

# Predicting dynamic expression patterns in budding yeast with a fungal DNA language model

Kuan-Hao Chao<sup>1,2,\*</sup>, Majed Mohamed Magzoub<sup>3</sup>, Emily Stoops<sup>3</sup>, Sean Hackett<sup>3</sup>, Johannes Linder<sup>3,\*</sup> and David R. Kelley<sup>3,\*</sup>

<sup>1</sup>Department of Computer Science, Johns Hopkins University, Baltimore, MD 21218, USA

<sup>2</sup>Center for Computational Biology, Johns Hopkins University, Baltimore, MD 21218, USA

<sup>3</sup>Calico Life Sciences LLC, South San Francisco, CA 94080, USA

\*Corresponding authors: [kuanhao.chao@gmail.com](mailto:kuanhao.chao@gmail.com), [jlinder@calicolabs.com](mailto:jlinder@calicolabs.com), [drk@calicolabs.com](mailto:drk@calicolabs.com)

## Abstract

Predicting gene expression from DNA sequence remains challenging due to complex regulatory codes. We introduce a masked DNA language model pretrained on 165 fungal genomes closely related to budding yeast that captures conserved regulatory grammar. Fine-tuning the LM on yeast RNA-seq data—including high-resolution transcriptional regulator induction time courses generated in this study—yielded Shorkie, a model that substantially improves gene expression prediction compared to baselines trained without self-supervision. Shorkie identified canonical transcription factor (TF) binding motifs and tracked their usage across induction experiments. Furthermore, Shorkie accurately predicted variant effects, outperforming leading sequence-to-expression models in *cis*-eQTL classification and achieving high concordance with massively parallel reporter assays. Interpretability analyses revealed Shorkie’s ability to resolve promoter dynamics, splicing signals, and temporal changes in regulatory motif usage. This framework demonstrates that evolutionary-scale pretraining combined with transfer learning substantially improves our ability to decode gene regulation from sequence, providing insights into noncoding variants and regulatory networks.

## 26 Introduction

27 Predicting gene expression levels from DNA sequence is a fundamental challenge in genomics with broad  
28 implications for understanding gene regulation and disease. *Saccharomyces cerevisiae* (budding yeast) has  
29 served as the premier model for eukaryotic gene regulation, with ~7,000 genes controlled by hundreds of  
30 transcription factors (TFs). Despite decades of work mapping *cis*-regulatory motifs and their regulators<sup>1-11</sup>,  
31 quantitative prediction of gene expression from regulatory sequences remains limited. Even sophisticated  
32 machine learning models explain at most ~73% of expression variance and rely on hand-crafted rules for  
33 motif spacing, orientation, and combinatorial logic<sup>12,13</sup>. This gap highlights the complexity of the regulatory  
34 code and motivates new computational approaches.

35 Supervised deep learning can learn directly from sequence without hand-crafted features<sup>14,15</sup>, but faces a  
36 fundamental limitation in yeast: the compact 12 Mb genome provides insufficient training examples, predis-  
37 posing models to overfitting. Self-supervised DNA language models (LMs) overcome this limitation by  
38 learning rich sequence representations from many unlabeled genomes. Models such as DNABERT<sup>16,17</sup>,  
39 Evo<sup>18</sup>, and others<sup>19-27</sup> demonstrate that masked-token prediction captures conserved regulatory syntax and  
40 the locations of genes. Parallel advances in protein LMs<sup>28-32</sup> further validate self-supervised pretraining for  
41 extracting functional patterns from sequence.

42 Despite hundreds of high-quality fungal genomes now available<sup>33</sup>, gene expression data exist for only a  
43 handful of species, precluding supervised pan-fungal training. Masked DNA LMs circumvent this limitation:  
44 by predicting masked bases, they capture major promoter motifs and the locations of genes without labels.  
45 Models pretrained on related species generalize better than those trained on single genomes<sup>17</sup>.

46 Here, we leverage this paradigm to improve yeast expression modeling. We first pretrained a bidirectional  
47 masked LM on diverse fungal genomes using a BERT-style objective. We then fine-tuned this model on  
48 high-resolution RNA-seq time courses from transcriptional regulator induction experiments in *S. cere-*  
49 *visiae*<sup>34</sup>, plus publicly available epigenomic and transcriptomic data, creating a high-quality expression pre-  
50 dictor called **Shorkie**.

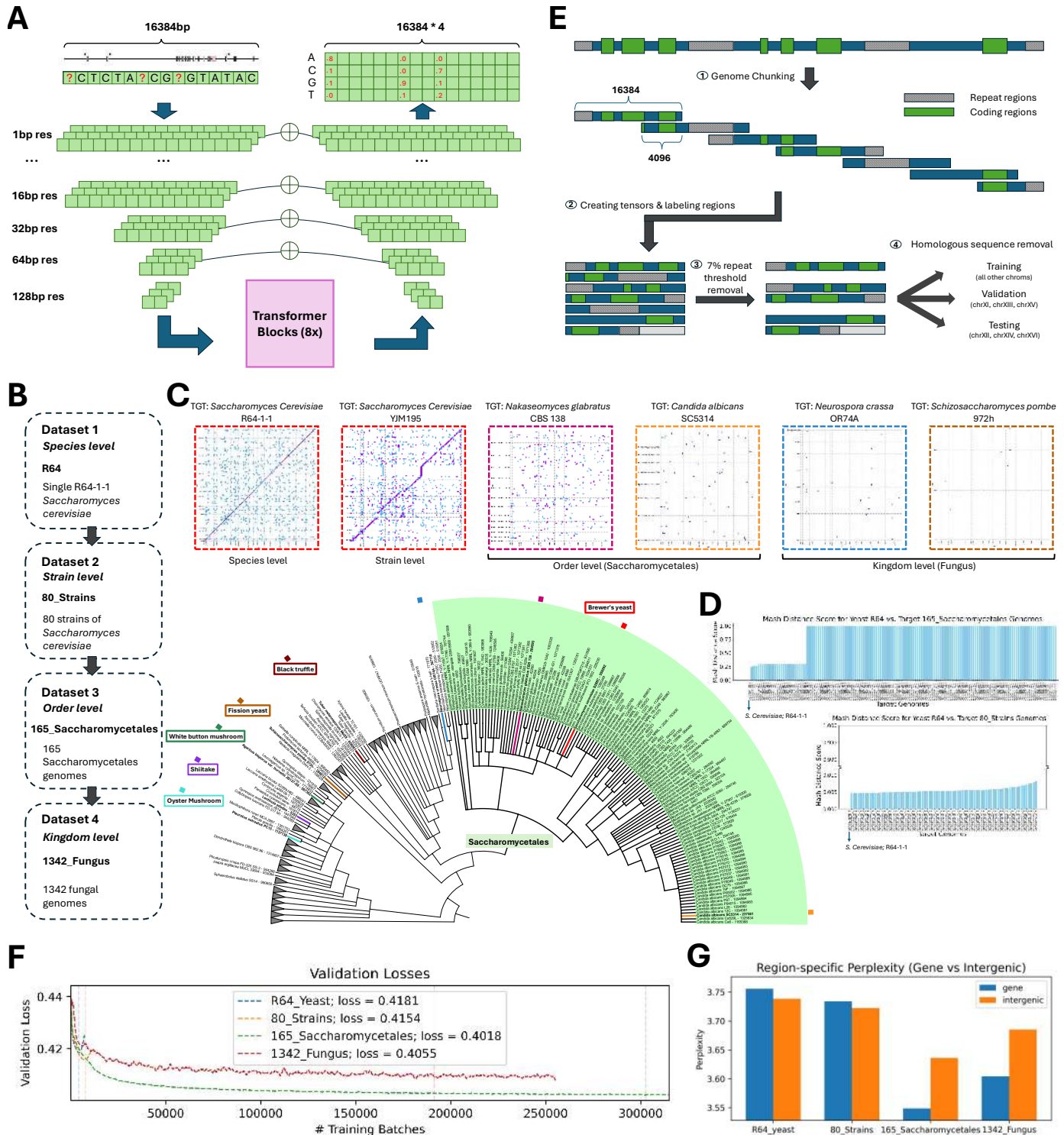
51 By combining evolutionary pretraining with transfer learning, Shorkie outperforms models trained without  
52 self-supervision in predicting expression and temporal dynamics of held-out genes. Furthermore, Shorkie  
53 delivers robust variant effect predictions in both *cis*-eQTL classification and massively parallel reporter as-  
54 says. These findings demonstrate that masked language modeling across diverse fungal genomes, coupled  
55 with transfer learning, provides a powerful framework for quantitative gene regulation modeling and  
56 noncoding variant interpretation in yeast.

## 57 Results

### 58 Yeast language model design and training across evolutionary divergences

59 We trained a masked DNA LM on more than 1,300 fungal genomes from Ensembl Fungi, which constitutes  
60 the same training data as the Species-aware LM developed by Karollus et al. (2024)<sup>22</sup> (Figure 1A). The  
61 model architecture integrates elements from Enformer<sup>35</sup> and Borzoi<sup>36</sup>, employing a convolutional tower with

62 subsampling followed by eight self-attention blocks operating at 128 bp resolution. We repeated Borzoi's  
 63 U-Net upsampling block<sup>37</sup> seven times to progressively restore single-nucleotide resolution for masked token  
 64 prediction (detailed model configurations in Figure S1; Methods). This flexible architecture enables fine-  
 65 tuning at coarser resolutions by removing U-net blocks.



