

Research

Evolution of transcription factor binding through sequence variations and turnover of binding sites

Gat Krieger,^{1,3} Offir Lupo,^{1,3} Patricia Wittkopp,² and Naama Barkai¹

¹Department of Molecular Genetics, Weizmann Institute of Science, Rehovot 76100, Israel; ²Department of Ecology and Evolutionary Biology, Department of Molecular, Cellular, and Developmental Biology, University of Michigan, Ann Arbor, Michigan 48109, USA

Variations in noncoding regulatory sequences play a central role in evolution. Interpreting such variations, however, remains difficult even in the context of defined attributes such as transcription factor (TF) binding sites. Here, we systematically link variations in *cis*-regulatory sequences to TF binding by profiling the allele-specific binding of 27 TFs expressed in a yeast hybrid, in which two related genomes are present within the same nucleus. TFs localize preferentially to sites containing their known consensus motifs but occupy only a small fraction of the motif-containing sites available within the genomes. Differential binding of TFs to the orthologous alleles was well explained by variations that alter motif sequence, whereas differences in chromatin accessibility between alleles were of little apparent effect. Motif variations that abolished binding when present in only one allele were still bound when present in both alleles, suggesting evolutionary compensation, with a potential role for sequence conservation at the motif's vicinity. At the level of the full promoter, we identify cases of binding-site turnover, in which binding sites are reciprocally gained and lost, yet most interspecific differences remained uncompensated. Our results show the flexibility of TFs to bind imprecise motifs and the fast evolution of TF binding sites between related species.

[Supplemental material is available for this article.]

Changes in gene expression play a key role in cellular adaptation, physiology, and development. Guiding these changes are transcription factors (TFs) that bind DNA at sequence motifs allowing activation or repression of gene transcription. Understanding how TF binding diverges between species is therefore central for understanding how gene regulation evolves.

TFs contain DNA-binding domains (DBDs) that bind with high affinity to short DNA sequence motifs (typically 6–12 base pairs). Sequence variations leading to the emergence or disappearance of binding motifs may therefore drive regulatory divergence by changing TF binding. Previous studies examined for such functional variations by comparing TF binding between related species (Borneman et al. 2007; Wilson et al. 2008; Bradley et al. 2010; Schmidt et al. 2010; Paris et al. 2013; Stefflova et al. 2013), between human individuals (Kasowski et al. 2010; Kilpinen et al. 2013; Maurano et al. 2015), or between alleles of heterozygous cells (Reddy et al. 2012). It was proven difficult, however, to relate the measured changes in TF binding with variations in motif sequence. In their analysis of allele-specific binding of 25 human TFs, Reddy et al. (2012) concluded that only 12% of differentially bound sites were associated with variations in known binding sequences. Similarly, studies comparing binding of six TFs between two *Drosophila* species revealed only modest correlation between interspecific differences in binding and sequence variations in known motifs (Bradley et al. 2010).

The difficulty of associating interspecies differences in TF binding with variations in *cis*-regulatory sequences mirrors the difficulty in predicting TF binding sites. Indeed, motif preference remains a poor indicator for TF binding *in vivo*, primarily because TFs typically bind at only a small subset of motif-con-

taining sites found in genomes. TF binding could therefore evolve in *cis* not only through the emergence or disappearance of binding motifs, but also through variations in DNA accessibility. Examples for such *cis* variations include changes in nucleosome positioning (Mirny 2010; Sun et al. 2015), variations affecting binding of a cooperating TF (Stefflova et al. 2013; Avsec et al. 2021), or variations in promoter regions surrounding the motif, perhaps recognized by TF regions outside the DBD (Brodsky et al. 2020).

In this work, we systematically associated variations in known TF binding motifs to changes in TF binding by mapping allele-specific binding of 27 TFs within an interspecific yeast hybrid. The hybrid's nucleus contains two related parental genomes. By applying allele-specific mapping, we could directly compare TF binding to the two genomes while ensuring a uniform *trans*-regulatory environment (Tirosh et al. 2009; Emerson et al. 2010; Metzger et al. 2016; Wong et al. 2017; Krieger et al. 2020; Floc'hlay et al. 2021; Hill et al. 2021; Lupo et al. 2021; Yang et al. 2021). Our analysis examined the contribution of two types of variations, in sequence motifs and in chromatin accessibility, to divergence of TF binding at individual binding sites. We further examined evolutionary changes of TF binding at the full promoter level by distinguishing cases of compensated binding-site turnover (in which loss of a binding site is compensated by gain of an adjacent binding site) from cases of an uncompensated gain/loss. Finally, by capitalizing on the hundreds of sequence variations in motif-containing sites between the genomes, we defined the effective cost of each binding site mutation *in vivo*, linking this cost with sequence conservation at the motif's vicinity. Our results

³These authors contributed equally to this work.

Corresponding author: naama.barkai@weizmann.ac.il

Article published online before print. Article, supplemental material, and publication date are at <https://www.genome.org/cgi/doi/10.1101/gr.276715.122>.

© 2022 Krieger et al. This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

highlight key aspects in the evolution of TF binding between closely related yeast species.

Results

Mapping allele-specific binding of 27 transcription factors within an interspecies hybrid

To examine systematically the effect of *cis* variation on TF binding, we generated F₁ hybrids by mating two closely related budding yeast species: *Saccharomyces cerevisiae* and *Saccharomyces paradoxus*. These two species diverged approximately five million years ago and largely retained gene identity and synteny. Sequence identity reaches ~90% in coding regions and ~75% in promoters (Scannell et al. 2011; Yue et al. 2017). Both species' genomes are highly compact with short intergenic regions (200–400 bp) that function primarily as gene promoters. We and others have previously used this hybrid as a model for studying the principles of regulatory evolution (Tirosh et al. 2009; Emerson et al. 2010; Artieri and Fraser 2014; Metzger et al. 2016; Weiss et al. 2018; Krieger et al. 2020; Lupo et al. 2021).

We selected 27 TFs of five protein families (Supplemental Table S1). All selected TFs are of known function, and their motif preferences were previously described through in vitro and in vivo experiments (Sandelin et al. 2004; De Boer and Hughes

2012). We mapped the localization of these TFs along the orthologous hybrid genomes using chromatin endonuclease cleavage followed by sequencing (ChEC-seq) (Zentner et al. 2015). For this, each TF was fused to a MNase, allowing us to trigger DNA cleavage at the close vicinity of the TF binding site using a short (30 sec) Ca²⁺ pulse. The short DNA fragments were extracted and sequenced. We found this method to give a highly reproducible and spatially resolved TF binding maps (Bar-Ziv et al. 2020; Brodsky et al. 2020; Gera et al. 2021; Lupo et al. 2021). We previously observed that orthologous TF proteins bind to similar locations in the hybrid genome, and at a similar level, by profiling both the *S. cerevisiae* TF ortholog and the *S. paradoxus* ortholog in separate experiments (Lupo et al. 2021). This observation was consistent with the generally slow evolution of TF preferences (Carroll 2005). We therefore profiled only the *S. cerevisiae* ortholog, examining how its binding differs between the two alleles (Fig. 1A). Notably, comparing our binding profile with six published data sets revealed high consistency in promoter binding, peak binding, and preferred motifs (Supplemental Note 1; Supplemental Figs. S16, S17; Supplemental Tables S4, S5). Binding signals defined by ChEC-seq were largely restricted to promoter regions, as expected (Supplemental Fig. S1A).

Next, we compared TF binding between the two orthologs. For this, we distinguished first the overall signal obtained

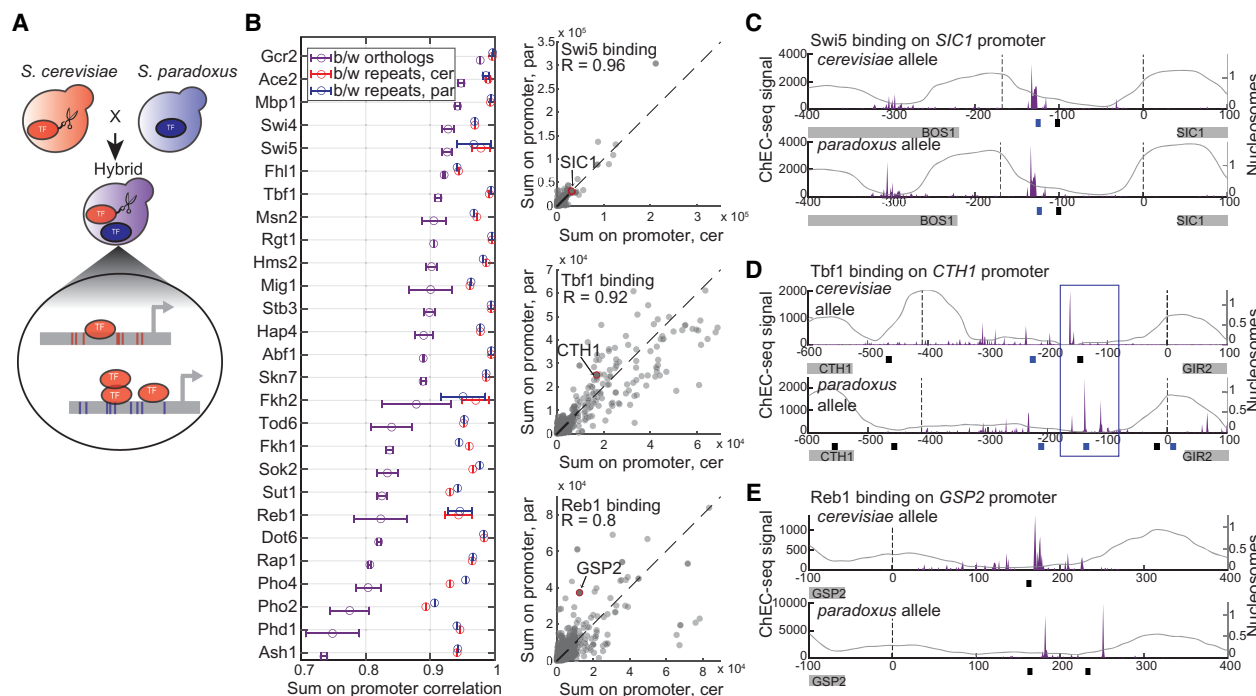


Figure 1. Experimental system to profile *cis* variation in transcription factor binding. (A) Scheme of the experimental system. *S. cerevisiae* strains where a TF was fused to a MNase (illustrated as scissors) were crossed with a WT *S. paradoxus* strain to form a hybrid, on which the ChEC-seq method was applied to profile in vivo TF binding. Orthologous promoters harbor sequence variations (red and blue lines) and differential binding levels. (B) Global similarity in orthologous binding of the 27 TFs examined here. Presented is the Pearson correlation coefficient of sum of signal on all yeast promoters (6701 promoters), between experimental replicates (*S. cerevisiae* genome in red, *S. paradoxus* genome in blue), and between orthologous genomes (purple). Data are the mean and standard deviation of two to five repeats. (Right) Promoter binding of three TFs; each point is the sum of signal on a specific gene promoter. (C–E) Examples for TF binding to orthologous promoters. (C) Conservation of Swi5 binding to *SIC1* promoter, *S. cerevisiae* ortholog (upper) and *S. paradoxus* ortholog (lower). ChEC-seq signal is the 5' end of reads, presented in purple. Nucleosome occupancy data of the hybrid (Tirosh et al. 2010) are presented as gray lines. Transcription start sites are presented in gray dashed lines (Pelechano et al. 2013; Park et al. 2014). For Swi5, CCAGC motif sequences are marked in blue (plus strand) and black (minus strand) boxes. ORFs are presented as gray boxes. (D) Binding-site turnover of Tbf1 to *GIR2/CTH1* promoter. The blue box marks the region of binding-site turnover, in which the Tbf1 motif appears on the plus strand in the *S. cerevisiae* allele (ACCTA), and the same motif realization appears on the minus strand in the *S. paradoxus* allele (TAGGT), where motif sequences partially overlap. Consensus motif of Tbf1 is [C/A]CCTA. (E) Divergence in Reb1 binding to *GSP2* promoter. Annotation as in D; Reb1 consensus motif is TTACCC[G/T].

throughout each promoter, and second, the locations of individual binding sites within promoters. Focusing first on the level of gene promoters, we find high conservation: in all TFs, promoter binding pattern was correlated between the two alleles (Pearson correlation coefficient, R , ranging from 0.75 to 0.96). Correlation between experimental replicates was significantly higher (R ranging from 0.9 to 0.99), supporting the reproducibility of our data and suggesting that some allelic differences do exist (Fig. 1B). Correlation between different TFs was much lower (average $R=0.12$) (Supplemental Fig. S1B). Other measures of correlation were also examined (Supplemental Note 3, Supplemental Fig. S19); however, we found the Pearson linear correlation coefficient to be most appropriate because most TFs bind a small number of targets.

