TABLE 3
Preferred codons in six species of *Saccharomyces*

Amino acid	<i>S. bayanus</i>	<i>S. castelli</i>	<i>S. kluyveri</i>	<i>S. mikatae</i>	<i>S. paradoxus</i>	<i>S. cerevisiae</i>
Ile	ATC 0.66	<i>ATT 0.51</i>	ATC 0.70	ATC 0.54	ATC 0.56	ATC 0.58
Asn	AAC 0.96	AAC 1.00	AAC 1.00	AAC 0.92	AAC 0.93	AAC 0.94
Asp	GAC 0.65	<i>GAT 0.50</i>	GAC 0.70	GAC 0.63	GAC 0.66	GAC 0.65
Gln	CAA 0.99	CAA 0.98	CAA 1.00	CAA 0.98	CAA 0.99	CAA 1.00
Ala	GCT 0.68	GCT 0.69	GCT 0.73	GCT 0.74	GCT 0.78	GCT 0.80
His	CAC 0.82	CAC 0.71	CAC 1.00	CAC 0.84	CAC 0.88	CAC 0.90
Thr	<i>ACT 0.52</i>	<i>ACC 0.50</i>	<i>ACC 0.60</i>	<i>ACT 0.51</i>	<i>ACC 0.52</i>	<i>ACT 0.50</i>
Tyr	TAC 0.91	TAC 0.90	TAC 1.00	TAC 0.85	TAC 0.91	TAC 0.92
Glu	GAA 0.97	GAA 0.99	GAA 0.94	GAA 0.95	GAA 0.98	GAA 0.98
Pro	CCA 0.91	CCA 0.92	CCA 0.96	CCA 0.90	CCA 0.92	CCA 0.94
Leu	TTG 0.83	TTG 0.74	TTG 0.95	TTG 0.80	TTG 0.83	TTG 0.89
Phe	TTC 0.84	TTC 0.85	TTC 0.86	TTC 0.84	TTC 0.83	TTC 0.81
Gly	GGT 0.94	GGT 0.96	GGT 0.95	GGT 0.96	GGT 0.96	GGT 0.96
Lys	AAG 0.88	AAG 0.89	AAG 1.00	AAG 0.84	AAG 0.84	AAG 0.85
Trm	TAA 0.90	TAA 1.00	TAA 1.00	TAA 0.62	TAA 0.80	TAA 0.90
Arg	AGA 0.84	AGA 0.89	AGA 0.93	AGA 0.88	AGA 0.84	AGA 0.85
Cys	TGT 0.86	TGT 0.83	TGT 1.00	TGT 1.00	TGT 1.00	TGT 1.00
Val	<i>GTC 0.51</i>	GTT 0.50	GTT 0.52	GTT 0.60	GTT 0.55	GTT 0.55
Ser	<i>TCC 0.51</i>	TCT 0.58	TCT 0.56	TCT 0.53	TCT 0.52	TCT 0.50

Ten highly expressed genes of *S. cerevisiae* (*RPL11A*, *ENO1*, *TDH1*, *RPL4A*, *RPL8A*, *RPL9A*, *RPL15A*, *RPS2*, *RPS3*, and *RPS5*) were selected, and their full-length closest homologs were identified in the other yeasts whenever possible. For each amino acid, the preferred codon and its frequency is listed for each yeast over the 10 selected proteins. Amino acids with only one codon (Met, Trp) are omitted. For most amino acids, all yeasts had the same preferred codon. For Ile, Asp, Thr, Val, and Ser, there were minor differences (indicated by italics). All of these minor differences occur when there are two commonly used codons, each with a frequency of close to 50%. For instance, in *S. castelli*, the preferred codon for Ile is ATT (frequency of 0.50), but the second-most preferred codon is ATC (frequency of 0.48). Similarly, in *S. castelli*, the preferred codon for Asp is GAT (frequency 0.50), but the second-most preferred codon is GAC (frequency 0.50). In the case of Thr, only two codons are substantially used, ACT and ACC, and these have nearly an equal frequency in each yeast. Similarly, for Ser, only TCC and TCT are substantially used, each at ~ 0.5 in each yeast.

Counterarguments: We have considered several alternative explanations for the unexpectedly high conservation between highly expressed *S. cerevisiae* gene pairs. One obvious alternative is that these gene pairs are not derived from the ancient genome-wide duplication event, but instead are the result of a much more recent chromosomal duplication. This is the case for several duplications found near telomeres, which we do not consider here. However, it appears not to be the case for the gene pairs that we consider here, because the unexpectedly high degree of homology among the gene pairs that we are considering ends abruptly at the boundary of the gene's open reading frame. The 5' and 3' noncoding regions of these genes do not show a strikingly high level of conservation. Furthermore, these genes are typically embedded in syntenic, duplicated regions, and there are typically poorly expressed, poorly conserved duplicated genes flanking the highly expressed, highly conserved genes. The exceptional genes in block 12, *LYS20-LYS21* and *PPH21-PPH22*, could be a recent duplication.