66  
 67 **Figure 1.** Overview of datasets, preprocessing pipeline, model architecture, and performance metrics for the fungal  
 68 language model (Shorkie LM). (A) Schematic of the Shorkie LM architecture. (B) Four datasets employed: single *S.*  
 69 *cerevisiae* genome (R64\_yeast, species-level), 80 *S. cerevisiae* strains (80\_strains, strain-level), 165

70 Saccharomycetales genomes (165\_Saccharomycetales, order-level) with Saccharomycetales highlighted in light green,  
71 and 1,341 fungal genomes spanning the kingdom (1341\_Fungal, kingdom-level), including common mushrooms such  
72 as oyster mushroom (*Pleurotus ostreatus*), shiitake (*Lentinula edodes*), white button mushroom (*Agaricus bisporus*),  
73 and black truffle (*Tuber melanosporum*), as well as fission yeast (*Schizosaccharomyces pombe*), and brewer's yeast  
74 (*Saccharomyces cerevisiae*). (C) Representative genome distance dot plots for selected genomes from each dataset,  
75 with the x-axis representing the R64 *S. cerevisiae* genome and the y-axis representing the comparison genome. (D)  
76 Mash distance between R64 *S. cerevisiae* genome and genomes in 80\_strains and 165\_Saccharomycetales. (E) Data  
77 preprocessing pipeline converting raw genomic data into tensors and labels for Shorkie LM training, validation, and  
78 testing. (F) Validation loss progression across training steps. (G) Comparison of test set perplexity in genic and inter-  
79 genic regions across four model variants.

80 To identify optimal training genomes for *S. cerevisiae*, we prepared four datasets with varying evolutionary  
81 divergence (Figure 1B): (1) R64: *S. cerevisiae* reference; (2) 80\_strains: 80 *S. cerevisiae* strains; (3)  
82 165\_Saccharomycetales: 165 genomes from the Saccharomycetales order; (4) 1341\_Fungus: 1,341 fungal  
83 kingdom genomes. We inferred phylogenetic relationships with ETE3 using NCBI Taxonomy<sup>38</sup> and visual-  
84 ized in iTOL<sup>39,40</sup> (Figure 1B; see Figure S2 for full tree; Methods). We quantified divergence from *S. cere-*  
85 *visiae* R64 using MUMmer dot plots<sup>41</sup> (Figure 1C) and Mash distances<sup>42</sup> (Figure 1D). Closely related strains  
86 (e.g., YJM195 at Mash  $\approx 0.01$ ) exhibited near-continuous synteny, whereas taxa such as *C. albicans* and *N.*  
87 *glabratus* showed fragmentation (Mash 0.25–1). Distant outgroups (*S. pombe*, *N. crassa*) yielded negligible  
88 alignments (Figure 1C).

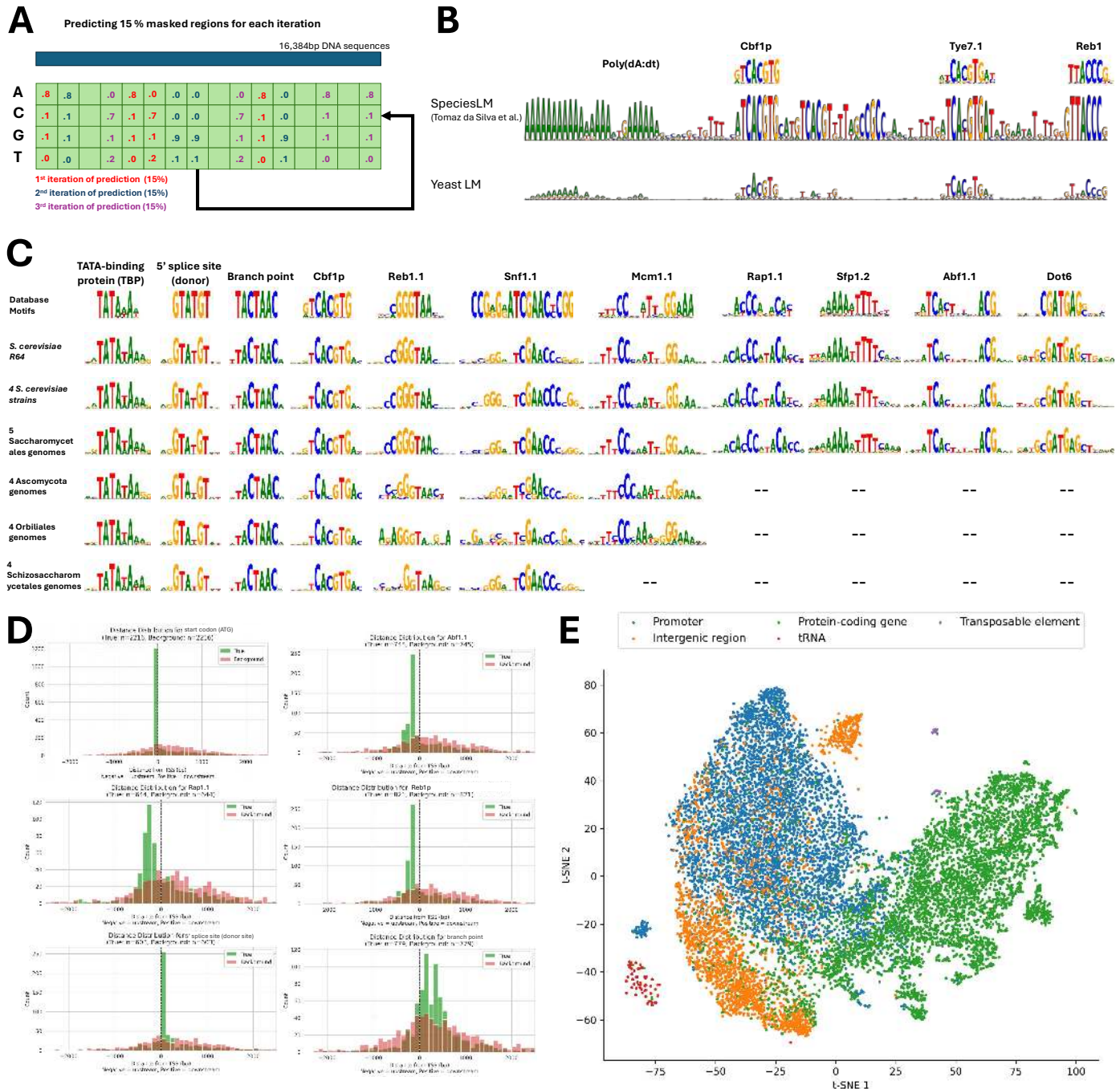
89 We prepared genomes for LM training by masking repetitive elements<sup>43–48</sup> (Figure S3A–B; Methods), seg-  
90 menting into overlapping 16,384 bp windows with 4,096 bp stride, and excluding windows with >7% repet-  
91 itive content. We assessed gene count per window (Figure S4A–B), coding-to-noncoding ratios per window  
92 (Figure S4C–D), repetitive region distribution (Figure S4E–F), and gene annotation completeness (Figure  
93 S4G–I; Methods). To focus learning on regulatory sequences, we down-weighted the loss function by 0.1 at  
94 coding (72% of *S. cerevisiae* R64) and repetitive (7.39%) positions (Figure S3C; Methods). Training/vali-  
95 dation/test sets were split by *S. cerevisiae* chromosomes. Validation/test sets included only *S. cerevisiae*; the  
96 training set included all genomes after removing sequences homologous to validation or test sets using min-  
97 imap2 at 20% divergence cutoff<sup>49,50</sup> (Figure 1E; Figure S3D–E; Methods).

98 Evolutionary divergence correlated with training complexity, reflected in higher training loss. The 165\_Sac-  
99 charomycetales model achieved the lowest validation loss, outperforming the more divergent 1341\_Fungus  
100 and avoiding the overfitting observed with 80\_strains and R64 (Figure 1F). This optimal performance ex-  
101 tended to the test set (Figure 1G) and held across alternative architectures (two residual CNN baselines and  
102 a larger U-Net-transformer) (Figure S5B–D, Figure S6). The 165\_Saccharomycetales-trained model, here-  
103 after “Shorkie LM”, thus represents the optimal evolutionary scale for *S. cerevisiae* generalization.

## 104 **Shorkie LM captures regulatory conservation and generalizes across diverse fungi**

105 Transcription factor (TF) binding motifs are fundamental regulatory units<sup>51</sup>, and previous studies demon-  
106 strate that masked DNA LMs learn subtle motif co-occurrence patterns<sup>22</sup>. We evaluated Shorkie LM's motif  
107 identification by segmenting the *S. cerevisiae* genome (Figure 1B), randomly masking 15% of bases, and  
108 iteratively imputing them to reconstruct position probability matrices (Figure 2A; Methods). In the SMT3

109 promoter, Shorkie LM identified canonical motifs including Poly(dA:dT), Cbf1, Tye7, and Reb1, consistent  
 110 with prior analyses<sup>52</sup> (Figure 2B). This alignment-free approach enables flexible sequence probability deri-  
 111 vation<sup>19</sup>.



