

# Polygenic and directional regulatory evolution across pathways in *Saccharomyces*

James H. Bullard<sup>a</sup>, Yulia Mostovoy<sup>b</sup>, Sandrine Dudoit<sup>a,c</sup>, and Rachel B. Brem<sup>b,1</sup>

<sup>a</sup>Division of Biostatistics, <sup>b</sup>Department of Molecular and Cell Biology, and <sup>c</sup>Department of Statistics, University of California, Berkeley, CA 94720

Edited\* by David Botstein, Lewis-Sigler Institute, Princeton, NJ, and approved February 1, 2010 (received for review November 10, 2009)

The search to understand how genomes innovate in response to selection dominates the field of evolutionary biology. Powerful molecular evolution approaches have been developed to test individual loci for signatures of selection. In many cases, however, an organism's response to changes in selective pressure may be mediated by multiple genes, whose products function together in a cellular process or pathway. Here we assess the prevalence of polygenic evolution in pathways in the yeasts *Saccharomyces cerevisiae* and *S. bayanus*. We first established short-read sequencing methods to detect *cis*-regulatory variation in a diploid hybrid between the species. We then tested for the scenario in which selective pressure in one species to increase or decrease the activity of a pathway has driven the accumulation of *cis*-regulatory variants that act in the same direction on gene expression. Application of this test revealed a variety of yeast pathways with evidence for directional regulatory evolution. In parallel, we also used population genomic sequencing data to compare protein and *cis*-regulatory variation within and between species. We identified pathways with evidence for divergence within *S. cerevisiae*, and we detected signatures of positive selection between *S. cerevisiae* and *S. bayanus*. Our results point to polygenic, pathway-level change as a common evolutionary mechanism among yeasts. We suggest that pathway analyses, including our test for directional regulatory evolution, will prove to be a relevant and powerful strategy in many evolutionary genomic applications.

gene regulation | adaptation | yeast

A main challenge of evolutionary biology is to understand the influence of selection on genetic variation within and between species. Classically, molecular evolution methods have targeted individual genes or loci for tests of selection. Their successes have uncovered both protein-coding and regulatory variants with signatures of non-neutral evolution (1–3). However, decades of quantitative genetic mapping studies indicate that in most cases, variation in phenotype between individuals is the result of multiple sequence changes at unlinked loci (4). The genetic response to changes in selective pressure is likely to follow the same pattern, but, to date, methods for identifying cases of polygenic evolution have been at a premium.

Some of the strongest evidence for polygenic adaptation has emerged from the study of protein-coding variants between phylogenetic lineages. Signatures of positive selection can be detected in the protein-coding sequences of groups of genes of related function (5–8), suggestive of a coherent series of genetic changes accumulated by a population in response to selection. Additionally, genetic variation in gene expression represents a rich data source for signatures of selection, both positive (9, 10) and purifying (11, 12), although the mechanisms that govern regulatory evolution are not fully understood. Regulatory variants can act in *cis*, to impact the expression of a neighboring gene, or in *trans*, targeting the expression of genes elsewhere in the genome. Hallmarks of positive selection have been observed at individually mapped *cis*-regulatory variants (1, 3); one appealing model predicts that such a locus is part of a suite of adaptive regulatory changes in downstream effectors of a pathway (13–16). More generally, in the face of either positive or

relaxed selection, a population may accumulate multiple independent changes at effector loci, in *cis*-regulatory sequences or in protein-coding regions, avoiding the potentially pleiotropic effects of variation in master regulators (1, 17).

In this work, we set out to test the polygenic evolutionary model in a simple eukaryotic system. We developed genome-scale sequencing methods to identify *cis*-regulatory variants between two *Saccharomyces* species, measuring allele-specific expression levels of a gene in a hybrid diploid (18, 19), and we used the results to uncover evidence of directional evolution of the expression of genes in pathways. To extend these findings, we independently tested for evidence of polygenic evolution using population resequencing data.