A second alternative is that the preferred codons are different in different species. In this case, the lack of

intra-*cerevisiae* divergence at high CAI would be explained by the need to maintain preferred codons, while the presence of interspecies divergence would be explained by highly expressed genes evolving to conform to a different, species-specific codon bias. However, three findings argue against this possibility. First, the preferred codons seem to be the same in each of the six species (Table 3). Second, a substantial proportion of the interspecies divergence (at least between *S. cerevisiae* and *S. bayanus*) is due to nonsynonymous base changes (Table 4), which of course is not explainable by differences in codon bias. Third, this argument does not explain gene families with a high proportion of correction nucleotide substitutions.

Patterns and properties of correction: Assuming that the unexpectedly high degree of conservation between pairs of highly expressed genes does reflect recombinational correction, we can draw some inferences about the properties of correction. First, correction ends abruptly at the boundaries of identity. This can be seen most easily at the beginning and end of each open reading frame. Within the open reading frame, the percentage identity can be 95% or more, but, immediately

TABLE 4
Synonymous and nonsynonymous nucleotide changes

Gene	Length	C-C Syn	C-C Non	C-B Syn	C-B Non	C-M Syn	C-M Non
<i>eno1</i>	1314	41	27	68	30	40	2
<i>rpl11a</i>	525	14	2	34	4	19	0
<i>rpl1a</i>	654	4	0	31	5	31	0
<i>rps8a</i>	603	11	0	25	7	10	1
<i>tef1</i>	1377	2	0	32	14	19	1
<i>tif51a</i>	474	28	15	30	10	16	6

Six highly expressed genes of *S. cerevisiae* were compared to their closest homolog in *S. cerevisiae* (“C-C” comparisons), in *S. bayanus* (“C-B” comparisons), or in *S. mikatae* (“C-M” comparisons). The number of synonymous (“Syn”) and nonsynonymous (“Non”) changes are tabulated.

outside the open reading frame, identity decays to essentially random levels (data not shown). The same effect can be seen in genes with multiple exons. Although identity may be high within each exon, the introns show roughly random levels of identity.

If correction ends at the end of a tract of high homology, then the correction of different exons of the same gene may be independent events. In this case, the frequency of correction should be proportional to the length of the exon (since longer exons have an increased chance of interacting with each other). Indeed, in pairs of genes with multiple exons we have found that the degree of sequence identity between the first exons of a pair of genes can be different from the sequence identity between the second exons of the same pair of genes. Furthermore, longer exons typically have higher degrees of identity than shorter exons (although there are exceptions; data not shown).

Exon length and RNA-DNA correction: If correction occurs via a cDNA intermediate, then this cDNA can correct the same gene that originally generated the cDNA, as well as any copy of the gene. Such self-correction would not have any effect on the nucleotide sequence of the open reading frame, since this would be identical between the cDNA and the gene. However, on (rare?) occasions when correction proceeded past a boundary of high sequence identity, it could remove an intron from the gene. Indeed, it has been proposed that this kind of self-correction is responsible for removing most of the (presumed) originally existing introns from the genome of *S. cerevisiae* (FINK 1987). The few introns that remain tend to be at the extreme 5'-end of the gene, suggesting that correction begins, or is more probable, at the 3'-end of the gene.

Interestingly, in many species, exons at the 3'-end of a gene tend to be longer than exons at the 5'-end of a gene (Table 5; see also XIA *et al.* 2003). This is consistent with the idea that, in many species, self-correction occurs via a cDNA intermediate, beginning at the 3'-end of the gene. Such self-correction would tend to remove 3' introns, thus generating abnormally long 3' exons.

DISCUSSION

For highly expressed gene pairs, within-species divergence is significantly less than between-species divergence, even though the within-species pairs have had a longer time to diverge (Figure 3). Furthermore, for the majority of highly expressed genes, the correction pattern of nucleotide substitution is much more common than the lineage pattern of substitution, while the opposite is true for poorly expressed genes (Figures 4 and 5; Tables 1 and 2). These observations are very difficult to explain by selection alone. We believe that selection and correction are synergistic with each other for highly expressed genes; selection for both function and codon bias tends to minimize the rate of drift, and the resulting high level of sequence identity keeps the gene pairs eligible for correction, which fully restores sequence identity between duplicates. There is less selection in poorly expressed genes, since codon bias is of little or no importance. In addition, there is less selection in the 5' or 3' regions of genes or within introns, since many base changes in these regions have little or no impact on gene function. Thus, poorly expressed genes, 5' and 3' untranslated regions, and introns drift more rapidly and soon diverge to the point where sequence identity is too low to allow a recombinational interaction. After this point, they are no longer eligible for correction and continue to drift apart with time.