112  
 113 **Figure 2. Shorkie LM identifies conserved transcription factor binding motifs across fungal genomes. (A)**  
 114 **Position probability matrix (PPM) reconstruction from DNA sequences using the fungal language model Shorkie LM. (B)**  
 115 **Comparative analysis of SMT3 promoter predictions (chrIV:1,469,090-1,469,198) by Shorkie LM and the Species-**  
 116 **aware DNA LM<sup>22,52</sup>, highlighting key motifs including poly(dA:dT), Cbf1, Tye7, and Reb1. (C) Summary of known**  
 117 **motifs detected by Shorkie LM across six datasets: (1) reference *S. cerevisiae* genome; (2) four randomly selected *S.***  
 118 ***cerevisiae* strains; (3) five genomes from the Saccharomycetales order; (4) four genomes from the Ascomycota phy-**  
 119 **lum; (5) four genomes from the Orbiliales order; and (6) four genomes from the Schizosaccharomycetales order.**

120 TF-MoDISco-identified motifs include TATA-binding protein, 5' splice site (donor), branch point, Cbflp, Reb1.1,  
121 Snf1.1, Mcm1.1, Rap1.1, Sfp1.2, Abf1.1 and Dot6. **(D)** Histograms depicting enrichment of TF-MoDISco-identified  
122 motifs upstream of transcription start sites (TSS) relative to background distributions in *S. cerevisiae*, and enrichment  
123 of 5' splice sites (donors) and branch points within genic regions. **(E)** t-SNE embeddings of different genomic elements  
124 from the first self-attention layer of Shorkie LM.

125 We assessed Shorkie LM across six fungal datasets spanning different evolutionary distances (see Methods).  
126 Following prediction, we employed TF-MoDISco-lite<sup>53,54</sup> for *de novo* motif clustering and matched clusters  
127 to yeast motif databases<sup>55–60</sup> (Methods). Motif conservation varied across evolutionary distance. Shorkie  
128 recovered core regulatory motifs and features, including TATA-binding protein (TBP)/TATA elements, start  
129 codons, splice sites, and TF binding sites such as Cbfl, Reb1, and Snf1. Mcm1.1 was absent in Schizosac-  
130 charomycetales, which lack a direct homolog; this role is instead served by the functionally analogous  
131 MADS-box transcription factor Map1<sup>61–64</sup>. Motif conservation declined beyond the Saccharomycetales order,  
132 consistent with model training metrics (Figure 1G; Figure 2C). See Supplemental Figures S7 and S8 for  
133 comprehensive motif discovery results.

134 To validate biological relevance, we mapped TF-MoDISco-derived motifs onto the *S. cerevisiae* genome,  
135 assigning motifs to their nearest genes and computing transcription start site (TSS) distances. Compared  
136 with random controls, TBP, Cbflp, Reb1.1, Mcm1.1, and Snf1.1 showed promoter enrichment (Figure 2D;  
137 Figure S7). The 5' splice donor site localized downstream of TSSs, while branch points distributed broadly  
138 within genes. Five randomly selected Saccharomycetales genomes produced similar enrichments (Figure  
139 S8). Additionally, Shorkie LM's first attention layer effectively differentiated genomic features by embed-  
140 ding patterns (Figure 2E), as did subsequent layers (Figure S9).

141 These results demonstrate that Shorkie LM captures conserved regulatory grammar, accurately identifies  
142 motifs, and generalizes across substantial evolutionary distances.

### 143 **Shorkie: LM transfer learning enables improved gene expression prediction**

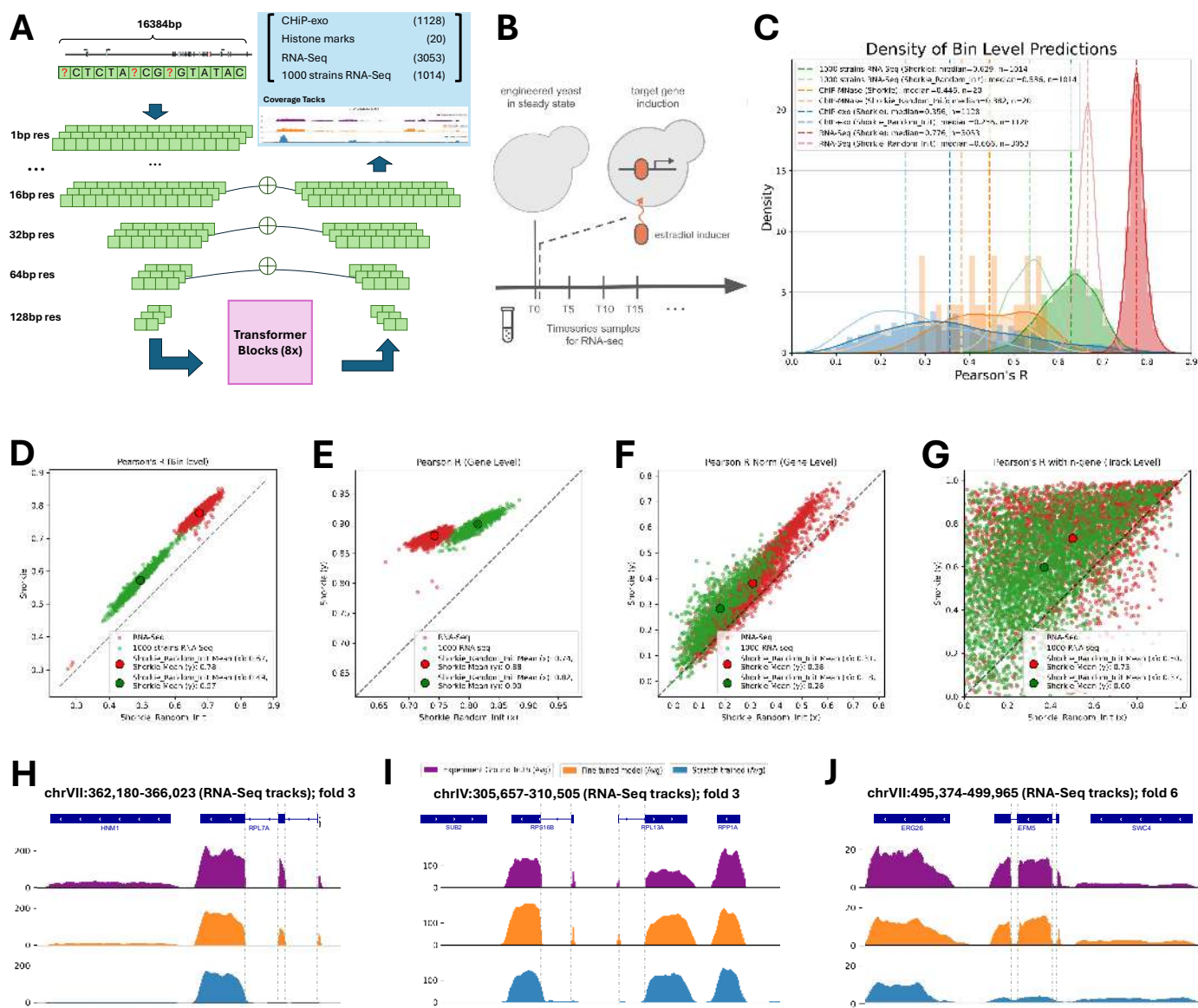
144 Building on the strong self-supervised sequence foundation, we developed Shorkie, a supervised model pre-  
145 dicting RNA-seq and ChIP-exo/MNase aligned coverage tracks at 16 bp resolution from DNA sequences.  
146 Starting with the LM architecture, we removed its final four upsampling layers and added task-specific out-  
147 put heads (Figure 3A; complete configurations in Figure S10; Methods). We curated 2,162 experimental  
148 tracks for training: 1,128 ChIP-exo<sup>65</sup>, 20 ChIP-MNase<sup>65</sup>, and 1,014 RNA-seq datasets from various yeast  
149 isolates<sup>66</sup>.

150 In addition, we generated 3,053 new high-resolution induction RNA-seq timepoints using a protocol adapted  
151 from Hackett et al.<sup>34</sup> (Figure 3B), bringing the total to 5,215 experimental tracks. Key protocols and descrip-  
152 tion of chemostats (ministat array), culture conditions, library preparation, data pre-processing and quality  
153 controls are detailed in Methods (Figure S11–S12).

154 To evaluate the impact of pretraining, we compared transfer learning from Shorkie LM against random ini-  
155 tialization (Shorkie\_Random\_Init) across eight-fold cross-validation (Figure S13). On the transcriptional  
156 regulator induction RNA-seq test data, Shorkie achieved median bin-level Pearson's R of 0.78 versus 0.67

157 for Shorkie\_Random\_Init, shifting the correlation distribution upward (Figure 3C) and boosting per-track  
 158 correlations (Figure 3D).

159 Gene-level aggregation across exon-overlapping bins (Methods) yielded mean Pearson's R of 0.88 for  
 160 Shorkie versus 0.74 for Shorkie\_Random\_Init (Figure 3E). Normalized gene-level correlations (quan-  
 161 tile-normalized per experiment and mean-centered per gene) confirmed this advantage (Figure 3F; Methods).  
 162 After averaging track-specific performance across tracks for each gene, Shorkie exceeded Shorkie\_Ran-  
 163 dom\_Init in 87.8% of genes, particularly at higher expression levels (Figure 3G; Figure S14L; Methods).



164  
 165 **Figure 3.** Shorkie architecture and RNA-seq prediction performance across multiple scales. **(A)** Shorkie architecture:  
 166 U-Net model with eight transformer blocks. All layers inherit pretrained Shorkie LM weights; task-specific output  
 167 heads (blue) predict perturbation timepoint RNA-seq (n = 3,053), 1000-strain RNA-seq (n = 1,014), ChIP-exo (n =  
 168 1,128) and ChIP-MNase histone marks (n = 20). **(B)** Yeast cells were grown to steady state, as determined by culture  
 169 density, prior to addition of b-estradiol to the culture and subsequent sampling. **(C)** Distribution of bin-level Pearson's  
 170 R on held-out test data for each track type, comparing Shorkie and Shorkie\_Random\_Init. **(D-G)** Scatter plots com-  
 171 paring Shorkie and Shorkie\_Random\_Init for RNA-seq tracks at **(D)** bin-level Pearson's R; **(E)** gene-level Pearson's

172 R; (F) quantile-normalized and mean-centered gene-level Pearson's R; (G) gene-by-gene, track-level Pearson's R.  
173 (H-J) RNA-seq coverage snapshots of *S. cerevisiae* test set gene loci: (H) chrVII:362,180–366,023 (RPL7A); (I)  
174 chrIV:305,657–310,505 (RPS16B and RPL13A); and (J) chrVII:495,374–499,965 (EFM5).