## Results

**Differential Allele-Specific Expression in an Interspecific Hybrid.** We sought to study directional evolution of *cis*-regulatory control of expression on a genomic scale in *Saccharomyces cerevisiae* and *S. bayanus*. As an experimental paradigm (20–22), we used genome-wide analysis of expression in a hybrid diploid formed from the mating of the two species. Because the orthologs of a given gene from both species are present in the same environment of *trans*-acting factors, differential expression of the orthologs reflects the presence of *cis*-acting differences between the species. Detecting expression variation within the hybrid requires techniques that recognize the alleles encoded by each ortholog in each species. For this purpose, we measured allele-specific expression with Illumina short-read sequencing (RNA-seq) (23–25). We mated the W303 strain of *S. cerevisiae* and the CBS 4001 isolate of *S. bayanus* to form a hybrid diploid and isolated RNA from two independent cultures of this strain as biological replicates. For each, we sequenced cDNA libraries, for a total of eight lanes. We used the set of sequencing reads mapping uniquely to either the *S. cerevisiae* or the *S. bayanus* genome to quantitate the expression of the allele of each gene from each species, yielding 4,238 ortholog pairs with observable allele-specific expression.

For a given ortholog pair, we used the average log-ratio of *S. bayanus* to *S. cerevisiae* per-base read counts as a statistic for differential allele-specific expression. However, such a measure reflects not only the biologically relevant difference in transcript abundance but also the inherent “sequenceability” of one allele relative to the other. In particular, GC-rich regions tend to be preferentially sequenced by the Illumina platform (26, 27), raising the possibility that comparison of raw allele-specific read counts from an ortholog pair might overestimate the abundance

Author contributions: R.B.B. designed research; R.B.B., J.H.B., and Y.M. performed research; J.H.B. and S.D. contributed new reagents/analytic tools; J.H.B., Y.M., S.D., and R.B.B. analyzed data; and J.H.B., S.D., and R.B.B. wrote the paper.

The authors declare no conflict of interest.

\*This Direct Submission article had a prearranged editor.

Freely available online through the PNAS open access option.

Data deposition: GEO data set GSE19837.

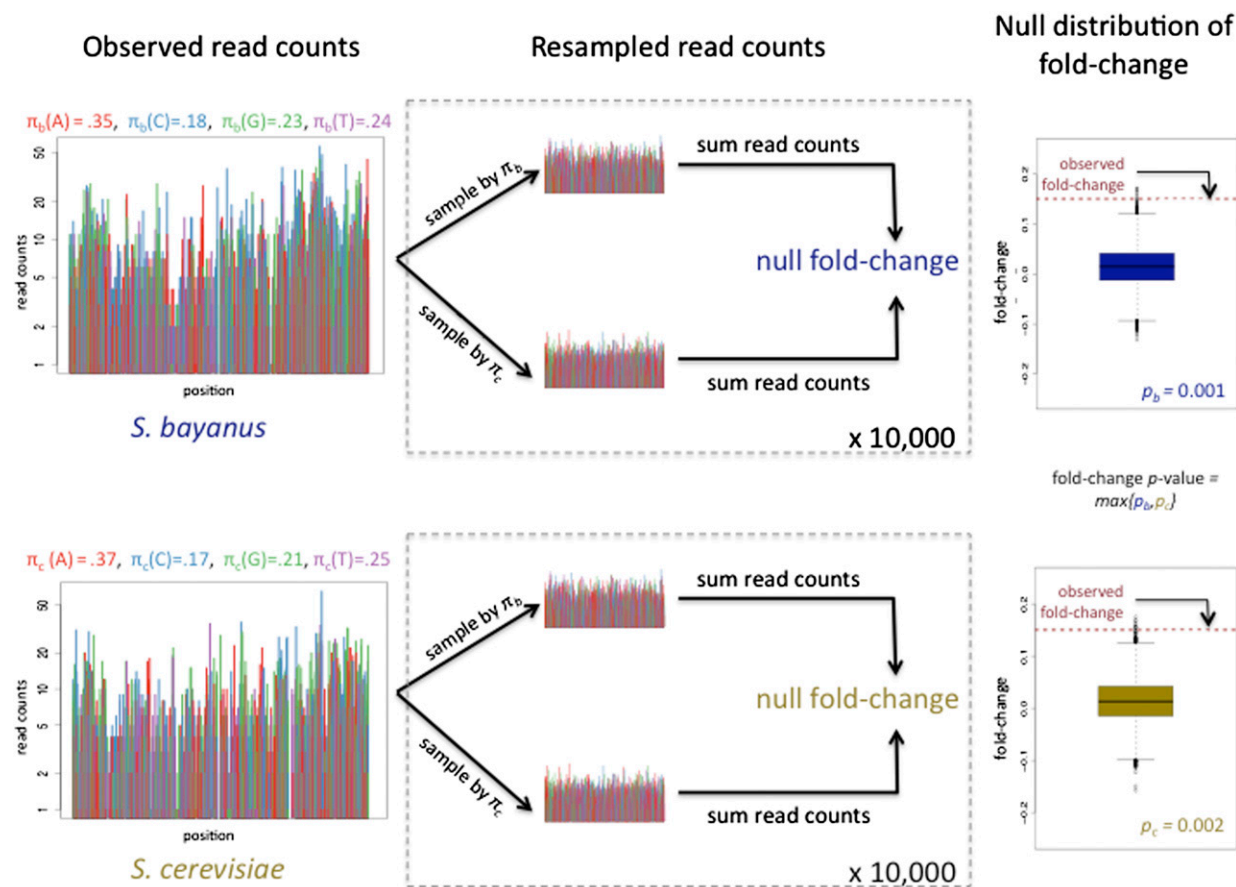
<sup>1</sup>To whom correspondence should be addressed. E-mail: rbrem@berkeley.edu.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0912959107/DCSupplemental](http://www.pnas.org/cgi/content/full/0912959107/DCSupplemental).