Surprisingly, we found a substantial minority of highly expressed genes that have a low C/L ratio, *i.e.*, that appear not to have corrected in the 20 MY since the split between *cerevisiae* and *bayanus* (Table 2). Why should some gene pairs fail to correct? In the context of the RNA-DNA correction model, one possibility is that some genes are more readily reverse transcribed than others. XU and BOEKE (1990) found that some cellular mRNAs copurified with Ty virus-like particles (VLPs; *TRP1*, *HIS3*, *RPS17a*), while other mRNAs did not (*ACT1*, *GALI*, *PYK1*). The mRNAs copurifying with Ty VLPs may actually have been packaged within the particles, and so these mRNAs would presumably be more likely to be reverse transcribed than mRNAs not

TABLE 5
Distribution of exon lengths

	1-exon genes:	2-exon genes		3-exon genes		
	Length of exon 1	Length of exon 1	Length of exon 2	Length of exon 1	Length of exon 2	Length of exon 3
<i>Drosophila</i>	934	403	643	271	454	458
<i>Arabidopsis</i>	970	441	489	346	284	402
<i>S. cerevisiae</i>	1419	313	1260	69	149	238
<i>S. pombe</i>	1462	301	946	206	279	721
<i>Caenorhabditis elegans</i>	611	229	327	178	267	244

Mean exon lengths for genes of 1, 2, and 3 exons are given.

so packaged. Interestingly, *RPS17a*, which copurifies with Ty VLPs, is highly similar to its paralog *RPS17b*, while *PYK1*, which has a CAI similar to *RPS17a*, but which does not copurify with Ty VLPs, is highly diverged from its paralog, *PYK2*. A second possibility is that the genes that fail to correct are those where the two paralogs diverged significantly by chance before the *cerevisiae-bayanus* split and, because of the divergence, were no longer eligible for gene conversion and correction. A third possibility (see below) is that the two copies have taken on somewhat different cellular roles, and so both genes are needed.

Correction could occur by a DNA-DNA interaction between the two genes of the pair or by a cDNA-DNA (the RNA-DNA model) interaction occurring after reverse transcription of an mRNA. Our evidence does not distinguish these two models. The RNA-DNA model has several appealing features. First, it gives a clear expectation that correction should be more prevalent for highly expressed genes. Second, it is widely believed that self-correction via a cDNA does occur in *S. cerevisiae*, and if self-correction can occur, then correction of a copy should also occur. Third, it explains why the introns of highly expressed genes are not conserved, whereas a DNA-DNA interaction between two chromosomal genes might tend to correct these introns as well as flanking exons. On the other hand, DNA-DNA events seem to be much more frequent than RNA-DNA events (DERR and STRATHERN 1993), and so one would expect correction to be dominated by the DNA-DNA mechanism, even if RNA-DNA events sometimes occurred.

It is unclear to what extent similar events may occur in other organisms. *S. cerevisiae* has a highly active system for homologous recombination and thus is especially suited to correction. However, most other organisms also have homologous recombination, and many or most other eukaryotes contain reverse transcriptases. We therefore imagine that correction could occur at some level in many or most other organisms. Table 5 shows that 3' exons are typically longer than 5' exons for many organisms, and this is consistent with correction. ZHANG *et al.* (2002) have shown that there are >2000

reverse-transcribed ribosomal protein pseudogenes in the human genome, showing that the conversion of a transcript to a cDNA is a reasonably common event in humans. ZHANG *et al.* (2002) have also shown the existence of a number of duplicate ribosomal protein genes. If any of these duplicates predate the divergence between, *e.g.*, humans and mice, then analysis of these duplicates, such as we have done here with *Saccharomyces*, may show whether these pairs are maintained by correction in mammals.

We note that correction could occur between a highly expressed gene and an unexpressed pseudogene. We have preliminary data from *S. cerevisiae* suggesting that, in a few cases, one member of a pair of ribosomal protein genes is expressed poorly; possibly such poorly expressed genes are maintained by correction from their highly expressed twin. Furthermore, this idea could explain the maintenance of the large number of nonmutated ribosomal protein pseudogenes that are present in the human genome (ZHANG *et al.* 2002).