175 To understand the model's regulatory focus, we analyzed self-attention weights centered on the three-exon  
176 gene EFM5 and housekeeping gene RPL7A (Methods). Attention from both the pretrained LM and Shorkie  
177 highlighted genic and discrete intergenic regulatory regions (putative promoters), whereas Shorkie\_Ran-  
178 dom\_Init produced diffuse attention patterns (Figure S15). Test loci visualization confirmed Shorkie's ac-  
179 curate prediction of intronic coverage drops and expression profiles, contrasting Shorkie\_Random\_Init less  
180 precise predictions (Figure 3H–J).

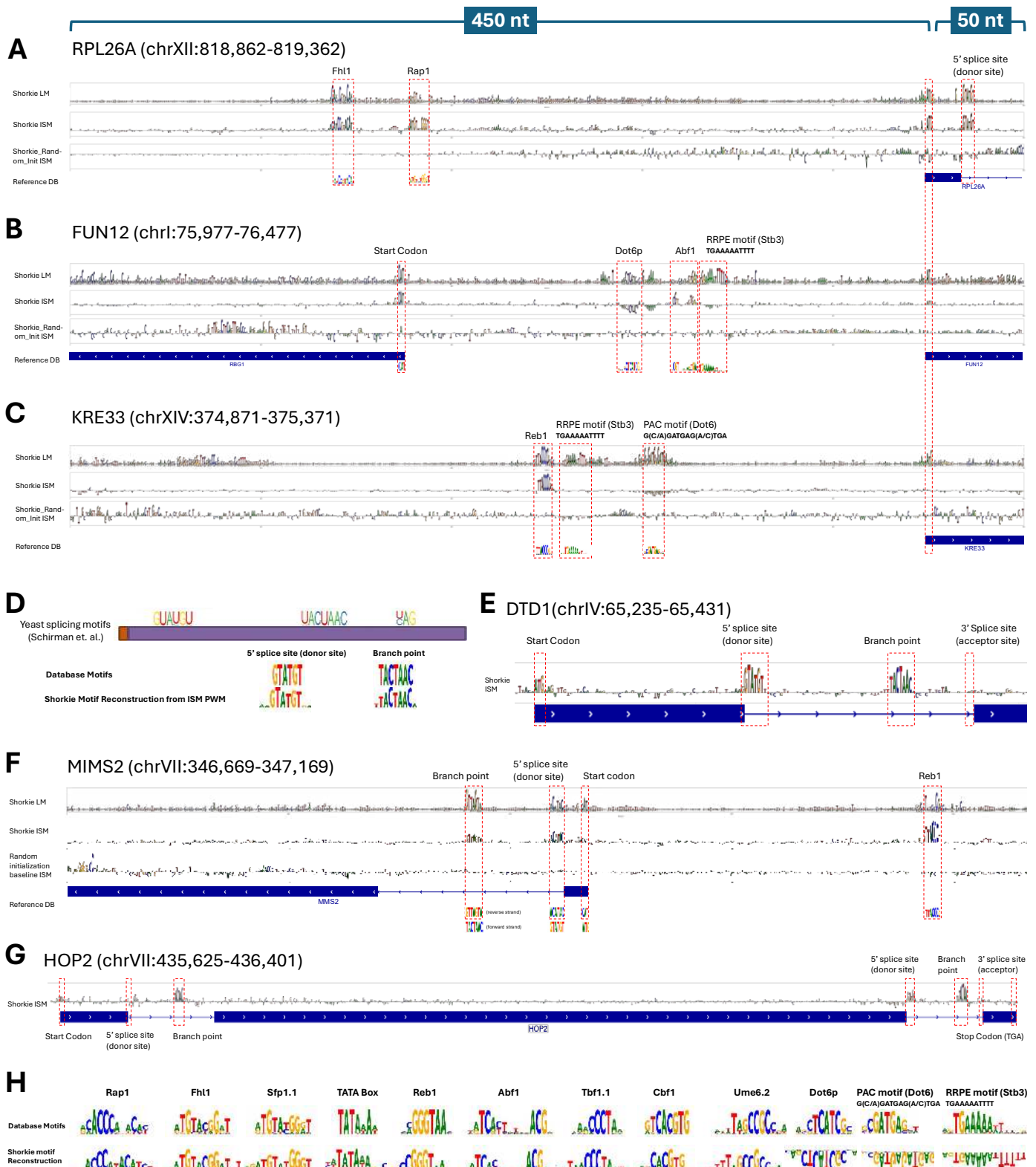
181 These results demonstrate that supervised transfer learning from pan-fungal self-supervised pretraining pro-  
182 vides robust, generalizable representations of exon-intron structure and regulatory grammar, yielding sub-  
183 stantial improvements in expression prediction.

### 184 **Shorkie transfer learning preserves regulatory motif recognition**

185 To identify sequence patterns utilized by Shorkie, we performed *in silico* saturation mutagenesis (ISM) on  
186 500 nt promoter windows (-450 to +50 relative to TSS) for three gene cohorts: 137 ribosomal protein (RP),  
187 64 ribosome and rRNA biosynthesis (RRB), and 3,258 additional protein-coding genes. For Shorkie and  
188 Shorkie\_Random\_Init, we computed ISM importance maps by averaging predictions across that model's  
189 eight cross-validation folds. We then compared these ISM maps to per-base information content derived  
190 from the pretrained Shorkie LM's predictive probability distribution relative to the genomic background.

191 In RP promoters like RPL26A (Figure 4A, Figure S16), both Shorkie LM and Shorkie recovered the fork-  
192 head-binding IFHL motif, typically located ~50–80 bp upstream of the TSS. This motif is recognized by the  
193 winged-helix domain of Fhl1, which upon phosphorylation recruits the coactivator Ifh1, linking TOR/PKA  
194 signaling to transcriptional activation<sup>67–69</sup>. Both models also identified the UASrpg element (Upstream Ac-  
195 tivation Sequence, ribosomal protein genes), bound by the pioneer factor Rap1, which recruits chromatin  
196 remodelers, scaffolds Fhl1–Ifh1 complexes, displaces the +1 nucleosome, and establishes nucleosome-de-  
197 pleted regions essential for preinitiation complex assembly<sup>70–73</sup>.

198 In RRB promoters, such as FUN12 (Figure 4B) and KRE33 (Figure 4C), Shorkie LM captured the RRPE  
199 motif (5'-TGAAAATTTT-3'), bound by the repressor Stb3, and the PAC motif (5'-GCGATGA-  
200 GATGAG-3'), recognized by Dot6/Tod6 repressors. These *cis*-elements coordinate rRNA processing and  
201 ribosome assembly genes during the cell cycle and stress responses<sup>74–79</sup>. Shorkie recapitulated RRPE and  
202 PAC motifs and additionally detected Abf1 and Reb1 binding motifs. See Figure S17 for four additional  
203 protein-coding genes.



204  
205  
206  
207  
208  
209

**Figure 4.** Shorkie uses promoter and splicing motifs learned during pretraining. (A-C) Promoter regions (-450 to +50 bp relative to the TSS; 500 bp total) of RPL26A (chrXII:818,862–819,362), FUN12 (chrI:75,977–76,477), and KRE33 (chrXIV:374,871–375,371). Rows 1-3 show DNA logos from Shorkie LM, and ISM maps from Shorkie (fine-tuned) and Shorkie\_Random\_Init (no self-supervision pretraining). Row 4 shows gene annotations from IGV JS<sup>80</sup>. The Shorkie LM PWMs were generated from PPMs (Figure 2A; Methods), whereas the Shorkie and Shorkie\_Random\_Init

210 ISM maps were produced via an ISM analysis that systematically substituted each nucleotide with the three alterna-  
211 tives. **(D)** Canonical *S. cerevisiae* splicing motifs<sup>81</sup>. **(E-G)** Shorkie ISM maps for splicing motifs in DTD1, MIMS2,  
212 and two-intron gene HOP2. **(H)** TF-MoDISco-identified motifs on Shorkie ISM maps: curated yeast database motifs  
213 (top) and Shorkie-derived motifs (bottom).

214 Within genes, Shorkie LM and Shorkie detected canonical splicing signals (Figure 4D)<sup>81</sup>: in DTD1 (Figure  
215 4E), MMS2 (Figure 4F), and the multi-exon HOP2 (Figure 4G), Shorkie ISM maps delineated the 5' splice  
216 donor and branch-point<sup>81-84</sup>. Models were insensitive to acceptor-site mutations—an observation mirrored  
217 by Tomaz da Silva et al.<sup>52</sup>, which likewise fails to reconstruct 3' acceptor motifs. TF-MoDISco analysis  
218 recovered additional motifs including TATA, Sfp1, Tbf1, Cbf1, and Ume6 (Figure 4H).

219 Across all ISM analyses, Shorkie's ISM maps preserved regulatory motif signatures acquired during lan-  
220 guage model pretraining, whereas Shorkie\_Random\_Init failed to recover key motifs. Thus, Shorkie lever-  
221 ages learned regulatory and genic features to improve predictions.