of the more GC-rich ortholog. We investigated the impact of allelic base composition on our RNA-seq read counts (Fig. S1) and found that, roughly, a 5% difference in GC content between alleles was associated with a 10% fold-change in their apparent expression. To account for such confounding, we developed a resampling-based method for analyzing differential allele-specific expression. In our approach, the significance of differences in read counts is evaluated by reference to a null distribution that incorporates differences in sequence composition. As illustrated in Fig. 1, we resampled the base-level read counts of the two orthologs to form “null” ortholog pairs with no differential expression and the same marginal nucleotide distributions as the original orthologs. In a given null pair, any apparent difference in the number of RNA-seq reads mapping to one ortholog rather than the other can be attributed to sequencing artifacts alone. When analyzing our actual data, to assert that the observed difference in RNA-seq reads between orthologs in a pair arose from variation in expression between the species and not from technical effects, we required that the observed RNA-seq fold-change lie outside the range of differences we expected under the null. This comparison between observed and null data resulted in a  $P$  value for each ortholog pair. The resampling strategy was stringent, identifying fewer cases of significant differential expres-

sion than did a standard statistical approach that does not account for GC bias (Table S1).

We applied the resampling method to each of the two biological replicates and identified ortholog pairs for which both  $P$  values fell below a given cutoff. The results, given in Table S2, yielded 2,124 ortholog pairs with maximum  $P$  value  $< 0.05$  (at which 212 genes would be expected under the null from sequenceability effects alone), 1,570 ortholog pairs with  $P < 0.005$  (21 would be expected under the null), and 1,176 ortholog pairs with  $P < 0.0005$  (2 would be expected under the null). To evaluate the resulting dataset, we designed quantitative RT-PCR assays to measure allele-specific expression for a subset of genes that spanned the range of  $P$  values and fold-changes from RNA-seq from among those significant at  $P < 0.05$  (Table S3). As shown in Fig. S2, we observed good agreement between RT-PCR and RNA-seq measures of allelic expression differences in the hybrid. Of immediate relevance for our downstream analysis (see below), the two methods agreed on the sign of the fold-change in 20/22 genes tested, confirming the ability of our RNA-seq procedure to determine for a given ortholog pair whether the *S. bayanus* or *S. cerevisiae* allele was associated with higher expression.