Examining individual promoters, we noted various patterns of conservation and divergence (Fig. 1). In the case of Swi5, for example, promoter binding profiles were highly similar between orthologs ($R=0.96$), and this similarity extended when examining binding peaks at highly bound promoters (e.g., *SIC1*) (Fig. 1C; experimental replicates in Supplemental Fig. S1C–E). In other cases, overall promoter binding was conserved between the two alleles, yet the distribution of binding peaks varied along specific promoters, implying on binding-site turnover (Ludwig et al. 2000; Moses et al. 2006; Dermitzakis and Clark 2009). For example, Tbf1 showed similar overall binding to the *CTH1/GIR2* promoter in the two orthologous alleles, yet the precise binding pattern differed, and this variation was linked to a change in the location of the Tbf1 motif within the two orthologous promoters (Fig. 1D, blue box). We also observed cases of divergence in which overall binding differed between the two orthologous promoters, as exemplified by the binding of Reb1 to the *GSP2* promoter, which was significantly stronger at the *S. paradoxus* allele. Also here, differential binding correlated with the presence of an additional Reb1 motif in *S. paradoxus* allele but not in the *S. cerevisiae* allele (Fig. 1E). We conclude that although TF binding remains largely invariant at the resolution of the full promoter, cases of *cis* divergence at the level of individual TF binding sites are readily identified.

Transcription factors bind a select subset of motif-containing sites within the two genomes

Promoter regions in budding yeast are typically 200–400 bp long (Supplemental Fig. S1F; Kristiansson et al. 2009), whereas individual binding sites contain only 6–12 bp. To examine whether our data can define TF binding at a resolution that is compatible with individual binding sites, we first observed the binding signal around motif-containing sites (Fig. 2A), referring to the known in vitro motif of each TF as curated in the YetFasco (De Boer and Hughes 2012) and JASPAR (Sandelin et al. 2004) databases. For most TFs, these in vitro defined motifs agreed well with de novo motifs defined from our data by either enrichment of 7-mer sequences around bound sites or the MEME-ChIP algorithm (Supplemental Table S1; Supplemental Note 1; Machanick and Bailey 2011).

Considering first the Reb1 TF, we find binding signal at locations containing its in vitro motif, as expected. Binding, however, was restricted to only ~30% of motif sites found within promoters (Fig. 2A). In this analysis, we estimated the significance of TF binding relative to a set of random sites within promoters and defined a binding threshold at 95% of random site distribution (Fig. 2B), resulting in 2063 Reb1-bound sites. Binding level at motif sites was moderately correlated ($R=0.25$) with the motif P -value as defined by FIMO (Grant et al. 2011). Consistent with the expected pattern

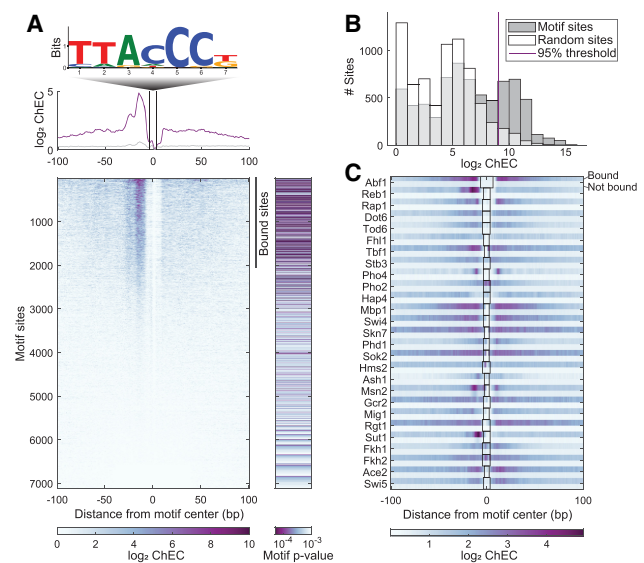


Figure 2. Transcription factors bind a select subset of motif-containing sites. (A) Reb1 binding signal to motif-containing sites (potential binding sites). (Top) In vitro motif of Reb1 (Fordyce et al. 2010). (Middle) Average ChEC-seq signal (5' end of reads) in a logarithmic scale; bound sites presented in purple, nonbound sites presented in gray. (Bottom left) Heatmap of ChEC-seq signal at 7115 sites that contain the Reb1 motif in both hybrid genomes; (bottom right) motif P -value according to FIMO (Grant et al. 2011). (B) Binding level distribution of Reb1 in Reb1 motif sites and in random sites. Binding level threshold was set as the 95% of random site distribution. This threshold defines the bound sites indicated in A. (C) Signal around motif sites of all examined TFs, at bound (top row) and non-bound (bottom row) sites. Boxes indicate motif size. Full profiles are presented in Supplemental Figure S2.

of this method, MNase-cleavage signal peaked at the motif boundaries and was depleted from within the motif itself; the latter is protected from cleavage by the bound TF (Fig. 2A). Results for other TFs were similar, although they varied in details depending on TF identity, such as the cleavage symmetry around the motif and the width of the cleavage-protected region (Fig. 2C; full profiles in Supplemental Fig. S2). These details perhaps reflect differences in TF mobility on the DNA (Suter 2020), motif-specificity, and the size of the protein or protein-complex bound to the DNA. An example for the latter factor may be Hap4, for which the protected area appears significantly larger (30 bases) than the known motif (7 bases) (Fig. 2C), perhaps indicating its binding as a subunit of the larger Hap2/3/4/5 transcriptional activation complex (McNabb 2005). We conclude that ChEC-seq allows mapping of individual binding locations with high resolution.

To map individual binding sites, we used an available peak-calling algorithm (Methods). Peak locations were largely consistent with previously published data (Supplemental Note 1; Supplemental Fig. S16). Notably, a considerable fraction (0.2–0.6) of reads were mapped to peaks (Supplemental Fig. S4). Overall, 28% of the peaks were associated with the known in vitro motif, here referred to as binding sites, in which peaks and motifs are less than 30 bases apart (Supplemental Fig. S5A). The percentage of peaks associated with the known in vitro motifs (all motif realizations with FIMO P -value < 0.001) ranged from 8% to 62% between the different TFs. This fraction was 2.5 times higher than the fraction of random sites that reside next to an in vitro motif (averaging over all TFs), resembling high motif association of peaks in our data. We observe high specificity of TFs to their

binding sites, with low overlap between TFs (2%–4%) (Supplemental Fig. S18D) and no typical binding pattern at binding sites of an unrelated TF (Supplemental Fig. S3).

To estimate the level of systematic noise in our data, namely, binding peaks that are not a result of TF binding, we tested the ChEC-seq profile of an endogenously expressed free-MNase (Supplemental Note 2; Supplemental Fig. S18). Only 2.5%–4% of TF binding peaks overlapped with free-MNase peaks, representing the false-positive rate of the method and agreeing with previous estimations (Zentner et al. 2015). Therefore, the high number of peaks that were not motif-associated were also not bound by a free-MNase and may indicate functional binding events. Such events might result, for example, from recruitment by interacting TF or from protein regions outside the DBD that interact with DNA. Because the sequence basis of these binding events is not characterized, we decided to focus our analysis on binding peaks containing the known *in vitro* motifs, representing 8%–62% of TF-specific peaks (Supplemental Fig. S5A).

Differential TF binding to the two hybrid alleles correlates with variations in motif sequence, but differences in motif accessibility play a minor role

TFs could gain new binding sites through at least two mechanisms (Fig. 3A). First, mutations could change the accessibility of the DNA in regions containing a motif site, for example, through sequence mutations causing a nucleosome-depleted region. Second, new motifs could emerge by mutations within accessible regions. As these two processes occur in parallel, their prevalence may vary depending on the motif type and the processes governing TF specificity, namely, its attraction to only a subset of its motif-containing sites.

To distinguish between these two mechanisms of divergence, we focused on TF binding to motif-containing sites. To enable comparison of orthologous binding sites, we locally aligned orthologous promoters and compared sequence and binding level over the aligned coordinates (Methods). We distinguish between sites where both orthologs contain the corresponding motif (conserved between species; common) and sites where the motif is found in only one species' genome (diverged between species; unique). Overall, 36% of sites were classified as common sites, 36% as *cerevisiae*-unique, and 28% as *paradoxus*-unique (Fig. 3B; Supplemental Fig. S6A). As

expected, the fraction of common sites bound by the respective TFs was, on average, twice as high as that of the unique sites (Supplemental Fig. S6B).