## 222 **Shorkie captures dynamic *cis*-regulatory motif changes across time-course TF inductions**

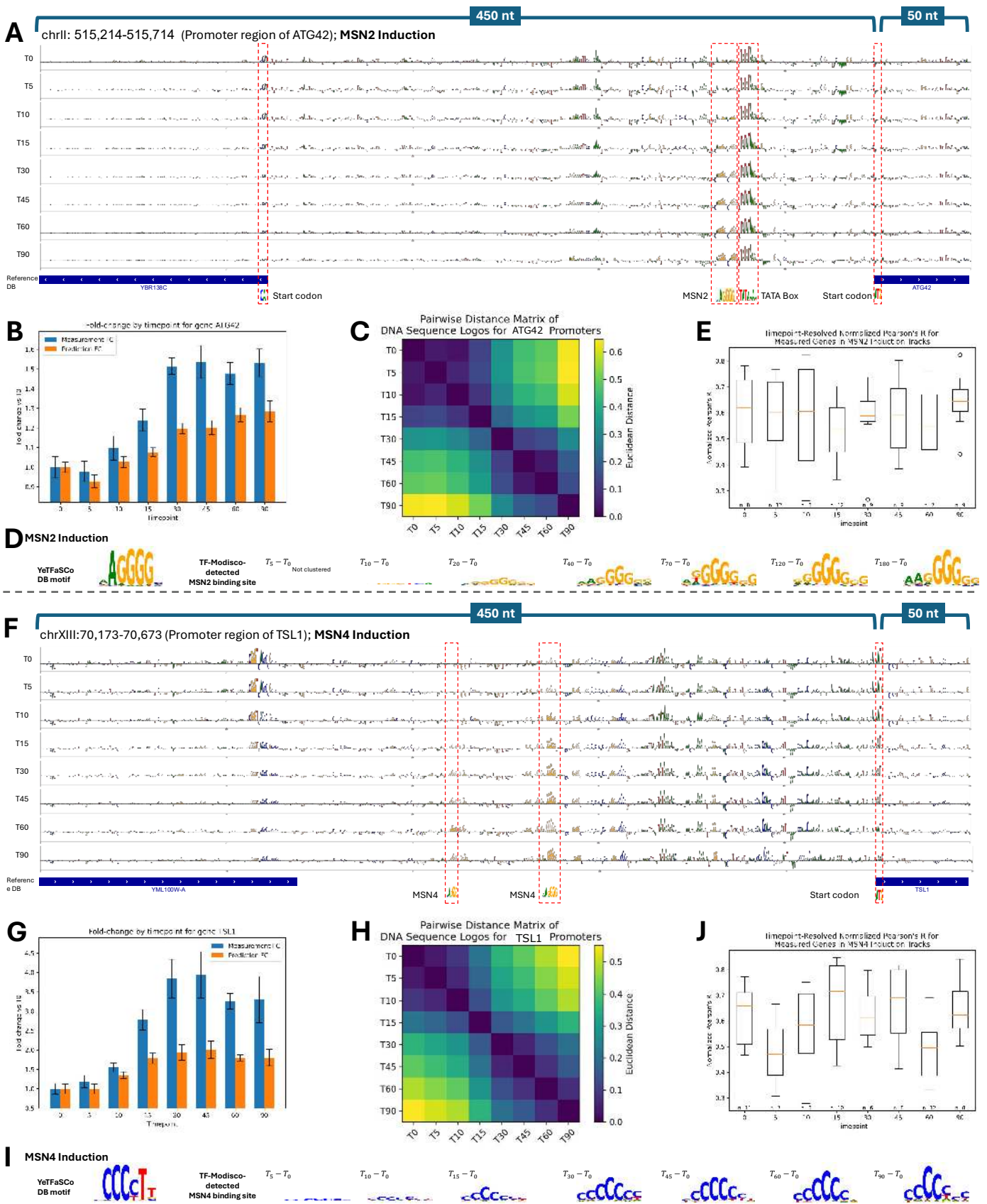
223 Building on Shorkie's ability to identify static *cis*-regulatory elements, we investigated temporal motif usage  
224 during TF induction. We performed ISM across promoters (-450 to +50 nt relative to TSSs) of Saccharo-  
225 myces Genome Database (SGD)-curated TF targets<sup>85</sup> to generate time-resolved maps (Methods).

226 We first examined MSN2, a C2H2 zinc-finger TF that activates ~200 stress-responsive genes via STRE  
227 motifs<sup>86-89</sup>. Shorkie's ISM maps at the ATG42 promoter revealed progressive STRE sharpening over 0-  
228 90 min (Figure 5A), mirroring RNA-seq fold-changes (Figure 5B; Methods). Euclidean distance heatmaps  
229 quantified temporal ISM divergence (Figure 5C), with normalized, mean-centered Pearson's R between ex-  
230 perimental and predicted RNA-seq across MSN2 perturbations ranging from 0.55 to 0.65 (Figure 5D). TF-  
231 MoDISco analysis of  $\Delta T$  ISM maps captured average motif kinetics (Figure 5E; Figure S18A-D show an-  
232 other example at the GLK1 promoter; Methods).

233 We next examined MSN4, an MSN2 paralog rapidly induced upon stress<sup>90</sup>. Shorkie's ISM analysis at the  
234 TSL1 promoter showed similar STRE dynamics, corresponding with RNA-seq fold-changes (Figure 5F-G)  
235 and reflected in Euclidean distance heatmaps (Figure 5H). TF-MoDISco analysis of  $\Delta T$  ISM maps quantified  
236 temporal motif changes, with normalized Pearson's R between experimental and predicted RNA-seq ranging  
237 0.45-0.70 (Figure 5I-J; Figure S18E-H, AYR1 promoter).

238 Finally, we investigated MET4, a bZIP co-activator recruited to E-box motifs (TCACGTG) by cofactors  
239 Cbf1 and Met31/Met32<sup>91,92</sup>. Shorkie's E-box ISM maps inversely correlated with TF induction and were  
240 attenuated by cofactor binding, suggesting capture of cofactor-mediated recruitment rather than direct  
241 MET4-DNA binding (Figure S19; see Discussion).

242 These results demonstrate that Shorkie dynamically tracks *cis*-regulatory motif usage across TF induction  
243 time courses, recapitulating activation kinetics and providing insights into temporal regulatory grammar.



244

245

246

**Figure 5.** Time-course analysis of stress-responsive transcription factor induction. (A-E) MSN2 induction at the ATG42 promoter region (−450 to +50 bp relative to the TSS; chrII:515,214–515,714), sampled at seven timepoints

247 labeled in minutes. **(A)** Shorkie ISM sequence logos: rows correspond to successive timepoints (top to bottom), with  
248 the bottom row showing the reference. Key TF-binding motifs are annotated. **(B)** Experimental fold-change in reads  
249 per million (RPM) (blue) versus Shorkie-predicted signal (orange) across the ATG42 locus at each timepoint. **(C)**  
250 Heatmap of pairwise Euclidean distances between ISM logos, illustrating temporal divergence in motif strength and  
251 composition. **(D)** TF-MoDISco-identified motifs extracted from  $\Delta T$  ISM matrices relative to  $T_0$ . **(E)** Boxplot of  
252 normalized Pearson's R between experimental and predicted profiles across all *S. cerevisiae* genes for MSN2  
253 induction at each timepoint. **(F–J)** MSN4 induction at the TSL1 promoter region (–450 to +50 bp relative to the TSS;  
254 chrXIII:70,173–70,673), with panels analogous to **(A–E)**.

## 255 **Shorkie predicts promoter variant effects validated by MPRAs**

256 Massively parallel reporter assays (MPRAs) provide high-throughput measurements of *cis*-regulatory activ-  
257 ity, enabling regulatory syntax hypothesis testing and variant interpretation. While MPRA data are not ideal  
258 for Shorkie due to its training on large endogenous sequences rather than short reporter constructs, we still  
259 expect reasonable concordance predicting MPRA sequences after marginalizing surrounding context.