**Fig. 1.** Schematic of RNA-seq method for inferring differential allele-specific expression in a hybrid diploid; resampling procedure for an example gene. (Left) Observed base-level read counts are displayed for the *S. bayanus* and *S. cerevisiae* orthologs, using color to represent the nucleotide at each base. The  $x$  axis gives the position of each base, and the  $y$  axis gives the number of allele-specific reads whose first nucleotide maps to a given position. Above each plot are the allele-specific marginal nucleotide frequencies  $\pi_b = [\pi_b(A), \pi_b(C), \pi_b(G), \pi_b(T)]$  and  $\pi_c = [\pi_c(A), \pi_c(C), \pi_c(G), \pi_c(T)]$ , for *S. bayanus* and *S. cerevisiae*, respectively. (Center) For each ortholog, “null” counts are created by resampling base-level read counts according to the *S. bayanus* and *S. cerevisiae* nucleotide frequencies  $\pi_b$  and  $\pi_c$ , respectively. Null expression log-fold-changes are computed by averaging (across lanes for each of the two biological replicates) log-ratios of *S. bayanus* to *S. cerevisiae* null per-base read counts. The resampling procedure is repeated 10,000 times for each ortholog. (Right) Boxplots of the 10,000 null expression log-fold-changes for each ortholog. The observed log-fold-change from the original read counts is represented by dark red dashed lines and is compared with each null distribution to obtain two-sided  $P$  values,  $p_b$  and  $p_c$ . The significance of the observed allele-specific expression difference is summarized by the maximum  $P$  value,  $\max\{p_b, p_c\}$ .



**Table 1. Directional evolution of gene expression in pathways**

| Name                                                                                                                | Statistic* | # Genes <sup>†</sup> | P <sup>‡</sup> | Annotation                                   |
|---------------------------------------------------------------------------------------------------------------------|------------|----------------------|----------------|----------------------------------------------|
| A. Gene groups defined by coregulation in <i>S. cerevisiae</i> (14); 1.08 false-positive groups expected            |            |                      |                |                                              |
| Cluster_Histidine                                                                                                   | -7         | 8                    | 9.00E-05       | Histidine biosynthesis                       |
| Node 73                                                                                                             | 71         | 349                  | 0.001          | Ribosome and translation                     |
| Cluster_NRG1                                                                                                        | -7         | 12                   | 0.002          | Stress-induced transport                     |
| Cluster_Lysine                                                                                                      | -5         | 8                    | 0.006          | Lysine biosynthesis                          |
| Node 45                                                                                                             | -9         | 45                   | 0.007          | Respiration                                  |
| B. Gene groups defined as the Biological Process categories from Gene Ontology; 0.53 false-positive groups expected |            |                      |                |                                              |
| GO:0422540                                                                                                          | 79         | 269                  | 1.0E-05        | Ribosome biogenesis                          |
| GO:0016070                                                                                                          | 124        | 827                  | 0.009          | RNA metabolic process                        |
| GO:0006275                                                                                                          | -8         | 56                   | 0.01           | Cellular aromatic compound metabolic process |

Listed are the top-scoring gene groups with sign imbalance in *cis*-regulatory variation. Results for all groups are given in Table S8.

\*Directional differential expression statistic measuring the imbalance in the signs of *cis*-regulatory variants within gene groups, defined as the sum of the directional scores of the group members (1 for a gene significantly up-regulated in *S. bayanus* for both biological replicates, -1 for a gene significantly up-regulated in *S. cerevisiae*, and 0 otherwise).

<sup>†</sup>Total number of genes in each group for which a directional score was available.

<sup>‡</sup>Resampling-based *P* value assessing the significance of directional *cis*-regulatory expression variation in each group.

high-scoring genes in McDonald–Kreitman-like tests. However, in McDonald–Kreitman analysis of coding regions, at nominal *P* < 0.005 we detected 94 genes (12.6 expected if all genes were evolving neutrally), including three genes with an excess of nonsynonymous changes between species: the fatty acid oxidative enzyme *TES1*, the transcription factor *YRMI*, and the cell wall component *CCW14*.

Excluding the latter candidate cases of positive selection from further analysis, we next sought to identify pathways enriched for genes with excess nonsynonymous polymorphism. Because analysis of Gene Ontology categories gave modest but appreciable power (Table S5), we focused on this class of gene groups, which revealed that the accumulation of coding changes among *S. cerevisiae* isolates can be pathway-dependent. Among the top-scoring results, shown in Table 2, were categories annotated in stress response, cell division, and DNA metabolism, each enriched for genes with an excess of protein-coding variants within *S. cerevisiae*. In a given pathway, such enrichment is consistent with a history of relaxed constraint or balancing selection across the species, or with pathway-level divergence between *S. cerevisiae* populations. Given the known preponderance of alleles segregating at low frequency within *S. cerevisiae* (30), we repeated McDonald–Kreitman tests and pathway analyses on a dataset in which singleton sites were filtered out (Tables S6 and S7). Top-scoring pathways from these filtered data mirrored pathways emerging from analysis of all sites, albeit