A significant fraction of common sites remained unbound in both genomes, implying that these sites are likely to be inaccessible for TF binding (Fig. 3B; Supplemental Fig. S6C). Indeed, for most TFs (e.g., Reb1), nucleosome occupancy at unbound sites

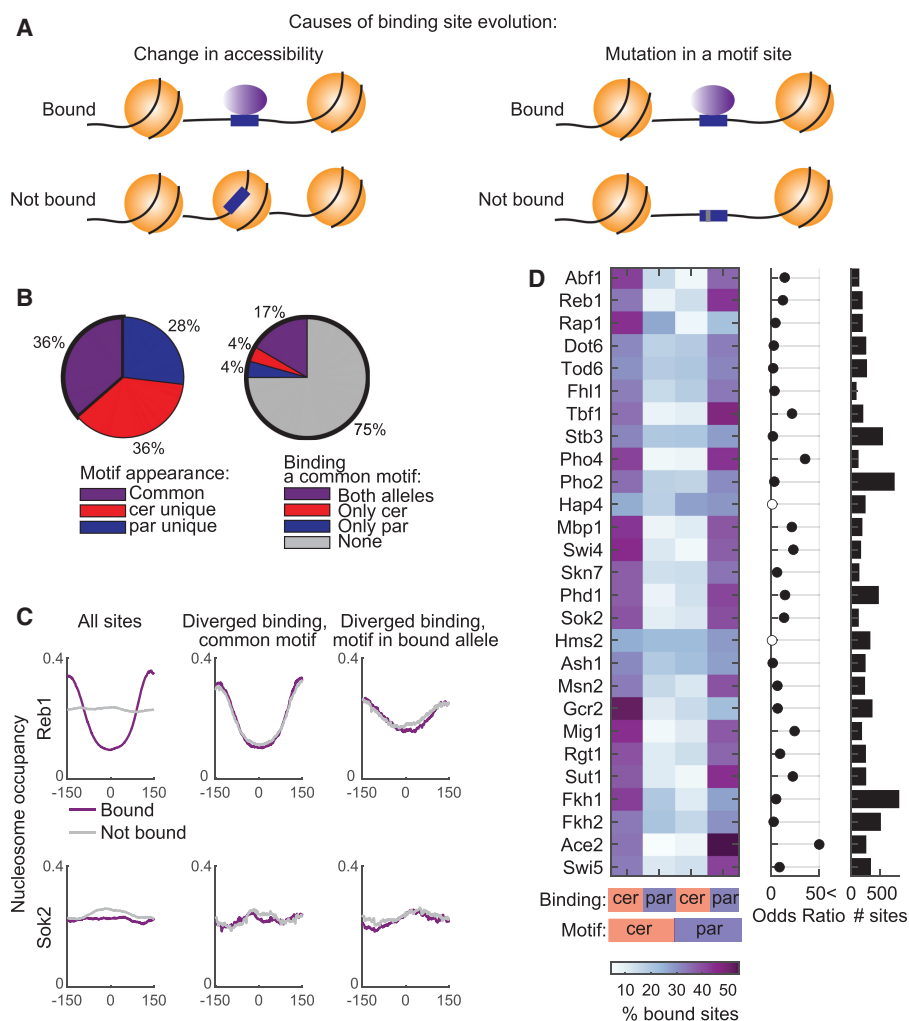


Figure 3. Differential TF binding to the two alleles correlates with variations in motif sequence, but differences in motif accessibility play a minor role. (A) Suggested mechanisms for TF binding evolution: (left) motif site is conserved, but in the unbound allele it is occupied by a nucleosome and therefore not accessible for TF binding; and (right) motif site is lost owing to a single-nucleotide variation. Nucleosomes are illustrated in orange; TF in purple oval; motif site in blue box; nucleotide variation as a gray stripe. (B) Proportion of motif sites (left) and proportion of bound sites among the common motif sites (right). (Left) Proportion of motif sites that are common to both orthologs, and sites that appear only in a certain ortholog (cer- or par-unique) among all *in vitro* defined motif sites of the full set of 27 TFs (62,970 are common, 63,455 cer-unique, 46,440 par-unique) (for TF-specific proportions, see Supplemental Fig. S6). (Right) Proportion of binding to common motif sites. (C) Nucleosome occupancy does not explain differential binding between orthologs. Presented are nucleosome occupancy profiles averaged over motif sites, centered at the binding motif of Reb1 (upper) and Sok2 (lower). (Left) All motif-containing sites, divided into bound sites (purple) and nonbound sites (gray). (Middle) Common motif sites that show diverged binding. Nucleosomes at the bound allele are in purple, and nucleosomes at the nonbound allele are in gray. (Right) Unique motif sites with diverged binding. Nucleosomes at the bound allele, which harbors a motif, are in purple, and nucleosomes at the nonbound and motifless allele are in gray (for profiles of all TFs, see Supplemental Fig. S7). (D) Binding to unique motif sites, with biased binding to the motif-containing allele. (Left) Percent of bound sites. (Middle) Odds ratio of Fisher's exact test; full black circles indicate significant comparisons (P-value < 0.05, FDR corrected). (Right) Number of unique motif sites.

was higher than that at bound ones (Fig. 3C; Supplemental Fig. S7), whereas other TFs (e.g., the stress-related TF Sok2) were preferentially bound at sites that are often nucleosome-occupied, as reported before (Kaplan et al. 2009). On average, most TFs were preferentially bound at sites of low nucleosome occupancy, prompting us to ask whether changes in nucleosome occupancy between alleles could also explain the divergence of binding. For this, we asked whether cases in which both orthologous alleles contain a motif (common motif), yet only one of these alleles is in fact bound, might result from differential nucleosome occupancy of the two alleles. We find that in these cases, and also in cases of unique motif sites, nucleosomes are equally positioned in the bound and nonbound orthologs (Fig. 3C; Supplemental Fig. S7). This suggests that the differences in DNA accessibility, at least as reflected by nucleosome occupancy, play a minor role in the divergence of TF binding preferences.

To examine whether differential TF binding correlates with the emergence or loss of the binding motif, we focused on cases of unique sites in which the motif is present in only one of the alleles. We asked whether, in these cases, the allele containing the motif is more likely to be bound than the one that lacks the motif. This was indeed the case (Fig. 3D): in 25/27 TFs in our set, we observed high correspondence between the allele containing the motif and the one bound by the respective TFs. Together, these data suggest that divergence in TF binding caused by changes in DNA accessibility are less frequent compared to these caused by the emergence or loss of a binding motif.

TF binding to an imperfect motif depends on the genomic context

Our analysis so far focused on cases in which TF binding was lost or gained in one of the genomes. Next, we considered also quantitative changes, in which TFs bound the two alleles but at different levels. Such quantitative differences in the allele-specific TF binding were in fact quite common and accounted for the majority of binding changes (Supplemental Fig. S8). We asked whether these quantitative differences could be explained by sequence variations within the binding motif. For this, we focused on binding peaks that contain the associated motif in at least one of the genomes, and considered the sequence variations within the motif site and in its immediate surroundings. To further focus the analysis, first we considered cases of unique alternative alleles, in which one ortholog has the consensus motif (as defined *in vitro*), whereas the second ortholog has a one-letter variant either within the motif itself or in its flanking region (five bases upstream/downstream from the core motif). Comparing TF binding occupancy at the two orthologs allowed quantifying the average cost (reduction in TF binding) of each deviation from consensus (Fig. 4A).

For Reb1, deviations from consensus in the core motif had a strong impact on binding (Fig. 4A). This sensitivity to deviation from the consensus motif differed between TFs and was further dependent on the position and the precise alternative (Fig. 4C). In fact, some TFs remained largely insensitive to single-letter variations (e.g., Skn7, Gcr2, Stb3), whereas others showed greatly reduced binding (e.g., Rap1, Tbf1, Pho4). In some cases, variations in sequences flanking the known motif were also of apparent consequences: in the case of Reb1, for example, a “T” at position –1 was associated with 100-fold reduction in binding (Fig. 4A). Reb1 protein was shown to bind the DNA base at the –1 position, and a “T” at that position was predicted to distort DNA shape (Jaiswal et al. 2016; Rossi et al. 2018). This effect outside of the core motif, however, was the exception; in most cases, variants of apparent ef-

fect were restricted to the motif itself, suggesting little contribution from the immediate motif-flanking region.

Our analysis therefore supports the notion that sequence variations within the known *cis*-regulatory motif reduce binding in a manner that depends on the TF and the precise sequence alternative. We next asked whether these same deviations from the consensus motif exert a similar cost on binding also when conserved in both species' genomes (common alternative). Here, we reasoned that deviations from the consensus that appear in both species' genomes have been preserved by selection, and may therefore reflect the need for lower binding, or, alternatively be compensated by contributions from adjacent sequences. For this analysis, we examined sites in which both alleles contain a motif variant that differs in the same one letter from the consensus motif (common alternative) and asked whether binding to these sites is weaker than binding to the consensus motif, as found in other locations in the genome. For Reb1, the apparent cost of common alternatives (Fig. 4B) was considerably lower than the cost of unique alternatives (the average binding fold change between consensus to alternative in common alternative sites is 1.54, but in unique alternative sites it is 3) (Fig. 4A). This same result extended to the majority of other TFs: the same alternative led to higher apparent cost when appearing in only one of the alleles (Fig. 4C) than when appearing in both alleles (Fig. 4D). The same effect was seen also when comparing unique alternative sites to consensus sites found elsewhere in the genome (Supplemental Fig. S9B). As a control, we validated that the consensus allele at unique sites is bound at the same level as sites of conserved consensus (Supplemental Fig. S9C).

Together, our results above support the notion that region-specific effects beyond the motif sequence act to modulate TF binding (Dror et al. 2015; Levo et al. 2015). This could occur through changes in motif accessibility, positioning of the motif within the promoter, or interaction with other DNA-bound cofactors. To test for such compensatory effects, we examined the sequence conservation between *S. cerevisiae* and *S. paradoxus* orthologs at the motif vicinity, postulating that if the surrounding region contributes to motif binding it will remain conserved between the similarly bound alleles (Fig. 4E). Focusing first on Reb1, we find that, as expected, common consensus and common alternative sites were in almost full conservation at the motif region (variation seen is a result of short INDELs). In contrast, sequence conservation decreased in the immediate vicinity of the motif in common consensus sites and in unique alternative sites but stayed relatively high in common alternative sites. The same pattern repeated when examining sequence conservation of seven yeast species (phastCons score) (Siepel et al. 2005), where common alternative sites showed higher conservation also at the motif region, and the conservation at their flanking region was higher than that of random-site background (Fig. 4F). Nucleosomes were equally depleted in both types of alternative sites, but depletion was deeper at common consensus sites (Supplemental Fig. S10). This pattern repeated also for Abf1, Rap1, and Tbf1 transcription factors, but was not apparent in other factors, in which conservation did not drop at the immediate motif vicinity (Supplemental Fig. S10). A related observation was reported before for CTCF binding in human, where motif mutations were associated with a strong reduction in binding when they appeared in a nonconserved genomic region, but the same mutations showed only a minor effect when they appeared in a highly conserved genomic region (Spivakov et al. 2012).

Overall, we find that when a weak motif appears in only one species, it diminishes binding, but when it is species-conserved, it allows a high level of binding. For specific TFs, the latter appear in