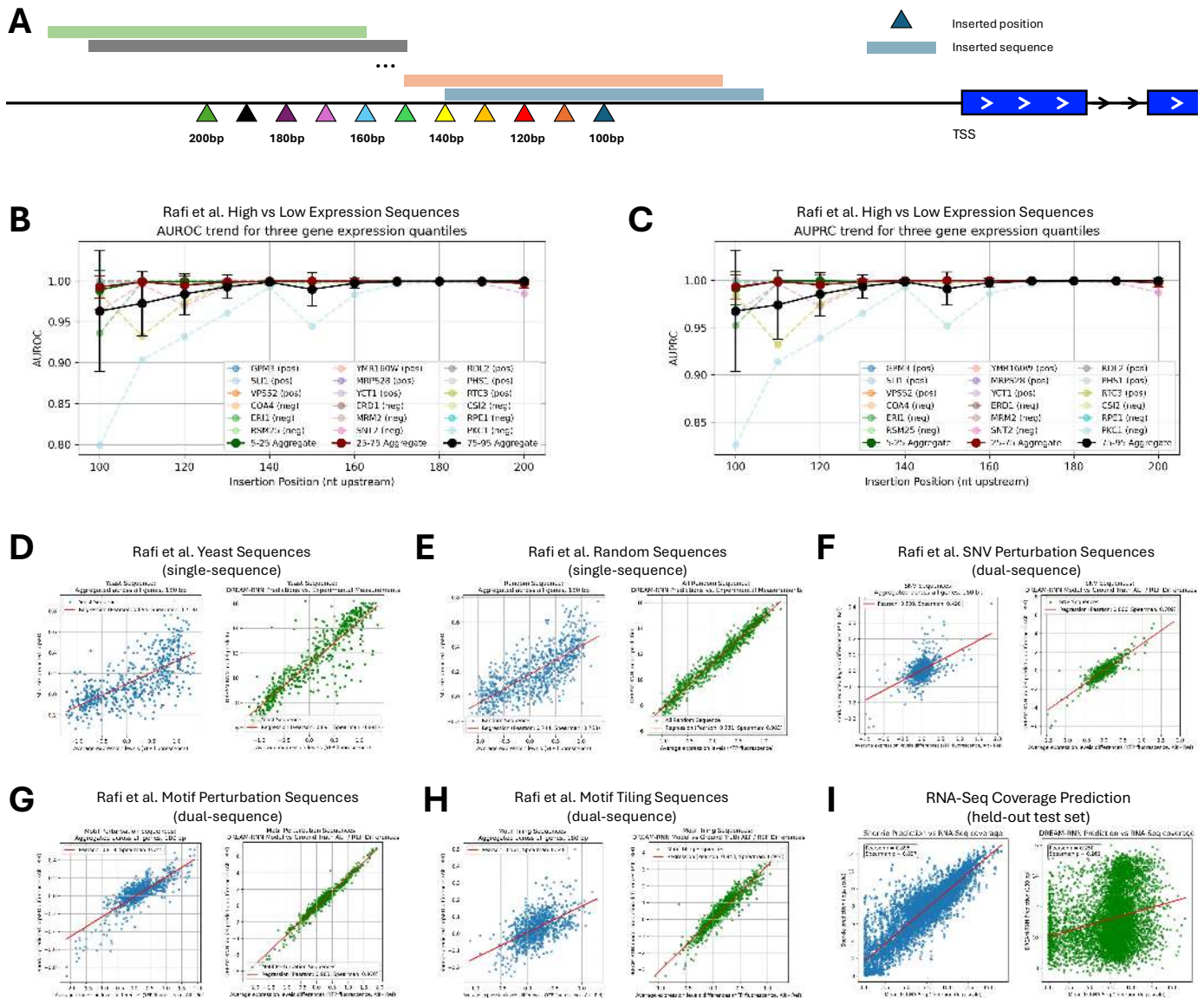
260 We assessed Shorkie using the Random Promoter DREAM Challenge MPRA dataset<sup>93</sup> containing a held-  
261 out set of 71,103 sequences across eight categories: native yeast promoters; random 80-bp oligonucleotides;  
262 high-expression sequences; low-expression sequences; sequences challenging prior models; single-nucleo-  
263 tide variant (SNV) perturbations; motif perturbations; and motif-tiling constructs. Each sequence was as-  
264 sayed in ~100 cells for precise expression estimates<sup>93</sup>. To make marginal predictions with Shorkie, we re-  
265 placed MPRA constructs upstream of TSS for selected “background” genes and averaged predictions across  
266 backgrounds.

267 To characterize positional effects, we selected 10 forward-strand and 12 reverse-strand genes representing  
268 low (5–25th percentile), medium (25–75th percentile), and high (75–95th percentile) pre-induction RNA-  
269 seq expression quantiles. We systematically inserted MPRA sequences at eleven positions, stepping every  
270 10 bp from 200 to 100 bp upstream of the TSS, and quantified regulatory impact as the log fold-change in  
271 downstream gene expression (Figure 6A; Methods).

272 Across tested genes, high-expression sequences yielded positive log fold-change scores, while low-expres-  
273 sion sequences typically produced negative scores (Figure S21A–B). Effect magnitudes modestly increased  
274 with greater TSS distance. Framing high- versus low-expression sequence prediction as binary classification,  
275 Shorkie achieved near-perfect discrimination at each insertion site (AUROC and AUPRC >0.95; Figure 6B–  
276 C; Figure S21C–D). The position 180 bp upstream was selected for subsequent analyses (Figure S20A).

277 Shorkie's marginalized predictions strongly correlated with experimental MPRA measurements: Pearson's  
278 R of 0.70 for native yeast promoters, 0.74 for random sequences, and 0.70 for challenging sequences (Figure  
279 6D–E, Figure S20B–D). For variant-specific categories, comparing reference-alternate differences, correla-  
280 tions were 0.54 for SNV perturbations, 0.82 for motif perturbations, and 0.56 for motif-tiling constructs  
281 (Figure 6F–H). While the best-performing MPRA-trained model, DREAM-RNN outperformed Shorkie on  
282 MPRA predictions, Shorkie's concordance is assuring given zero-shot generalization from endogenous ge-  
283 nome training. Moreover, Shorkie outperformed DREAM models when predicting endogenous gene

284 expression, highlighting that context drives performance differences (Figure 6I). In sum, these results  
 285 demonstrate Shorkie's robust capability to predict promoter-driven expression.



286  
 287 **Figure 6.** Evaluation of Shorkie's predictions of promoter variant effects using MPRA data. (A) Experimental schematic showing MPRA sequences inserted at positions 100–200 bp upstream of the TSS in 10 bp increments for selected yeast genes. (B–C) Classification performance distinguishing high- versus low-expression sequences across upstream insertion sites assessed by (B) AUROC and (C) AUPRC. Genes were stratified into three RNA-seq expression quantiles (5–25%, 25–75%, and 75–95%); dashed colored lines represent individual genes, and black lines depict mean  $\pm$  standard error. (D–E) Comparison between Shorkie predictions (log fold-change scores) and DREAM-RNN model predictions with experimentally measured expression for (D) native yeast sequences and (E) challenging sequences. (F–H) Model performance evaluated for specific regulatory variant sets: (F) single-nucleotide variants (SNV), (G) motif perturbations, and (H) motif tiling constructs. (I) RNA-seq coverage predictions comparing Shorkie and DREAM-RNN against experimentally measured coverage.

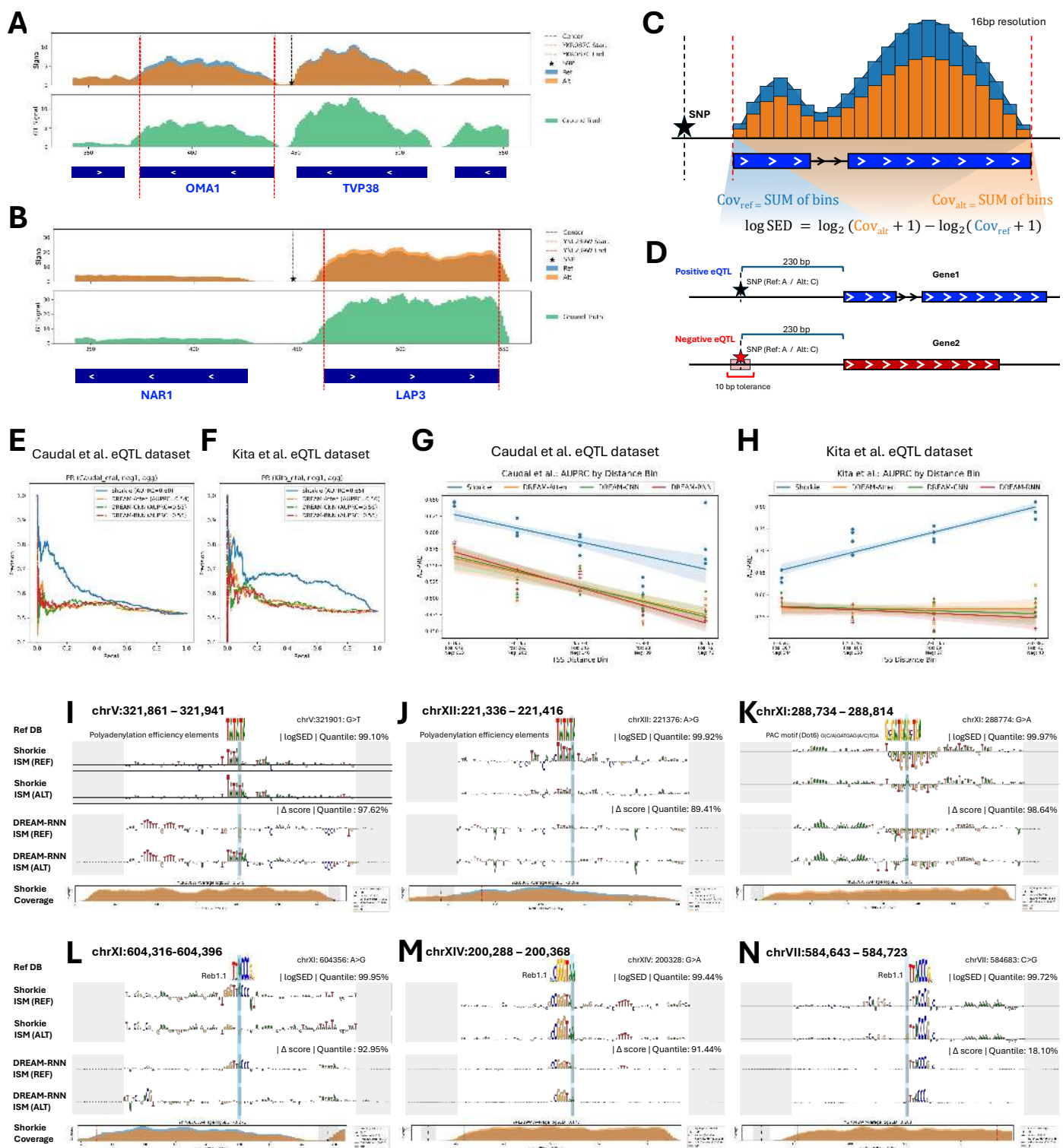
## 297 **Shorkie predicts *cis*-eQTL regulatory impacts**

298 Predicting how regulatory variants alter gene expression is critical for dissecting genetic association mech-  
299 anisms. We used Shorkie to interpret yeast *cis*-eQTL effects<sup>66,94,95</sup>. For instance, at the OMA1 locus, the  
300 eQTL alternate allele reduced Shorkie’s predicted RNA-seq coverage relative to reference (Figure 7A)<sup>66,95</sup>,  
301 while at LAP3, the lead eQTL alternate allele increased predicted coverage (Figure 7B). We quantified var-  
302 iant effects via gene expression log fold-change (Figure 7C; Methods).