with weaker *P* values for enrichment of polymorphism (Table S7); thus, although singleton alleles drive much of the statistical power for this analysis, more common variants likely exhibit similar trends of enrichment in pathways. We conclude that an excess of fixed sequence differences between *S. cerevisiae* and *S. bayanus*, indicative of positive selection between the species, can be detected at a handful of individual loci and that genes of common function share patterns of sequence polymorphism within *S. cerevisiae*.

## Discussion

Decades of work based on the study of individual genes with signatures of selection have shed light on the mechanisms of evolutionary change (1–3). However, any such locus may represent only one component of a suite of changes among genes of related function that have arisen in response to the same selective pressure over evolutionary time. The prevalence and mechanisms of polygenic evolution are unknown, and addressing the question requires genome-scale data sources and analytic methods that in many cases remain to be developed. For genomic analysis of the evolution of gene expression, *cis*-acting differences in a heterozygote diploid can be mapped at high resolution (20–22), but to date, methods for this purpose have required custom microarray tools. The RNA-seq platform we pioneer here enables study of *cis*-acting variation in expression in any organism and any strain (23–25). Likewise, we have demonstrated the power of extending the classical sequence-based McDonald–Kreitman test for selection to pathways, an approach that has increasing utility with the advent of population genomic sequencing data (5–8, 32).

We have developed an analytic strategy to detect nonneutral evolution in pathways, finding *cis*-regulatory changes with effects in the same direction that have accumulated in one phylogenetic lineage relative to another. Any case of directional evolution of gene expression in a pathway can be explained with either of two evolutionary models. A lineage can accumulate novel alleles that predominantly up-regulate or predominantly down-regulate the genes of a pathway because these changes result in a fitness advantage in a particular niche. Alternatively, coherent purifying selection can keep levels of pathway genes predominantly low or predominantly high in a particular lineage; if this selective force is relaxed in a second lineage, novel alleles accumulated by drift will change expression levels in a direction opposing the original selective constraint. Tests for directional selection have their origin in analyses of quantitative trait loci mapped to macroscopic phenotypes (28), but studies of a given trait rarely have sufficient power for the test. By

**Table 2. Pathway enrichment of within-species sequence variation**

| Name       | MK enrichment* | P <sup>†</sup>       | Annotation            |
|------------|----------------|----------------------|-----------------------|
| GO:0007124 | 6.0            | 9 × 10 <sup>-5</sup> | Pseudohyphal growth   |
| GO:0006950 | 2.2            | 0.0005               | Response to stress    |
| GO:0006259 | 2.1            | 0.006                | DNA metabolic process |
| GO:0007126 | 2.5            | 0.03                 | Meiosis               |

Listed are the top-scoring pathways in a test for enrichment among genes with evidence for protein coding variation within *S. cerevisiae* according to the McDonald–Kreitman test. Results for all pathways are given in Table S5. Gene groups were defined as the Biological Process categories from Gene Ontology, with 1.2 false-positive groups expected.

\*The ratio between the number of genes in the pathway scored as significant in the McDonald–Kreitman test and the number of genes in all pathways scored as significant, divided by the analogous ratio for all tested genes regardless of significance.

<sup>†</sup>Hypergeometric *P* value for pathway enrichment of genes significant in the McDonald–Kreitman test.

contrast, genome-scale datasets of *cis*-regulatory variants of the kind reported here represent the result of hundreds to thousands of independent genetic events in the history of the two parent genomes. Harnessing these data, our test for coherent regulatory change complements previous strategies to observe the evolution of regulatory motifs in pathways (13–16), providing a significance estimate and a direction of inferred effects on pathway output for each candidate case of nonneutral evolution. Interestingly, among the hits were a number of metabolic pathways, including respiration and amino acid biosynthesis genes. Although the latter could be the result of adaptation to laboratory conditions, it is tempting to speculate that changes in nutrient availability may underlie cases of pathway evolution in Saccharomycetes, as would be expected from findings in other organisms (33–35).