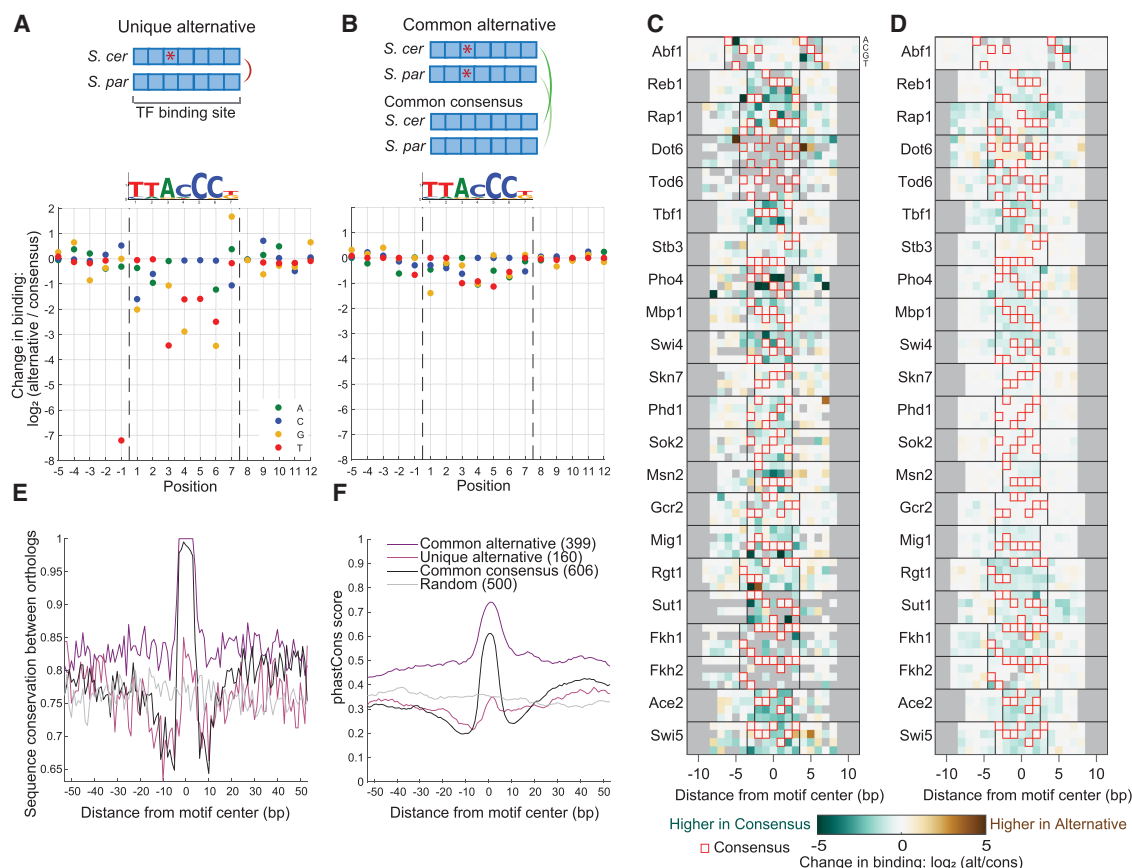


Figure 4. The cost of *cis*-regulatory mutations on TF binding. (A) The mutation cost of Reb1 measured at unique alternative sites, where one ortholog has the consensus motif (as defined *in vitro*) and the second ortholog has a one-letter variant (alternative). Each dot represents the mean of at least two sites (for the number of sites, see Supplemental Fig. S9B). Sequence logo of *in vitro* motif of Reb1 (Fordyce et al. 2010) is presented on top. (B) Mutation cost of Reb1, measured for common alternative sites, where both orthologs have the same one-letter variant. These sites are compared to common consensus sites found elsewhere in the genome. (C) Costs of unique alternatives for 22 TFs. The heatmap represents the change in binding as in A; here, the four rows stand for the four nucleotides A, C, G, T. Red box represents the consensus allele. Minimal two sites, gray color represent missing data. Bases flanking the motif have no consensus sequence, therefore the computation was performed relative to the most common nucleotide. (D) Cost of common alternatives for 22 TFs, as in C. (E) Common alternative binding sites are found at conserved genomic regions. Shown is the sequence conservation between *S. cerevisiae* and *S. paradoxus* orthologs (same nucleotide = 1, different nucleotide/ INDEL = 0) at Reb1 binding sites of type: common alternative, unique alternative, and common consensus, as well as in random sites at promoters. Shown is the mean signal per group. The number of sites in each group is indicated in parenthesis. (F) Common alternative sites are conserved through the yeast lineage. The phastCons conservation score (Siepel et al. 2005) is shown for the three Reb1 site groups as in E.

highly conserved genomic regions that perhaps compensate for the motif weakness.

Sequence variation in the motif predicts TF binding variation

Our analysis revealed that, for most TFs, interspecific variations in the core motif reduce DNA binding. To examine whether these differences in sequence are sufficient for predicting binding variations, we devised three linear predictive models. The first two models use a motif score as a single predictor: The first is based on the known *in vitro*-derived motif (PWM score), and the second is based on 7-mer sequence enrichment from our data (7-mer motif score). As seen in Figure 5A, variation in the 7-mer motif score was highly correlated with binding variation of Reb1 between the two alleles. The third model is a multivariate predictor that combines the two aforementioned motif scores and additional features: GC content at 15 bases flanking the motif, nucleosome occupancy at 300 bp centered at the motif, sequence conservation score (phastCons) (Siepel et al. 2005), and distance of the peak from

the closest transcription start site (Fig. 5B). To distinguish between cases of differential binding that resulted from motif variation and cases that resulted from other changes, we applied the models on different subsets of the data: all peaks, peaks associated with no motif, peaks associated with a motif, and finally, peaks associated with a nonconserved motif (Supplemental Fig. S11A).

Using the multivariate model, we could explain, on average, 35% of the variance in binding to the two alleles at peaks associated with a nonconserved motif (Fig. 5B). Percentage of variability explained ranged between TFs, with $R^2=0.03$ for Hms2 to $R^2=0.75$ for Ace2. In the majority of the TFs we examined (18/27), variability in binding to the two orthologs was well explained by sequence variation in the motif and its immediate surroundings ($R^2>0.3$), where in three cases sequence-based prediction exceeded R^2 of 0.5. Expanding the prediction to all motif-containing peaks (motif is either conserved or nonconserved) resulted in a somewhat lower predicting power ($0.02<R^2<0.66$, median = 0.32) (Supplemental Fig. S11A). Only 19% of the differentially bound peaks (with more than twofold change between alleles)

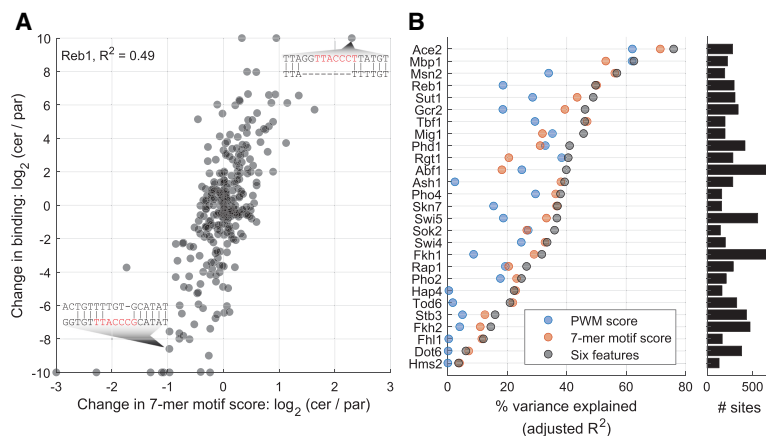


Figure 5. Sequence variation predicts DNA binding variation. (A) Change in motif score predicts variation in Reb1 binding. Shown is the log₂-ratio of 7-mer motif scores (x-axis) and ChEC-seq signal (y-axis) between *S. cerevisiae* and *S. paradoxus* orthologs, at peaks associated with a nonconserved Reb1 motif (i.e., a different motif sequence appears in each ortholog). Sequence alignment at two specific sites is presented: *S. cerevisiae* ortholog on the upper row, *S. paradoxus* ortholog in the lower row, and the Reb1 motif written in red. (B) Linear models predict binding variation at peaks associated with nonconserved motifs. Shown is the percentage of explained variability ($R^2 \times 100$) for each TF, using three models: in vitro PWM score, 7-mer motif score derived from our data, and a compilation of these with another four predictors (see text). Predictions for other peak categories are presented in Supplemental Figure S11A.

were associated with a nonconserved motif, consistent with previous results (Reddy et al. 2012). However, when accounting only for motif-associated peaks, most (60%) of differentially bound peaks were associated with a nonconserved motif (Supplemental Fig. S11C). Limiting the prediction to peaks with no motif resulted in no predicting power ($0.04 < R^2 < 0.26$, median = 0.11); hence, the additional features added to the model beyond the motif could not explain binding variability in the absence of a motif (Supplemental Fig. S11A).

In most cases (21/27 TFs), the log-ratio of motif score derived from the data was the best predictor of binding variation. Notable exceptions include Abf1 and Rgt1, for which the 7-mer score did not capture the full motif (Supplemental Table S1), likely because of the large gap between the two parts of Abf1 motif and the long A-stretch of Rgt1 motif. Difference in GC content was predictive only for Skn7 peaks in correlation with the GC-rich motif of this TF (Supplemental Fig. S11B). Other features had no predicting power (Supplemental Fig. S11B). Transforming the features and the predicted change in binding to absolute level (i.e., predicting how big the change is, regardless of its direction) resulted in lower R^2 values and therefore lower prediction power (Supplemental Fig. S11D).

Based on the studied factors, we conclude that the variation in motif sequence within binding sites is a strong predictor of binding variation for the majority of TFs, reaching 35% explained variation on average.

Gains and losses of binding sites are more common than binding-site turnover

Our analysis so far focused on TF binding at individual binding peaks. We next revisited the integration of binding peaks within the context of the full promoter. Specifically, we wished to define the prevalence of binding-site turnover, whereby, for example, a loss of a binding site in one location along a promoter is compensated by the gain of a binding site at an adjacent location within the same promoter.

To characterize cases of binding-site turnover we classified promoters into four classes: (1) conserved promoters: cases in which all binding sites are common to both alleles (53% of promoters) (Fig. 6A); (2) turnover promoters: cases in which binding sites appear in both orthologs, but on different locations along the promoter, suggesting reciprocal gain or loss of binding motifs (7% of promoters) (Fig. 6B); (3) unbalanced promoters: cases in which one or more binding site is allele-specific, but other binding sites remain conserved (9% of promoters) (Fig. 6C); and (4) fully unbalanced promoters: cases in which only one of the alleles is bound by the TF (27% of promoters) (Fig. 6D). Cases of conserved motifs that are bound at only one allele were considered as “not defined” (5% of promoters). As in previous analyses, we also considered only peaks that reside next to a strong binding motif (with FIMO P -value < 0.001).

We built a custom algorithm for promoter classification, which takes a list of peaks as input, and classifies motifs, binding sites, and promoters into the aforementioned classes (Supplemental Fig. S12A). To assess the algorithm performance, we manually defined 220 promoters, including up to 20 cases of each class for three TFs, and compared our manual classification with the algorithm output per promoter class. We observe mean sensitivity of 82% and mean specificity of 94% across the different classes (Supplemental Fig. S12B,C).

We find that the conserved and fully unbalanced classes were the largest promoter classes, consisting of 53% and 27% of all examined promoters, respectively (Fig. 6E, left). Conserved promoters are bound to a higher level (59% of the signal, summing across all TFs) as compared to fully unbalanced promoters (15% of the signal; considering the more highly bound allele) (Fig. 6E, right). Among the different TFs there was little variation in the proportion of promoter classes and their binding levels (Supplemental Fig. S13A). These trends generally repeated also when elevating the minimal peak threshold, although the proportion of conserved promoters increased with increased threshold (Supplemental Fig. S13B). We note that the different promoter classes are bound to different extents, on average, although our classifier does not take the total promoter binding level into account: the turnover and unbalanced classes were bound at levels twice as high as the conserved promoters (Supplemental Fig. S13C). This higher binding reflected a larger number of binding sites in these classes, while binding at individual sites was at a similar level (Supplemental Fig. S13C).