303 We benchmarked Shorkie on eQTL datasets from Caudal et al. (1,901 local *cis*-eQTLs from ~1,000 yeast  
304 isolates; Figure S22A–B; Methods)<sup>66,95</sup> and Kita et al. (683 eQTLs categorized as “Promoter,” “UTR5,”  
305 “UTR3,” and “ORF”; Methods)<sup>94</sup>. Negative controls were generated by randomly selecting noncoding SNPs  
306 matched by reference/alternate alleles, TSS distance, and minor allele frequency  $\geq 5\%$  (Figure 7D; Figure  
307 S22C; Figure S23A; Methods).

308 Compared to top DREAM Challenge models (DREAM-CNN, DREAM-RNN, DREAM-Atten)<sup>93</sup>, Shorkie  
309 achieved superior ROC and PR metrics for both Caudal et al. (Figure 7E; Figure S22E) and Kita et al. da-  
310 tasetes (Figure 7F; Figure S23), with higher AUROC and AUPRC scores across all TSS distance bins (Figure  
311 7G–H).

312 ISM analyses of eQTL SNP loci revealed allele-specific remodeling of key *cis*-regulatory elements: creation  
313 (Figure 7I) or loss (Figure 7J) of polyadenylation efficiency motifs; disruption of PAC motif bound by  
314 Dot6/Tod6 repressors, correlating with increased expression (Figure 7K); and Reb1 motif alterations either  
315 weakening binding (Figure 7L) or strengthening it (Figure 7M–N). See Figure S24 for additional examples.



316

317 **Figure 7.** Shorkie accurately predicts cis-eQTL variant effects. (A) Positive eQTL example demonstrating  
 318 expression associated with the alternate allele at OMA1 (chrXI:603,195–604,232). (B) Positive eQTL showing  
 319 increased expression associated with the alternate allele at LAP3 (chrXIV:200,569–201,933). (C) Computation of log  
 320 fold-change expression scores for evaluating variant effects. (D) Generation of negative eQTL controls matched by  
 321 genomic characteristics from ~1,000 natural yeast isolates. (E,F) Precision-recall (PR) curves comparing Shorkie and  
 322 DREAM models for Caudal et al. (E) and Kita et al. (F) datasets. (G,H) AUPRC scores by TSS distance bins in Caudal

323 et al. (G) and Kita et al. (H) datasets. (I–N) ISM maps centered on eQTL SNPs, highlighting regulatory motifs  
324 identified by Shorkie (SNP in light-blue) and DREAM-RNN, with adaptor sequences in gray.

## 325 Discussion

326 In this study, we demonstrate that pretraining a language model on Saccharomycetales order genomes strikes  
327 an optimal balance between dataset size and diversity with regulatory conservation for learning the *S. cere-*  
328 *visiae* cis-regulatory grammar. Integrating convolutions, transformers, and residual connections within a U-  
329 Net architecture, Shorkie LM attains low perplexity on held-out *S. cerevisiae* chromosomes and reconstructs  
330 canonical TF binding motifs *de novo* across diverse fungal species. Transfer learning from this foundation  
331 into Shorkie substantially improves RNA-seq and ChIP-exo/MNase coverage prediction, raising median bin-  
332 level Pearson’s R from 0.67 (random initialization baseline) to 0.78 and gene-level Pearson’s R from 0.74  
333 to 0.88.

334 TF-MoDISco and ISM analyses show that Shorkie captures both static and dynamic motif signatures. The  
335 model recovers static motifs including Poly(dA:dT) tracts, Cbf1, Reb1 and canonical splice-site signals,  
336 confirming that evolutionary pretraining internalizes core regulatory elements. During transcription factor  
337 induction time courses, Shorkie dynamically tracks regulatory changes: STRE motif signals sharpen during  
338 MSN2 and MSN4 responses, and E-box engagement reflects cofactor-mediated MET4 activation.

339 While ISM and TF-MoDISco correctly identify E-box motifs at MET4 target promoters, they also highlight  
340 secondary signals with unclear mechanisms. These features may reflect persistent SCF<sup>Met30</sup>-mediated ubi-  
341 quitination of Met4<sup>96</sup>, transcriptional squelching via Mediator/SAGA cofactor sequestration<sup>97,98</sup> or feedback  
342 from intracellular sulfur metabolites<sup>99</sup>. Disentangling these regulatory layers demands targeted experi-  
343 ments—e.g. Met4 and cofactor ChIP-seq, ubiquitination-deficient Met4 variants and time-resolved metabo-  
344 lomics. Such studies could distinguish direct DNA binding from indirect regulatory cascades and guide fu-  
345 ture Shorkie applications to transcriptional networks.

346 Shorkie’s ability to learn regulatory patterns from perturbation experiments opens new avenues for network  
347 inference. Genome-wide perturbation screens could systematically map regulator-target relationships  
348 through sequence interpretation, while time-series experiments could reveal kinetic parameters governing  
349 transcription, mRNA degradation, and splicing. These applications extend Shorkie beyond prediction toward  
350 mechanistic discovery. Current limitations point to specific technical challenges. ChIP-exo was the most  
351 challenging assay: predictions track broad trends but systematically underestimate extreme, narrow peaks.  
352 This reflects the zero-inflated, heavy-tailed nature of ChIP-exo data and suggests that variance-stabilizing  
353 target transforms (e.g., square-root) or distribution-aware losses (e.g., focal loss) could improve performance  
354 (Figure S25). Similarly, Shorkie RNA-seq time course predictions show compressed dynamic range relative  
355 to experimental measurements (Figure 5B,G).

356 Yeast exemplifies where self-supervised pretraining excels. With a ~12 Mb compact genome, *S. cerevisiae*  
357 cannot support large-scale supervised learning from scratch. MPRA approaches offer synthetic training data  
358 but sacrifice native genome context, creating a fundamental tradeoff between scale and biological authen-  
359 ticity. Our evolutionary pretraining strategy complements these efforts, allowing Shorkie to learn

360 representations grounded in native promoters. Future work may explore clever combinations to exploit the  
361 strengths of the different approaches.

362 We propose that genome size and label abundance jointly determine when pretraining yields benefit. Small-  
363 genome organisms with limited experimental data (like yeast) benefit most from phylogenetically informed  
364 pretraining. Large-genome species (like mammals) with extensive experimental resources may show smaller  
365 pretraining gains, though benefits could persist in data-sparse contexts and cross-species transfer. Systematic  
366 studies varying taxonomic scope, phylogenetic distance, and experimental data volume will define the  
367 boundaries of this pretraining regime.

368 Our results establish masked nucleotide pretraining as a powerful foundation for regulatory genomics. By  
369 training at optimal evolutionary scales, models learn transferable representations that substantially improve  
370 sequence-to-expression prediction. While technical challenges remain, this framework provides a scalable  
371 path toward mechanistic understanding of gene regulation across diverse biological systems.

## 372 **Online Methods**

### 373 **Fungal genome and annotation download, filtering, and biotype partitioning**

374 We downloaded the main species table containing accession IDs, species names, and assembly and database  
375 identifiers from Ensembl Fungi release 59 ([https://ftp.ebi.ac.uk/ensemblgenomes/pub/release-59/fungi/species\\_EnsemblFungi.txt](https://ftp.ebi.ac.uk/ensemblgenomes/pub/release-59/fungi/species_EnsemblFungi.txt)) and used to retrieve all genomic and annotation data. For each taxon, we retrieved  
376 the unmasked genome FASTA file and its XML metadata, plus the corresponding GTF annotation file.  
377

378 We filtered each genome FASTA was filtered by parsing its XML to determine assembly level (“chromo-  
379 some” vs. “scaffold”). We retained both chromosome- and scaffold-level assemblies, discarding contigs  
380 shorter than 32,768 bp. We renamed remaining contigs according to a unified convention (chrI through  
381 chrXVI for yeast, original names otherwise). We produced a cleaned species manifest listing assembly lev-  
382 els, chromosome counts, and total base counts for all retained genomes. We partitioned each annotation file  
383 into subsets based on gene biotype (e.g. protein-coding, non-coding, rRNA, etc.).

### 384 **Shorkie LM model architecture**

385 Shorkie is a U-Net transformer-based model that processes 16,384-bp genomic windows and outputs a prob-  
386 ability distribution over the four nucleotides (A, C, G, T) at each position (Figure S1). It contains 13,665,828  
387 parameters (13,651,812 trainable; 14,016 non-trainable).