About 90% of the genes originally duplicated in the ancient duplication event have since been deleted, while ~10% remain as duplicates. Why do these 10% remain? Our results and the recent results of KELLIS *et al.* (2004) allow us to point to two kinds of reasons. First, for poorly expressed genes, one member of the gene pair seems to evolve quickly, gaining many substitutions rapidly and acquiring a new biological role (KELLIS *et al.* 2004). Thus, for poorly expressed genes, sequence divergence is favorable for maintaining the copy. Second, for highly expressed genes, we now argue that the duplication aids in making large amounts of protein in cases where large amounts of protein are needed. Thus, for highly expressed genes, sequence conservation is favorable for maintaining the copy. Figure 3, b and c, supports this view, because it shows that for poorly expressed genes, nonsynonymous substitutions are relatively favored (*i.e.*, promoting divergence of protein function), while for highly expressed genes, synonymous substitutions are relatively favored (*i.e.*, conserving protein function). Correction fits into this scheme well, since correction seems to work only on highly expressed genes, which

are precisely the genes where sequence conservation, and not divergence, leads to preservation of the copy. It is interesting to speculate that the rare, highly expressed genes not showing copy correction are highly expressed genes that have nevertheless evolved to take on new cellular roles. It is interesting to note that deletion of either copy of *RPS1a/b* (the most extremely diverged high-bias gene pair) leads to severe growth defects (Saccharomyces Genome Database), suggesting that the two copies may have nonoverlapping roles.

This work was sponsored by National Institutes of Health grants GM39978 and GM648131 to B.F. and National Science Foundation grant EIA0325123 to S.S.

LITERATURE CITED

- BALTIMORE, D., 1985 Retroviruses and retrotransposons: the role of reverse transcription in shaping the eukaryotic genome. *Cell* **40**: 481–482.
- CHEN, W., and S. JINKS-ROBERTSON, 1998 Mismatch repair proteins regulate heteroduplex formation during mitotic recombination in yeast. *Mol. Cell. Biol.* **18**: 6525–6537.
- CLIFTEN, P., P. SUDARSANAM, A. DESIKAN, L. FULTON, B. FULTON *et al.*, 2003 Finding functional features in Saccharomyces genomes by phylogenetic footprinting. *Science* **301**: 71–76.
- DATTA, A., M. HENDRIX, M. LIPSITCH and S. JINKS-ROBERTSON, 1997 Dual roles for DNA sequence identity and the mismatch repair system in the regulation of mitotic crossing-over in yeast. *Proc. Natl. Acad. Sci. USA* **94**: 9757–9762.
- DERR, L. K., and J. N. STRATHERN, 1993 A role for reverse transcripts in gene conversion. *Nature* **361**: 170–173.
- DERR, L. K., J. N. STRATHERN and D. J. GARFINKEL, 1991 RNA-mediated recombination in *S. cerevisiae*. *Cell* **67**: 355–364.
- FINK, G. R., 1987 Pseudogenes in yeast? *Cell* **49**: 5–6.
- FUTCHER, B., G. I. LATTER, P. MONARDO, C. S. McLAUGHLIN and J. I. GARRELS, 1999 A sampling of the yeast proteome. *Mol. Cell. Biol.* **19**: 7357–7368.
- KELLIS, M., N. PATTERSON, M. ENDRIZZI, B. BIRREN and E. S. LANDER, 2003 Sequencing and comparison of yeast species to identify genes and regulatory elements. *Nature* **423**: 241–254.
- KELLIS, M., B. W. BIRREN and E. S. LANDER, 2004 Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617–624.
- MODRICH, P., and R. LAHUE, 1996 Mismatch repair in replication fidelity, genetic recombination, and cancer biology. *Annu. Rev. Biochem.* **65**: 101–133.
- NEI, M., and T. GOJOBORI, 1986 Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol. Biol. Evol.* **3**: 418–426.
- POWELL, J. R., and E. N. MORIYAMA, 1997 Evolution of codon usage bias in *Drosophila*. *Proc. Natl. Acad. Sci. USA* **94**: 7784–7790.
- SAXE, D., A. DATTA and S. JINKS-ROBERTSON, 2000 Stimulation of mitotic recombination events by high levels of RNA polymerase II transcription in yeast. *Mol. Cell. Biol.* **20**: 5404–5414.
- SEOIGHE, C., and K. H. WOLFE, 1999 Updated map of duplicated regions in the yeast genome. *Gene* **238**: 253–261.
- SHARP, P. M., and W. H. LI, 1987 The codon adaptation index—a measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.* **15**: 1281–1295.
- WOLFE, K. H., and D. C. SHIELDS, 1997 Molecular evidence for an ancient duplication of the entire yeast genome. *Nature* **387**: 708–713.
- XIA, X., Z. XIE and W. H. LI, 2003 Effects of GC content and mutational pressure on the lengths of exons and coding sequences. *J. Mol. Evol.* **56**: 362–370.
- XU, H., and J. D. BOEKE, 1990 Localization of sequences required in cis for yeast Ty1 element transposition near the long terminal repeats: analysis of mini-Ty1 elements. *Mol. Cell. Biol.* **10**: 2695–2702.
- ZHANG, Z., P. HARRISON and M. GERSTEIN, 2002 Identification and analysis of over 2000 ribosomal protein pseudogenes in the human genome. *Genome Res.* **12**: 1466–1482.

Communicating editor: M. JOHNSTON