Binding-site turnover is a result of reciprocal gains and losses of binding sites. The distance between the turning-over binding sites on the aligned sequence coordinates could be long, as in *CDC5* promoter (35 bp), or short, as in *YBL055C* promoter, where the two motifs overlap but appear on different strands (Fig. 6B,F). Examining the full set of TFs, we find that most turning-over binding sites appear in close proximity (median distance = 20 bp), and in 37% of these, the distance is 10 base pairs or less (Supplemental Fig. S14). Specific examples of short-distance binding-site turnover are presented in Figure 6F for Reb1-bound and Ace2-bound

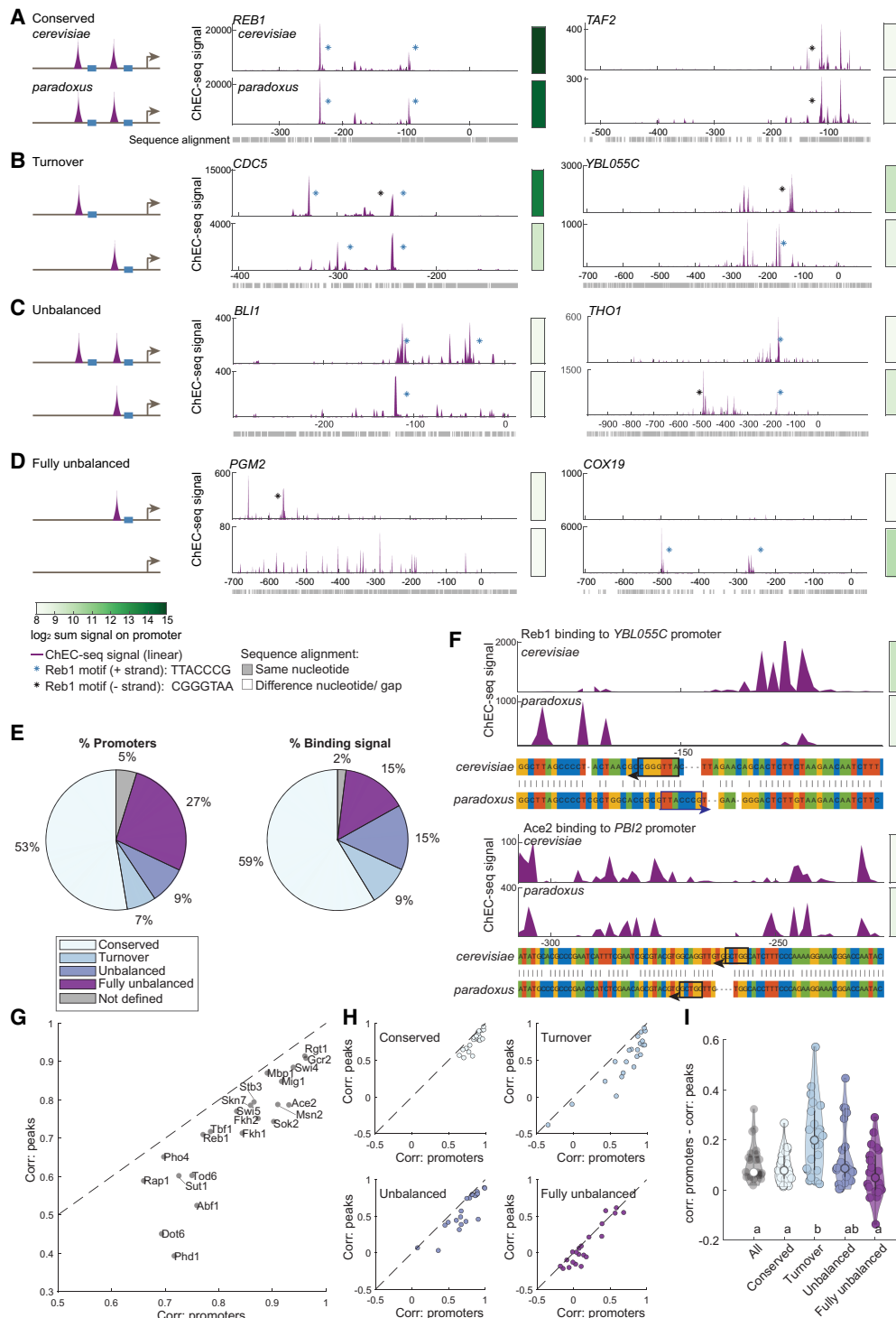


Figure 6. Promoter evolution and binding-site turnover. (A–D) Four evolutionary classes of TF binding variation. Schemes are shown in the *left* panels; genome browser snapshots of Reb1-bound promoters are shown as examples in the *middle* and *right* panels. Sequence alignment between orthologs is presented *below* each example, x-axis represents the location on the promoter relative to TSS (figure legend *below* D). (A) Conserved: all binding sites are species-conserved. (B) Turnover: reciprocal gain and loss of binding sites. (C) Unbalanced: species-specific sites along with conserved sites. (D) Fully unbalanced: binding sites appear in only one ortholog. (E) Proportion of number of promoters and binding signal per promoter class. Shown is the distribution of the full set of TFs; proportion per TF is presented in Supplemental Figure S13A. Binding signal refers to the ChEC-seq signal in the higher-bound ortholog, normalized by the total signal in that ortholog. (F) Examples of short-distance binding-site turnover. Shown are the binding signal in both alleles, and the sequence alignment in the *lower* panel. Boxes mark motifs, and arrows mark the motif's strand. (G) Correlation on peaks when comparing orthologs. Shown are correlation coefficients between orthologs, over motif-associated peaks (y-axis) and promoters (x-axis), among all promoters. Here, we summed the binding signal only on peaks within peak-containing promoters, but the correlation coefficients were quantitatively similar to those obtained from the more simplistic approach of summing up the signal over the full promoter, as presented in Figure 1B (cf. Fig. S15A). (H) The shift in promoter correlation versus peak correlation is more apparent at turnover and unbalanced promoters. Shown is the correlation between orthologs, summing over promoters (x-axis) and over motif-associated peaks (y-axis) as in G, per promoter class. (I) Turnover promoters show a higher promoter similarity despite lower peak similarity. Shown are the differences between correlation on promoters to correlation on peaks for the different promoter classes. Each dot represents a TF; letters represent statistically distinguished groups after Tukey's honestly significant difference test.

promoters. In the case of Reb1, peaks appear upstream of the motif in divergent directions corresponding to the appearance of the motif on different strands; in the case of Ace2, the signal-depleted area aligns with the motif location of each ortholog. This resembles evolutionary conservation of binding site location in the presence of sequence divergence.

Another aspect of binding-site turnover would be buffering of differential binding: although binding sites appear on different locations in the two orthologs, the total promoter binding should remain similar. To test that we plotted the correlation between orthologs, summing either on total promoters or on individual peaks (Fig. 6G). Indeed, correlation coefficients were higher on promoters than on peaks. For a control, we summed the binding signal along increasing genomic bins (30, 100, 300, and 1000 bp) and found a strong shift in correlation between the bins of 30 and 100 bp (individual peaks are 20 bp wide), resembling higher divergence in peak binding relative to promoter binding (Supplemental Fig. S15C,D). Experimental repeats showed high correlation ($R > 0.9$) in all examined bins, in most TFs (23 of 27), and along motif-associated peaks (Supplemental Fig. S15B). To examine this shift in correlation in more detail, we repeated the same analysis but separately for the four promoter classes (Fig. 6H). Namely, for each TF, we considered each time only a single promoter class and examined the shift in orthologous correlation at the promoter level to the correlation on the individual peak level. Indeed, turnover promoters showed the largest shift in promoter-peak correlation (significantly different than the shift in conserved promoters, P -value = 0.003, Tukey's honestly significant difference test), regardless of these being a relatively small fraction of total promoters (Fig. 6I). The shift observed in unbalanced promoters results from the presence of an additional conserved motif, which increases the correlation at the promoter level.

To conclude, we find that half (53%) of bound promoters retained a fully conserved set of binding sites. Among diverged promoters, a large fraction is bound in only one of the alleles (27%), and a small fraction (7%) shows compensation through binding-site turnover, likely indicating functional conservation despite sequence divergence.

Discussion

Understanding the sequence determinants of transcription factor binding in living cells is a major challenge. In this study, we profiled *in vivo* TF binding in a yeast hybrid that contains two related genomes within the same nucleus. The high sequence divergence in regulatory regions between these two genomes (~25%) provides a wide range of sequence variation that can be examined in parallel. This enabled us to measure how multitudes of sequence variations affect TF binding within a genomic context.

Differential binding can result from differential use of a common pool of potential binding sites containing the binding motif, or from gain/loss of sequence motifs. Our data support the second scenario, because we found that most conserved motif sites were either bound or unbound in both genomes, whereas differentially bound sites were associated with a sequence variation within the binding motifs. Further, we detected little, if any, differences in nucleosome positioning at sites that are differentially bound. Therefore, in yeast, sequence evolution in regulatory regions appears to occur more readily than changes in chromatin accessibility at TF binding sites, in agreement with previous studies (Tirosh et al. 2010; Tsankov et al. 2010). The ability of a TF to bind its motif-containing site was also shown to depend on DNA shape fea-

tures (Abe et al. 2015; Zhou et al. 2015); however, these features did not separate between *in vivo* bound and unbound motif-containing sites (Zentner et al. 2015) and therefore were not examined in this study.

TFs bind strongly at sites containing their consensus motif, compatible with sequence motifs defined *in vitro*. We expected to observe binding also at sites containing imprecise motifs, for example, sites containing one alternative base, but that this binding would be on average lower than binding to the consensus motif. This was indeed the case when an alternative base was present in one allele only, while the second allele carried the consensus motif. Notably, however, we find that genomic sites containing the same alternative allele in both orthologous genomes were bound at almost the same extent as sites containing the consensus motif. A perfect consensus sequence might not always be the best for the organism in terms of fitness at each site, however. Nonconsensus or low-affinity binding sites are in fact widespread in the yeast genome (Tanay 2006) and were shown to be important in fly and mouse development (Scardigli et al. 2003; Rowan et al. 2010; Crocker et al. 2016).

We have further shown that imprecise (weak) motifs of certain TFs, which appear in both orthologous alleles, commonly reside in regions of exceptionally high sequence conservation. This observation is in contrast to the cases of consensus sites, which are often found as islands of conservation within sequence-diverged regions. This may be related to a previous report showing that polymorphisms in CTCF motifs have greater effects on TF binding when they appear in sequence-diverged regions than in sequence-conserved regions (Spivakov et al. 2012). The investigators attribute this effect to cofactors that allow CTCF to bind at imprecise motifs. In our case, such binding partners are known for Rap1 (Tornow et al. 1993) but have not been described for the other TFs. We find it more likely that this effect is related to chromatin, because the TFs showing this effect (Abf1, Reb1, and Rap1) all act in the regulation of nucleosome positioning (Rhee and Pugh 2011). Therefore, we speculate that the presence of other TFs at nearby sites stabilizes the binding of these TFs to nonoptimal motif sites (Mirny 2010). Another possibility involves DNA interactions through the non-DNA-binding domain parts of the protein (Brodsky et al. 2020). To conclude, we suggest that nonconsensus sites are bound to a high level and are species-conserved owing to a local sequence compensation.

In the majority of the TFs we examined, variability in binding to the two orthologs was well explained by sequence variation in the motif and its immediate surroundings. Previous studies addressing the problem of predicting binding variation from sequence variation reported on a generally limited predicting power for differential TF binding (Bradley et al. 2010; Zheng et al. 2010; He et al. 2011b; Reddy et al. 2012; Stefflova et al. 2013; Halow et al. 2021). In an influential study for the field, Reddy et al. (2012) measured TF binding of in human heterozygous cell lines, and reported that only 12% of differentially bound sites were associated with sequence variations in known binding motifs. Here, we report a similar fraction (19%); however, when considering only motif-associated peaks, we find that most (60%) of the differentially bound sites were associated with a sequence variation in the motif. Further, using quantitative models, we show that variability in motif score is the best predictor for variability in TF binding, whereby other features had limited contribution. Our improved prediction may result from the use of the ChEC-seq method, which provides high-resolution mapping of TF binding. In addition, the use of F1 hybrids for this work allows

profiling both orthologous genomes in the same cell and thus reduces both technical and *trans*-driven variations that can reduce power and hamper interspecific comparisons.