However, our test for directional regulatory evolution is not designed to detect cases of selection on amino acid changes, and it focuses on the divergence between a pair of isolates rather than divergence among many strains. Indeed, sequence-based tests detected cases of polygenic evolution in a set of pathways different from those identified by the expression-based approach, as expected given the distinctions between the two methods. In particular, McDonald–Kreitman tests revealed patterns of protein-coding polymorphism at the pathway level between *S. cerevisiae* isolates. Previous analyses of rare alleles in *S. cerevisiae* have suggested relaxed selection as a major evolutionary force in this system (30), and relaxation resulting from environmental change would be expected to drive an accumulation of derived alleles in particular response pathways. Thus, pathways enriched for nonsynonymous polymorphism across *S. cerevisiae* may shed light on the changes in the life history of this species since its divergence from *S. bayanus*. In addition, some cases of pathway-level polymorphism may result from divergence between *S. cerevisiae* populations, as the result of adaptation (36, 37) or of relaxed selection in particular niches. Improved sequence sampling of the well-defined populations in this species (30) ultimately will enable in-depth study of the changes in selective pressure at play during their divergence.

Emerging from our work and that of others (5–7, 13–16, 32) is a picture of evolutionary change within and between species mediated by polygenic suites of variants in pathways. With respect to adaptation, a more parsimonious mechanism could involve a single advantageous change in an upstream regulator that up- or down-regulates a pathway in one mutational step (2). However, in experimental organisms, changes in *trans*-acting transcriptional regulators are more prevalent in lines that have accumulated mutations than in wild isolates (12) and are more common within species than between species (17), suggesting that upstream effects are often deleterious and subject to purifying selection. As such, the time taken by a mutational search for an adaptive variant in an upstream regulator often may be comparable to that required for the acquisition of multiple changes in downstream effectors (1). By the same token, relaxed selection on the function of a pathway frequently would give rise to suites of downstream changes but less often to a major upstream variant. The prevalence of polygenic evolution may explain why detecting evidence for strong selection on single loci has proven challenging to date and will serve as continued motivation for the study of mechanisms by which genes in pathways coevolve.

## Materials and Methods

**Strains, RNA Preparation, and Sequencing.** Strain OZY27, a hybrid of *S. cerevisiae* (isogenic to W303; *his3 leu2 lys2 trp1 ura3*) and *S. bayanus* (derived from CBS 4001; *ade2 his3 lys2*) (38), was kindly provided by O. Zill, University of California, Berkeley. For each of two biological replicates, an OZY27 culture was grown to log phase at 25 °C in YPD medium (39). Total RNA was isolated by the hot acid phenol method (39) and treated with Turbo DNase (Ambion) according to the manufacturer's instructions. A Solexa/Illumina 1G Genome Analyzer library was constructed for each biological replicate as described (40). Each biological replicate was sequenced at both the Vincent J. Coates Genomic Sequencing Laboratory at the University

of California, Berkeley, and the Core Instrumentation Facility at the University of California, Riverside. The former provided two lanes of 37-bp reads for each biological replicate; the latter provided three lanes of 40-bp reads for biological replicate 1 and one lane of 40-bp reads for replicate 2. All resulting reads were trimmed to 30 bp.

**RNA-seq Preprocessing and Mapping.** We downloaded 4,853 ortholog pairs (<http://www.broad.mit.edu/reggev/orthogroups/orthologs.html>) and in each pair converted the *S. cerevisiae* ortholog sequence to that of W303 using sequence from (30). For each lane, 30-bp reads were mapped to either strand of the ortholog pair set using Bowtie (41), allowing no mismatches between the read and the reference sequence. All further analyses were done in R (42). For each position in the combined orfeome, we determined whether the read starting at that location was uniquely-mappable (2), i.e., had an edit distance of 2 or greater to any other position in the combined (double-stranded) reference orfeome or genome. The combined orfeome consists of all 4,853 *S. bayanus* and *S. cerevisiae* ortholog pairs; the combined genome consists of the concatenation of the *S. bayanus* ([http://downloads.yeastgenome.org/sequence/fungal\\_genomes/S\\_bayanus/WashU/contig/](http://downloads.yeastgenome.org/sequence/fungal_genomes/S_bayanus/WashU/contig/)) and *S. cerevisiae* ([ftp://ftp.yeastgenome.org/yeast/data\\_download/sequence/NCBI\\_genome\\_source/](ftp://ftp.yeastgenome.org/yeast/data_download/sequence/NCBI_genome_source/)) reference genomes, converting the latter to that of the W303 strain as above. The ensuing analysis considered only uniquely-mappable (2) reads and the 4,238 ortholog pairs with a minimum of 200 uniquely mappable (2) bases in both species and a maximum difference in length of 100 bp. The correlation in read counts was greater between replicate lanes, sequencing centers, and biological replicates than between *S. bayanus* and *S. cerevisiae* orthologs (Fig. S3).