When classifying evolutionary changes in TF binding at the promoter level, we find that most of the bound promoters involve unbalanced gain or loss of binding sites, whereas only 7% of the bound promoters show evidence of compensation by binding site turnover. This result is in agreement with reports from other model organisms, including the Zeste TF in *Drosophila* (Moses et al. 2006), liver-specific TFs in mice (Stefflova et al. 2013), and individual TFs in yeast (Borneman et al. 2007). Furthermore, higher rates of positive and purifying selection compared to compensatory neutral evolution were modeled for *Drosophila* enhancers (He et al. 2011a), but not for yeast promoters (Mustonen et al. 2008). The flexibility and high rate of binding site evolution suggests that in many cases binding site loss or addition is not deleterious. Overall, we find that TF binding evolves through gains and losses of binding sites, with quantitative changes in binding level being highly predictable from sequence variation within the motif.

To conclude, in this study we report on two linked observations: (1) imprecise but conserved motifs are bound to a high level by TFs, and (2) the observation of short-distance binding sites turnover, where binding localization is conserved despite of sequence divergence. These observations show the fast and flexible evolution of TF binding sites between related species, and we expect to see these phenomena in other organisms as well.

Methods

Yeast strains

Yeast strains in this study were constructed on the background of *S. cerevisiae* BY4741 and *S. paradoxus* CBS432 (OS142) and their hybrid. For ChEC-seq, transcription factors were tagged with MNase on their C terminus, by amplifying the MNase-KanMX cassette from the pGZ108 plasmid, a gift from Steven Henikoff. Strains that were previously generated in our laboratory were based on transformation of BY4741 with MNase-KanMX cassette, with an ORF-MNase linker of 33 amino acids (Bar-Ziv et al. 2020; Brodsky et al. 2020; Lupo et al. 2021). In this study, strains were generated on the background of the C-SWAT library (Meurer et al. 2018) a gift from Maya Schuldiner. In these strains, the MNase-KanMX cassette was inserted between L3 and L4 linkers, with an ORF-MNase linker of 15 amino acids. Free-MNase strain contains MNase from the pGZ108 plasmid, without any linker, under the *TDH3* promoter, integrated into the *MSN2* genomic locus.

Transformations to *S. cerevisiae* were performed using the traditional LiAc/SS DNA/PEG method (Gietz et al. 1995). Transformations to *S. paradoxus* were performed using SORB-competent cells (Bleuven et al. 2019). Strains are listed in [Supplemental Table S2](#); primers are listed in [Supplemental Table S3](#).

ChEC-seq

ChEC-seq experiments were performed as described previously (Zentner et al. 2015) with modifications. In this study, replicates are biological replicates, starting from separate overnight starters of the same strain. Each TF was profiled in at least two replicates. Cultures were grown overnight to saturation in YPD media and diluted into 5 mL of fresh YPD media to reach OD₆₀₀ of 4 the following morning after ~10 divisions. Cultures were pelleted at 1500g and resuspended in 1 mL Buffer A (15 mM at Tris pH 7.5, 80 mM KCl, 0.1 mM EGTA, 0.2 mM spermine, 0.5 mM spermidine, 1× Roche cOmplete EDTA-free mini protease inhibitors, 1 mM

PMSF) and then transferred to DNA low-bind tubes (Eppendorf 022431021). Cells were washed twice more in 500 µL Buffer A, pelleted, and resuspended in 150 µL Buffer A containing 0.1% digitonin. Then, cells were transferred to an Eppendorf 96-well plate (Eppendorf 951020401) for permeabilization (for 5 min at 30°C). CaCl₂ was added to a final concentration of 2 mM for 30 sec. Next, 100 µL of stop buffer (400 mM NaCl, 20 mM EDTA, 4 mM EGTA, and 1% SDS) was mixed with 100 µL of sample. Proteinase K was then added (5 µL of 20 mg/mL) and incubated for 30 min at 55°C. Nucleic acids were extracted with an equal volume (200 µL) of ultrapure phenol/chloroform/isoamyl alcohol, and ethanol-precipitated (for >1 h at −80°C) with 2.5 volumes of cold EtOH 96%, 45 µg Glycoblue, and sodium acetate to a final concentration of 20 mM. DNA was centrifuged (for 10 min at 4°C), washed with EtOH 70%, and treated with RNase A in a final concentration of 2.5 mg/mL (for 20 min at 37°C), followed by another round of DNA cleanup and ethanol precipitation. To enrich for small DNA fragments, reverse 0.8× SPRI cleanup (right-side size selection) was performed, followed by isopropanol SPRI (left-side size selection) of 1.8× SPRI and 5.4× isopropanol. DNA samples were eluted in 20 µL 0.1× TE.

Library preparation was performed similarly to a published protocol (Skene and Henikoff 2017) with specific modifications. End-repair and A-tailing (ERA) of the small DNA fragments was performed by the following: [T4 DNA ligase buffer (10×), dNTPs (10 mM), ATP (10 mM), 50% PEG 4000, T4 PNK (6 units), T4 DNA Pol (0.3 units), Taq DNA Pol (0.1 µL per sample) with 14.6 µL sample] with the PCR protocol: lid heated to 60°C, for 15 min at 12°C, 15 min at 37°C, and 45 min at 58°C. Samples were cleaned with reverse 0.5× SPRI followed by left-side isopropanol SPRI: 1.3× SPRI (with the previous step reaches to 1.8× SPRI) and 5.4× isopropanol. Indexed adaptors (Blecher-Gonen et al. 2013) were ligated to the DNA using quick ligase (2000 units/µL, 2 µL per sample) and quick ligase buffer (2×), for 15 min at 20°C. Cleanup was performed: 1.2× SPRI (left side) followed by addition of 1.2× HXN buffer (24 µL 5 M NaCl, 19.2 µL 50% PEG 8000, and 4.8 µL H₂O per sample), reaching 1.6× SPRI. Library amplification was performed with library-specific enrichment primers (23 µL sample DNA, 2 µL enrichment primers, 25 µL KAPA Hifi PCR mix) with the following PCR protocol: 45 sec in 98°C, 16 cycles of 15 sec at 98°C and 15 sec at 60°C, and a final elongation step of 1 min at 72°C. PCR products were cleaned with left-side 1× SPRI. Library concentration was measured with Q-bit, and library size distribution was measured with TapeStation. Libraries were sequenced on Illumina NovaSeq and NextSeq500 machines, with 51-base paired end reads.

Computational analysis

Programs

Programs used for read alignment are indicated below. Downstream analyzes were originally implemented in MATLAB 2019 and in R 3.6.3 (R Core Team 2013). Online programs of the MEME suite were used as well.

Read alignment

FASTQ reads were trimmed from adaptors with cutadapt (Martin 2011), then aligned to a the hybrid genome using Bowtie 2 (Langmead and Salzberg 2012) with the parameters: [−p8 −local −very-sensitive −trim-to 40 −dovetail −score-min G,16,8]. The hybrid genome is a concatenation of the genomes of *S. cerevisiae* S288c (R64-1-1/sacCer3) and *S. paradoxus* CBS432 (Yue et al. 2017), including the mitochondrial genomes. Bowtie 2 reports on one (or zero) alignments per read, therefore a given read was

mapped only once to one of the parental genomes. Reads with zero mismatches were 94%–95% of total mapped reads in three representative samples. Genome coverage of the 5' end of reads was generated using SAMtools (Li et al. 2009) and BEDTools (Quinlan and Hall 2010), with the `genomcov` parameters: [-5 -fs 1].

ChEC-seq data normalization and processing

Raw genome coverage counts were divided by the total number of reads and multiplied by 10^7 . Gene promoters were defined using two published data sets of *S. cerevisiae* transcription start sites (Pelechano et al. 2013; Park et al. 2014), where the version with the shorter 5' UTR, in which the TSS is upstream of the start codon, was selected per gene. *S. paradoxus* TSS were defined for orthologous genes using the 5' UTR lengths defined for *S. cerevisiae*. For both genomes, location of specific TSS was manually edited based on functional genomic data. Promoters were defined as intergenic regions 400 bps upstream to the TSS or to the position where a promoter meets another transcript. Promoters were defined for 5105 of 6701 genes.

Motif enrichment

Motif enrichment was performed using two methods.

The first method was the Motif score. As in Brodsky et al. (2020), all possible sequences of length k (k -mers) were given a numerical index (16,384 possibilities of 7-mers), where each nucleotide in the hybrid genome was indexed accordingly. To compute the motif score of a given sample, ChEC-seq signal was smoothed (moving average of 20 nt), and the averaged signal for each k -mer was then calculated across all of its occurrences in all promoters. Motif scores of TFs were based on 7-mer sequences.

The second method was MEME-ChIP. Sequences of 60 bp centered at top peaks (98% bootstrap level) were extracted per TF and were used as an input for MEME-ChIP (Machanic and Bailey 2011), with YEASTRACT (Teixeira et al. 2006) and JASPAR (Sandelin et al. 2004) as reference databases.

Probability weight matrices (PWMs)

In vitro PWMs were collected from the public databases YeTFaSCO (De Boer and Hughes 2012) and JASPAR (Sandelin et al. 2004) and are listed in Supplemental Table S1. To allocate significant realizations of these motifs in the hybrid genome, we used FIMO (Grant et al. 2011), with the in vitro PWM and aligned hybrid genome as input, with significance threshold of P -value < 0.001 .

Data-driven PWMs of the different TFs were generated based on the top 20 7-mer sequences of each factor, as in Brodsky et al. (2020). Sequence logos were generated with LogoMaker (Tareen and Kinney 2020).

Aligned genome coordinates

To directly compare ChEC-seq signal and sequence variation between the hybrid alleles, we aligned orthologous gene promoters and ordered the genomic data accordingly, as done previously (Venkataram and Fay 2010). Specifically, we extracted 5105 orthologous and locally aligned their upstream intergenic region with MATLAB function (`swalign`) with a gap-opening penalty of 10, gap-extension penalty of 0.5, and "NT" alphabet. This resulted in a reduced, comparable genome of 2,544,708 million base pairs.

Peak calling

Peaks were called from smoothed ChEC-seq profiles (5' end of reads, 20 bases moving average) using MATLAB (`findpeaks`) function with the following parameters: "MinPeakHeight" was defined

from the data, "MinPeakProminence" was equal to "MinPeakHeight," "MinPeakDistance" was 20 bases, "MinPeakWidth" was 10 bases. Because the basal signal level was higher in promoters with high peaks, only peaks that exceeded the 90th percentile of their promoter signal were selected. "MinPeakHeight" definition was the 95th percentile of signal at random sites on promoters. Peak tables are provided as Supplemental Table S6.

Peak-motif association

The highest motif score was located at a range of 60 bp centered at the peak. Orthologous peaks that were separated by less than 10 bases were unified into a single peak location. Peaks further than 800 from any TSS were filtered out from further analysis.

Position-specific mutation cost

To measure the binding cost owing to mutations at specific positions of the motif, the following analysis was performed: peaks were aligned relative to the location of their maximal motif score. Then, for each peak, the motif sequence of the better-scored allele was aligned to the motif PWM for sequence comparison. Alignment to the PWM was based on the product of probabilities (P) of all positions. To allow flexibility, the minimal PWM score allowed sequence variability at positions with maximal probability (P) < 0.7 . This way Reb1 sites of TTACCCG and TTACCCT were both allowed. Sequence substitution was analyzed relative to the motif consensus sequence, which is the maximal PWM-scoring sequence. To find sequences with an alternative allele, an iterated algorithm was implemented: in each iteration a certain position of the motif is "mutated" so the nucleotide probabilities in that position equals 0.25. Figure 4 summarizes this analysis, where the average \log_2 ratio of alternative to consensus is shown.

Prediction of TF binding variation

Multiple linear models were analyzed in R 3.6.3 (R Core Team 2013). Relative feature importance was analyzed using `RelaImpo` package (Grömping 2006).

Data access

All raw and processed sequencing data generated in this study have been submitted to the NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE196451. The FASTQ data generated in this study have been submitted to the NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA700498. Source code used in this study is available as Supplemental Code and at GitHub (<https://github.com/GatKrieger/TFhybrid>).

Competing interest statement

The authors declare no competing interests.

Acknowledgments

We thank the Barkai laboratory members, especially Felix Jonas, Michal Chapal, and Sagie Brodsky for their help with experiments, useful discussions, and constructive comments on the manuscript. This project was supported by the U.S. National Science Foundation–U.S. Israel Binational Science Foundation–Molecular and Cellular Biosciences (NSF-BSF-MCB) (2019625), the Israel

Science Foundation (ISF) (1738/15), and the Minerva Center (AZ 57 46 9407 65).

Author contributions: G.K. and N.B. designed the research. G.K. and O.L. performed experiments. G.K. analyzed the data. All authors contributed to the writing of the paper.

References

- Abe N, Dror I, Yang L, Slattery M, Zhou T, Bussemaker HJ, Rohs R, Mann RS. 2015. Deconvolving the recognition of DNA shape from sequence. *Cell* **161**: 307–318. doi:10.1016/j.cell.2015.02.008
- Artieri CG, Fraser HB. 2014. Evolution at two levels of gene expression in yeast. *Genome Res* **24**: 411–421. doi:10.1101/gr.165522.113
- Avsec Z, Weilert M, Shrikumar A, Krueger S, Alexandari A, Dalal K, Fropf R, McAnany C, Gagneur J, Kundaje A, et al. 2021. Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat Genet* **53**: 354–366. doi:10.1038/s41588-021-00782-6
- Bar-Ziv R, Brodsky S, Chapal M, Barkai N. 2020. Transcription factor binding to replicated DNA. *Cell Rep* **30**: 3989–3995.e4. doi:10.1016/j.celrep.2020.02.114
- Blecher-Gonen R, Barnett-Itzhaki Z, Jaitin D, Amann-Zalcenstein D, Lara-Astiaso D, Amit I. 2013. High-throughput chromatin immunoprecipitation for genome-wide mapping of *in vivo* protein-DNA interactions and epigenomic states. *Nat Protoc* **8**: 539–554. doi:10.1038/nprot.2013.023
- Blouven C, Dubé AK, Nguyen GQ, Gagnon-Arsenault I, Martin H, Landry CR. 2019. A collection of barcoded natural isolates of *Saccharomyces paradoxus* to study microbial evolutionary ecology. *Microbiologyopen* **8**: e773. doi:10.1002/mbo3.773
- Borneman AR, Gianoulis TA, Zhang ZD, Yu H, Rozowsky J, Seringhaus MR, Wang LY, Gerstein M, Snyder M. 2007. Divergence of transcription factor binding sites across related yeast species. *Science* **317**: 815–819. doi:10.1126/science.1140748
- Bradley RK, Li XY, Trapnell C, Davidson S, Pachter L, Chu HC, Tonkin LA, Biggin MD, Eisen MB. 2010. Binding site turnover produces pervasive quantitative changes in transcription factor binding between closely related *Drosophila* species. *PLoS Biol* **8**: e1000343. doi:10.1371/journal.pbio.1000343
- Brodsky S, Jana T, Mittelman K, Chapal M, Kumar DK, Carmi M, Barkai N. 2020. Intrinsically disordered regions direct transcription factor *in vivo* binding specificity. *Mol Cell* **79**: 459–471.e4. doi:10.1016/j.molcel.2020.05.032
- Carroll SB. 2005. Evolution at two levels: on genes and form. *PLoS Biol* **3**: e245. doi:10.1371/journal.pbio.0030245
- Crocker J, Preger-Ben Noon E, Stern DL. 2016. The soft touch: low-affinity transcription factor binding sites in development and evolution. *Curr Top Dev Biol* **117**: 455–469. doi:10.1016/bs.ctdb.2015.11.018
- De Boer CG, Hughes TR. 2012. YeTFaSCo: a database of evaluated yeast transcription factor sequence specificities. *Nucleic Acids Res* **40**: D169–D179. doi:10.1093/nar/gkr993
- Dermitzakis ET, Clark AG. 2009. Evolution of transcription factor binding sites in mammalian gene regulatory regions: handling counterintuitive results. *J Mol Evol* **68**: 654–664. doi:10.1007/s00239-009-9238-1
- Dror I, Golan T, Levy C, Rohs R, Mandel-Gutfreund Y. 2015. A widespread role of the motif environment in transcription factor binding across diverse protein families. *Genome Res* **25**: 1268–1280. doi:10.1101/gr.184671.114
- Emerson JJ, Hsieh LC, Sung HM, Wang TY, Huang CJ, Lu HHS, Lu MYJ, Wu SH, Li WH. 2010. Natural selection on *cis* and *trans* regulation in yeasts. *Genome Res* **20**: 826–836. doi:10.1101/gr.101576.109
- Floc'hlay S, Wong ES, Zhao B, Viales RR, Thomas-Chollier M, Thieffry D, Garfield DA, Furlong EEM. 2021. *Cis*-acting variation is common across regulatory layers but is often buffered during embryonic development. *Genome Res* **31**: 211–224. doi:10.1101/gr.266338.120
- Fordyce PM, Gerber D, Tran D, Zheng J, Li H, DeRisi JL, Quake SR. 2010. *De novo* identification and biophysical characterization of transcription-factor binding sites with microfluidic affinity analysis. *Nat Biotechnol* **28**: 970–975. doi:10.1038/nbt.1675
- Gera T, Jonas F, More R, Barkai N. 2021. Evolution of binding preferences among whole-genome duplicated transcription factors. *eLife* **11**: e73225. doi:10.7554/eLife.73225
- Gietz RD, Schiestl RH, Willems AR, Woods RA. 1995. Studies on the transformation of intact yeast cells by the LiAc/SS-DNA/PEG procedure. *Yeast* **11**: 355–360. doi:10.1002/yea.320110408
- Grant CE, Bailey TL, Noble WS. 2011. FIMO: scanning for occurrences of a given motif. *Bioinformatics* **27**: 1017–1018. doi:10.1093/bioinformatics/btr064
- Grömping U. 2006. Relative importance for linear regression in R: the package relaimpo. *J Stat Softw* **17**: 1–27. doi:10.18637/jss.v017.i01
- Halow JM, Byron R, Hogan MS, Ordoñez R, Groudine M, Bender MA, Stamatoyannopoulos JA, Maurano MT. 2021. Tissue context determines the penetrance of regulatory DNA variation. *Nat Commun* **12**: 2850. doi:10.1038/s41467-021-23139-3
- He BZ, Holloway AK, Maerkl SJ, Kreitman M. 2011a. Does positive selection drive transcription factor binding site turnover? A test with *Drosophila cis*-regulatory modules. *PLoS Genet* **7**: e1002053. doi:10.1371/journal.pgen.1002053
- He Q, Bardet AF, Patton B, Purvis J, Johnston J, Paulson A, Gogol M, Stark A, Zeitlinger J. 2011b. High conservation of transcription factor binding and evidence for combinatorial regulation across six *Drosophila* species. *Nat Genet* **43**: 414–420. doi:10.1038/ng.808
- Hill MS, Vande Zande P, Wittkopp PJ. 2021. Molecular and evolutionary processes generating variation in gene expression. *Nat Rev Genet* **22**: 203–215. doi:10.1038/s41576-020-00304-w
- Jaiswal R, Choudhury M, Zaman S, Singh S, Santosh V, Bastia D, Escalante CR. 2016. Functional architecture of the Reb1-Ter complex of *Schizosaccharomyces pombe*. *Proc Natl Acad Sci* **113**: E2267–E2276. doi:10.1073/pnas.1525465113
- Kaplan N, Moore IK, Fondufe-Mittendorf Y, Gossett AJ, Tillo D, Field Y, LeProust EM, Hughes TR, Lieb JD, Widom J, et al. 2009. The DNA-encoded nucleosome organization of a eukaryotic genome. *Nature* **458**: 362–366. doi:10.1038/nature07667
- Kasowski M, Grubert F, Heffelfinger C, Hariharan M, Asabere A, Waszak SM, Habegger L, Rozowsky J, Shi M, Urban AE, et al. 2010. Variation in transcription factor binding among humans. *Science* **328**: 232–235. doi:10.1126/science.1183621
- Kilpinen H, Waszak SM, Gschwind AR, Raghav SK, Witwicki RM, Orioli A, Migliavacca E, Wiederkehr M, Gutierrez-arcelus M, Panousis NI, et al. 2013. Coordinated effects of sequence variation on DNA binding, chromatin structure, and transcription. *Science* **342**: 744–747. doi:10.1126/science.1242463
- Krieger G, Lupo O, Levy AA, Barkai N. 2020. Independent evolution of transcript abundance and gene regulatory dynamics. *Genome Res* **30**: 1000–1011. doi:10.1101/gr.261537.120
- Kristiansson E, Thorsen M, Tamás MJ, Nerman O. 2009. Evolutionary forces act on promoter length: identification of enriched *cis*-regulatory elements. *Mol Biol Evol* **26**: 1299–1307. doi:10.1093/molbev/msp040
- Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**: 357–359. doi:10.1038/nmeth.1923
- Levo M, Zalckvar E, Sharon E, Machado ACD, Kalma Y, Lotan-Pompan M, Weinberger A, Yakhini Z, Rohs R, Segal E. 2015. Unraveling determinants of transcription factor binding outside the core binding site. *Genome Res* **25**: 1018–1029. doi:10.1101/gr.185033.114
- Li H, Handsaker B, Wysoker J, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079. doi:10.1093/bioinformatics/btp352
- Ludwig MZ, Bergman C, Patel NH, Kreitman M. 2000. Evidence for stabilizing selection in a eukaryotic enhancer element. *Nature* **403**: 564–567. doi:10.1038/35000615
- Lupo O, Krieger G, Jonas F, Barkai N. 2021. Accumulation of *cis*- and *trans*-regulatory variations is associated with phenotypic divergence of a complex trait between yeast species. *G3 (Bethesda)* **11**: jkab016. doi:10.1093/g3journal/jkab016
- Machanic P, Bailey TL. 2011. MEME-ChIP: motif analysis of large DNA datasets. *Bioinformatics* **27**: 1696–1697. doi:10.1093/bioinformatics/btr189
- Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J* **17**: 10–12. doi:10.14806/ej.17.1.200
- Maurano MT, Haugen E, Sandstrom R, Vierstra J, Shafer A, Kaul R, Stamatoyannopoulos JA. 2015. Large-scale identification of sequence variants influencing human transcription factor occupancy *in vivo*. *Nat Genet* **47**: 1393–1401. doi:10.1038/ng.3432
- McNabb DS. 2005. Assembly of the Hap2p/Hap3p/Hap4p/Hap5p-DNA complex in *Saccharomyces cerevisiae*. *Eukaryot Cell* **4**: 1829–1839. doi:10.1128/EC.4.11.1829-1839.2005
- Metzger BPH, Duveau F, Yuan DC, Tryban S, Yang B, Wittkopp PJ. 2016. Contrasting frequencies and effects of *cis*- and *trans*-regulatory mutations affecting gene expression. *Mol Biol Evol* **33**: 1131–1146. doi:10.1093/molbev/msw011
- Meurer M, Duan Y, Sass E, Kats I, Herbst K, Buchmuller BC, Dederer V, Huber F, Kirmmaier D, Štefl M, et al. 2018. Genome-wide C-SWAT library for high-throughput yeast genome tagging. *Nat Methods* **15**: 598–600. doi:10.1038/s41592-018-0045-8
- Mirny LA. 2010. Nucleosome-mediated cooperativity between transcription factors. *Proc Natl Acad Sci* **107**: 22534–22539. doi:10.1073/pnas.0913805107
- Moses AM, Pollard DA, Nix DA, Iyer VN, Li XY, Biggin MD, Eisen MB. 2006. Large-scale turnover of functional transcription factor binding sites in