**Identifying Ortholog Pairs with Differential Allele-Specific Expression.** Our method for detecting significant differential expression between orthologs in a pair is described in detail in *SI Text*. Briefly, for a given ortholog pair we sampled the base-level read counts of the *S. bayanus* ortholog to create null *S. bayanus* data; to create null *S. cerevisiae* counts, we resampled, from the *S. bayanus* ortholog, base-level read counts according to the nucleotide frequencies in the *S. cerevisiae* ortholog. Computing differential expression statistics for each of 10,000 such null orthologs yielded a null distribution to which we compared the observed differential expression statistic to obtain a *P* value. Repeating the procedure with the base-level read counts of the *S. cerevisiae* ortholog yielded a second *P* value, and we conservatively retained the maximum of the two.

### Testing for Directional Imbalance in *cis*-Regulatory Effects Across Pathways.

Our test for directional regulatory evolution is described in detail in *SI Text*. Briefly, for pathway analyses we used the regulons defined in (14) and GO\_slim Biological Process categories ([http://downloads.yeastgenome.org/literature\\_curation/go\\_slim\\_mapping.tab](http://downloads.yeastgenome.org/literature_curation/go_slim_mapping.tab)). For a given pathway, we added directional differential expression scores across all genes in the group (1 if the expression in the *S. bayanus* allele was significantly higher than in the *S. cerevisiae* allele, –1 if the expression in the *S. cerevisiae* allele was significantly higher than in the *S. bayanus* allele, and 0 otherwise) and obtained a *P* value by comparing the resulting statistic with a null distribution of 100,000 such statistics obtained by resampling genes at random, with replacement. At a given *P* value threshold  $P_0$ , we estimated the expected number of false-positive pathways as the product of  $P_0$  and the number of tests; we set  $P_0$  to attain as close to one false-positive pathway as possible. For comparison, results from the Benjamini–Hochberg multiple testing procedure (43), which controls the false-discovery rate, are given in Table S8.

**McDonald–Kreitman Analyses.** We downloaded *S. cerevisiae* strain sequences from (30), eliminating the laboratory strains W303, S288C, and SK1 from further analysis. For each gene, we aligned the amino acid sequences of *S. cerevisiae* strains and the reference sequence of *S. bayanus* ([http://downloads.yeastgenome.org/sequence/fungal\\_genomes/S\\_bayanus/other](http://downloads.yeastgenome.org/sequence/fungal_genomes/S_bayanus/other)) using MUSCLE (44) and regenerated DNA alignments with tralign (<http://emboss.sourceforge.net/>). For a total of 2,511 genes, we applied the McDonald–Kreitman test using software from (32). For upstream regions, we used alignments from (30) for the region starting 200 bp upstream of each gene and ending at coding start, and tabulated polymorphic and divergent sites; for polymorphic and divergent silent sites, we used ORF alignments of nonlaboratory strains from (30). Given the  $2 \times 2$  table from these data, for a total of 1,910 genes, we calculated single-gene McDonald–Kreitman-like *P* values using Fisher's exact test. For filtering of singletons in both coding and upstream analyses, we identified polymorphic sites at which exactly one *S. cerevisiae* strain harbored the minor allele and eliminated these sites from the analysis, using  $P < 0.005$  as a cutoff for significance in the McDonald–Kreitman test.

For pathway analyses, we used 38 GO\_slim categories as above, as well as gene regulons from (14), filtered as described in *SI Text*; to minimize the number of the latter with poor sampling, we considered a regulon only if  $\geq 10$