- Drosophila*. *PLoS Comput Biol* **2**: e130. doi:10.1371/journal.pcbi.0020130
- Mustonen V, Kinney J, Callan CG, Lässig M. 2008. Energy-dependent fitness: a quantitative model for the evolution of yeast transcription factor binding sites. *Proc Natl Acad Sci* **105**: 12376–12381. doi:10.1073/pnas.0805909105
- Paris M, Kaplan T, Li XY, Villalta JE, Lott SE, Eisen MB. 2013. Extensive divergence of transcription factor binding in *Drosophila* embryos with highly conserved gene expression. *PLoS Genet* **9**: e1003748. doi:10.1371/journal.pgen.1003748
- Park D, Morris AR, Battenhouse A, Iyer VR. 2014. Simultaneous mapping of transcript ends at single-nucleotide resolution and identification of widespread promoter-associated non-coding RNA governed by TATA elements. *Nucleic Acids Res* **42**: 3736–3749. doi:10.1093/nar/gkt1366
- Pelechano V, Wei W, Steinmetz LM. 2013. Extensive transcriptional heterogeneity revealed by isoform profiling. *Nature* **497**: 127–131. doi:10.1038/nature12121
- Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**: 841–842. doi:10.1093/bioinformatics/btq033
- R Core Team. 2013. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna. <https://www.R-project.org/>.
- Reddy TE, Gertz J, Pauli F, Kucera KS, Varley KE, Newberry KM, Marinov GK, Mortazavi A, Williams BA, Song L, et al. 2012. Effects of sequence variation on differential allelic transcription factor occupancy and gene expression. *Genome Res* **22**: 860–869. doi:10.1101/gr.131201.111
- Rhee HS, Pugh BF. 2011. Comprehensive genome-wide protein-DNA interactions detected at single-nucleotide resolution. *Cell* **147**: 1408–1419. doi:10.1016/j.cell.2011.11.013
- Rossi MJ, Lai WKM, Pugh BF. 2018. Genome-wide determinants of sequence-specific DNA binding of general regulatory factors. *Genome Res* **28**: 497–508. doi:10.1101/gr.229518.117
- Rowan S, Siggers T, Lachke SA, Yue Y, Bulyk ML, Maas RL. 2010. Precise temporal control of the eye regulatory gene *Pax6* via enhancer-binding site affinity. *Genes Dev* **24**: 980–985. doi:10.1101/gad.1890410
- Sandelin A, Alkema W, Engström P, Wasserman WW, Lenhard B. 2004. JASPAR: an open-access database for eukaryotic transcription factor binding profiles. *Nucleic Acids Res* **32**: D91–D94. doi:10.1093/nar/gkh012
- Scannell DR, Zill OA, Rokas A, Payen C, Dunham MJ, Eisen MB, Rine J, Johnston M, Hittinger CT. 2011. The awesome power of yeast evolutionary genetics: new genome sequences and strain resources for the *Saccharomyces sensu stricto* genus. *G3 (Bethesda)* **1**: 11–25. doi:10.1534/g3.111.000273
- Scardigli R, Bäumer N, Gruss P, Guillemot F, Le Roux I. 2003. Direct and concentration-dependent regulation of the proneural gene *Neurogenin2* by *Pax6*. *Development* **130**: 3269–3281. doi:10.1242/dev.00539
- Schmidt D, Wilson MD, Ballester B, Schwalie PC, Brown GD, Marshall A, Kutter C, Watt S, Martinez-Jimenez CP, Mackay S, et al. 2010. Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription factor binding. *Science* **328**: 1036–1040. doi:10.1126/science.1186176
- Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier LW, Richards S. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res* **15**: 1034–1050. doi:10.1101/gr.3715005
- Skene PJ, Henikoff S. 2017. An efficient targeted nuclease strategy for high-resolution mapping of DNA binding sites. *eLife* **6**: e21856. doi:10.7554/eLife.21856
- Spivakov M, Akhtar J, Kheradpour P, Beal K, Girardot C, Koscielny G, Herrero J, Kellis M, Furlong EEM, Birney E. 2012. Analysis of variation at transcription factor binding sites in *Drosophila* and humans. *Genome Biol* **13**: R49. doi:10.1186/gb-2012-13-9-r49
- Stefflova K, Thybert D, Wilson MD, Streeter I, Aleksic J, Karagianni P, Brazma A, Adams DJ, Talianidis I, Marioni JC, et al. 2013. Cooperativity and rapid evolution of cobound transcription factors in closely related mammals. *Cell* **154**: 530–540. doi:10.1016/j.cell.2013.07.007
- Sun Y, Nien CY, Chen K, Liu HY, Johnston J, Zeitlinger J, Rushlow C. 2015. Zelda overcomes the high intrinsic nucleosome barrier at enhancers during *Drosophila* zygotic genome activation. *Genome Res* **25**: 1703–1714. doi:10.1101/gr.192542.115
- Suter DM. 2020. Transcription factors and DNA play hide and seek. *Trends Cell Biol* **30**: 491–500. doi:10.1016/j.tcb.2020.03.003
- Tanay A. 2006. Extensive low-affinity transcriptional interactions in the yeast genome. *Genome Res* **16**: 962–972. doi:10.1101/gr.5113606
- Tareen A, Kinney JB. 2020. Logomaker: beautiful sequence logos in Python. *Bioinformatics* **36**: 2272–2274. doi:10.1093/bioinformatics/btz921
- Teixeira MC, Monteiro P, Jain P, Tenreiro S, Fernandes AR, Mira NP, Alenquer M, Freitas AT, Oliveira AL, Sá-Correia I. 2006. The YEASTRACT database: a tool for the analysis of transcription regulatory associations in *Saccharomyces cerevisiae*. *Nucleic Acids Res* **34**: D446–D451. doi:10.1093/nar/gkj013
- Tirosh I, Reikavav S, Levy AA, Barkai N. 2009. A yeast hybrid provides insight into the evolution of gene expression regulation. *Science* **324**: 659–662. doi:10.1126/science.1169766
- Tirosh I, Sigal N, Barkai N. 2010. Divergence of nucleosome positioning between two closely related yeast species: genetic basis and functional consequences. *Mol Syst Biol* **6**: 365. doi:10.1038/msb.2010.20
- Tornow J, Zeng X, Gao W, Santangelo GM. 1993. GCR1, a transcriptional activator in *Saccharomyces cerevisiae*, complexes with RAP1 and can function without its DNA binding domain. *EMBO J* **12**: 2431–2437. doi:10.1002/j.1460-2075.1993.tb05897.x
- Tsankov AM, Thompson DA, Socha A, Regev A, Rando OJ. 2010. The role of nucleosome positioning in the evolution of gene regulation. *PLoS Biol* **8**: e1000414. doi:10.1371/journal.pbio.1000414
- Venkataram S, Fay JC. 2010. Is transcription factor binding site turnover a sufficient explanation for cis-regulatory sequence divergence? *Genome Biol Evol* **2**: 851–858. doi:10.1093/gbe/evq066
- Weiss CV, Roop JJ, Hackley RK, Chuong JN, Grigoriev IV, Arkin AP, Skerker JM, Brem RB. 2018. Genetic dissection of interspecific differences in yeast thermotolerance. *Nat Genet* **50**: 1501–1504. doi:10.1038/s41588-018-0243-4
- Wilson MD, Barbosa-Morais NL, Schmidt D, Conboy CM, Vanes L, Tybulewicz VJ, Fisher EMC, Tavaré S, Odom DT. 2008. Species-specific transcription in mice carrying human chromosome 21. *Science* **322**: 434–438. doi:10.1126/science.1160930
- Wong ES, Schmitt BM, Kazachenka A, Thybert D, Redmond A, Connor F, Rayner TF, Feig C, Ferguson-Smith AC, Marioni JC, et al. 2017. Interplay of cis and trans mechanisms driving transcription factor binding and gene expression evolution. *Nat Commun* **8**: 1092. doi:10.1038/s41467-017-01037-x
- Yang MG, Ling E, Cowley CJ, Greenberg ME, Vierbuchen T. 2021. Characterization of sequence determinants of enhancer function using natural genetic variation. *bioRxiv* doi:10.1101/2021.12.17.473050
- Yue JX, Li J, Aigrain L, Hallin J, Persson K, Oliver K, Bergström A, Coupland P, Warringer J, Lagomarsino MC, et al. 2017. Contrasting evolutionary genome dynamics between domesticated and wild yeasts. *Nat Genet* **49**: 913–924. doi:10.1038/ng.3847
- Zentner GE, Kasinathan S, Xin B, Rohs R, Henikoff S. 2015. ChEC-seq kinetics discriminates transcription factor binding sites by DNA sequence and shape *in vivo*. *Nat Commun* **6**: 8733. doi:10.1038/ncomms9733
- Zheng W, Zhao H, Mancera E, Steinmetz LM, Snyder M. 2010. Genetic analysis of variation in transcription factor binding in yeast. *Nature* **464**: 1187–1191. doi:10.1038/nature08934
- Zhou T, Shen N, Yang L, Abe N, Horton J, Mann RS, Bussemaker HJ, Gordán R, Rohs R. 2015. Quantitative modeling of transcription factor binding specificities using DNA shape. *Proc Natl Acad Sci* **112**: 4654–4659. doi:10.1073/pnas.1422023112

Received February 28, 2022; accepted in revised form May 20, 2022.



Evolution of transcription factor binding through sequence variations and turnover of binding sites

Gat Krieger, Offir Lupo, Patricia Wittkopp, et al.

Genome Res. 2022 32: 1099-1111 originally published online May 26, 2022

Access the most recent version at doi:[10.1101/gr.276715.122](https://doi.org/10.1101/gr.276715.122)

Supplemental Material <http://genome.cshlp.org/content/suppl/2022/06/14/gr.276715.122.DC1>

References This article cites 79 articles, 26 of which can be accessed free at:
<http://genome.cshlp.org/content/32/6/1099.full.html#ref-list-1>

Creative Commons License This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>