### 388 **Encoder (down-sampling) path**

389 The encoder begins with a 1D convolution (kernel size = 11, filters = 96), projecting the input tensor to a  
390 feature map of shape (16,384 × 96). Seven successive residual down-sampling stages follow. At stage  $i$ , the  
391 feature map is passed through a residual convolutional block: batch normalization → Gaussian Error Linear  
392 Unit (GELU) activation<sup>100</sup> → Conv1D (kernel size = 5; filters =  $C_i$ ) → dropout ( $p = 0.05$ ) → learned residual  
393 scaling. The block output is added to its input, and a MaxPooling1D layer (pool size = 2) reduces the se-  
394 quence length in half. Channel widths increase across stages as  $C_i \in [96, 128, 160, 192, 256, 320, 384]$ , while  
395 the sequence length decreases by  $2 \times$  each stage from 16,384 (1-bp) to 128 (128-bp).

## 396 **Transformer bottleneck**

397 At 128 bp resolution, the (128 × 384) feature map was feed into a stack of eight Transformer blocks to  
398 capture long-range dependencies. Each layer uses LayerNorm, then multi-head self-attention (model dim =  
399 384; 8 heads; key dim = 64) with residual dropout (p = 0.05), followed by a two-layer position-wise feed-  
400 forward network (Dense → Dropout → ReLU → Dense → Dropout) with a residual connection. This bot-  
401 tleneck integrates context across the full input span.

## 402 **Decoder (up-sampling) path and output**

403 The decoder mirrors the encoder in seven up-sampling stages. Each stage applies BatchNorm and GELU, a  
404 channel-preserving projection (Dense 384→384), and UpSampling1D (size = 2) to double the sequence  
405 length (e.g., 128 → 256 → ... → 16,384). The up-sampled features are merged with the corresponding  
406 encoder features via U-Net–style skip connections to restore fine-grained detail, then refined by a depthwise-  
407 separable convolution (kernel size = 3; filters = 384), BatchNorm, and GELU. A final 1×1 Conv1D (filters  
408 = 4) followed by softmax yields a per-position probability distribution over {A, C, G, T}. (Nearest-neighbor  
409 upsampling follows the UpSampling1D definition.)

## 410 **Phylogenetic tree reconstruction and fungal genomes distance estimation**

### 411 **Phylogenetic tree creation**

412 We converted all species names to NCBI Taxonomy identifiers (TaxIDs) with the ETE Toolkit Python API  
413 <sup>38</sup>. To reconstruct a minimal spanning tree that includes only our taxa of interest, we ran the ete-ncbiquery  
414 command-line module with our TaxID list and requested Newick-formatted output. This step pruned all non-  
415 target lineages while preserving branch lengths and hierarchical relationships. We loaded the resulting  
416 Newick file into ETE to verify correct monophyletic groupings and confirm presence of every target species  
417 in NCBI. For visualization, we uploaded the pruned tree to the Interactive Tree Of Life (iTOL v5) web  
418 server<sup>39,40</sup> and overlaid custom annotations (e.g., colored clade highlights) using iTOL dataset templates  
419 (Figure S2).

### 420 **Pairwise genomic distance estimation using alignment- and sketching-based methods**

421 We quantified genomic divergence of the 1,341 fungal assemblies relative to the *S. cerevisiae* R64 reference  
422 using complementary alignment-based and sketching-based methods. For the alignment-based approach, we  
423 ran MUMmer4’s nucmer<sup>41,101</sup> with 40 CPU threads to align each cleaned FASTA against the R64 reference,  
424 producing a delta file of all maximal unique matches. We then extracted detailed alignment statistics with  
425 `show-coords -lcr`, yielding tables of reference and query start-end positions, block lengths, percent  
426 identity, and coverage for every alignment block. To visualize large-scale synteny and structural variation,  
427 we generated dot plots for each genome pair using MUMmerplot. These plots enable rapid inspection of  
428 collinearity breaks, inversions, and translocations across the fungal assemblies.

429 For rapid, alignment-free genomic distance estimation, we applied two sketching-based tools: Dashing2<sup>102</sup>  
430 and Mash<sup>42</sup> on the same genome pairs. Dashing 2 leverages the SetSketch data structure<sup>103</sup> (using a truncated  
431 logarithm of hashed k-mers) together with ProbMinHash<sup>104</sup> for multiplicity-aware sketching. Mash<sup>42</sup> em-  
432 ploys classical MinHash<sup>105</sup> on each genome’s k-mer set (default k = 21, sketch size s = 1,000) to estimate  
433 Jaccard similarity<sup>106</sup>.

## 434 **Shorkie LM data preprocessing**

### 435 **Repetitive region detection and masking**

436 Most fungal assemblies in the Ensembl database lack adequate soft-masking of repetitive elements. To ad-  
437 dress this, we retrieved 1,501 genomes from the Ensembl FTP site and implemented a two-tiered repeat-  
438 masking pipeline. First, we generated a *de novo* repeat library for each genome using RepeatModeler v2.0<sup>44</sup>,  
439 which integrates multiple discovery algorithms, RepeatScout, RECON, and LTR\_retriever<sup>44,47,48</sup>, to capture  
440 both interspersed and structural elements unique to each assembly, and then merged the resulting consensus  
441 sequences with curated entries from Dfam<sup>46,107</sup>. Next, we ran RepeatMasker against this custom library (plus  
442 standard repeat databases), employing RMBlast, a RepeatMasker compatible version of the standard NCBI  
443 blastn program (<https://www.repeatmasker.org/rmbblast/>), for sensitive and high-throughput alignments. We  
444 enabled the “-xsmall” option to soft-mask repeats in lowercase thereby preserving original sequence  
445 length and coordinates and the “-gff” flag to output annotations in GFF3 format. Finally, we applied the  
446 DUST algorithm (<https://meme-suite.org/meme/doc/dust.html>) via the MEME suite<sup>43</sup> (default cutoff score  
447 threshold 20) to soft-mask residual low-complexity regions (Figure S3A). After the pipeline, 1,341 genomes  
448 were successfully masked and used for further preprocessing. We validated the workflow by comparing  
449 Ensembl’s original soft-masked regions with our custom masks in six representative assemblies and ob-  
450 served strong concordance (Figure S3B). We then partitioned each genome into 16,384-bp windows with a  
451 4,096-bp stride and applied a 7% repeat-content threshold for quality control, excluding ~20% of windows  
452 from the training, validation, and test sets (Figure S3C).

### 453 **Homologous and paralogous sequence removal between training, validation and test sets**

454 To guard against data leakage and ensure truly independent evaluation splits, we developed a three-step  
455 homology-filtering pipeline that removes any training sequences sharing appreciable similarity with those  
456 in our validation or test sets<sup>50</sup>. First, we aligned every training sequence against both the validation and test  
457 sets using Minimap2 (v2.28-r1209; assembly-to-assembly mode, “-x asm”)<sup>49</sup> to produce PAF files that  
458 report, for each query, the coordinates and match statistics of all detected alignments. From each PAF file,  
459 we extracted two metrics for every query-target pair:

- 460 • Coverage = matching bases (PAF field 10) / query length (PAF field 2)
- 461 • Identity = matching bases (PAF field 10) / alignment block length (PAF field 11)

462 These ratios quantify how much of a training sequence overlaps with other splits and how similar those  
463 overlaps are. Finally, we removed any training sequence for which both coverage  $\geq 5\%$  and identity  $\geq 30\%$   
464 in any alignment, yielding a leakage-free training corpus (Figure S3D). The alignment scatter plots for train-  
465 test and train-validation splits for each dataset are shown in Figure S3E.

### 466 **Evaluating fungal genome annotation completeness using BUSCO**

467 We first retrieved and unpacked the BUSCO v5 fungal lineage dataset (fungi\_odb10, 758 orthologs; 2024-  
468 01-08; [https://busco-data.ezlab.org/v5/data/lineages/fungi\\_odb10.2024-01-08.tar.gz](https://busco-data.ezlab.org/v5/data/lineages/fungi_odb10.2024-01-08.tar.gz)) to establish a con-  
469 sistent benchmark. For each of the 1,341 assemblies, we extracted all proteins from the Ensembl Fungi  
470 release 59 GFF annotation and its corresponding genome using gffread<sup>108</sup>. We then ran BUSCO<sup>109</sup> in protein

471 mode (`-m proteins`) on each proteome. BUSCO classifies each ortholog as “Complete (single-copy or  
472 duplicated)”, “Fragmented”, or “Missing”, providing a standardized completeness score across all annota-  
473 tions (Figure S4G–I).

#### 474 **TFRecord generation of one-hot encoded genomic windows with exon and repeat masks**

475 We loaded the repetitive, homologous, and paralogous-cleaned 16,384 bp windows with `pysam`  
476 (<http://code.google.com/p/pysam/>). For each window we: (1) one-hot encoded the DNA sequence; (2) com-  
477 puted exon masks by projecting GTF-derived transcript models onto the window and trimming a 2-bp flank-  
478 ing “chew” region; (3) derived repeat masks by flagging lowercase bases; and (4) assigned a species index.  
479 We flattened the resulting arrays (“sequence”, “exon\_mask”, “repeat\_mask”, “species”) to bytes and serial-  
480 ized them as TensorFlow Example protocol buffers. Finally, we wrote ZLIB-compressed TFRecord shards  
481 containing 32 examples each.

#### 482 **Shorkie LM training**

483 We trained and evaluated Shorkie LM using the ZLIB-compressed TFRecord shards described above via a  
484 custom version of the baskerville API, called `baskerville-yeast`. We split shards 80:20 into train:validation.  
485 Each epoch sampled up to 150 minibatches (batch size = 8) from an in-memory shuffle buffer of 256 records,  
486 with full reshuffling between epochs.

487 At each step, we masked 15% of positions per sequence ( $m = \lfloor 0.15 \times 16,384 \rfloor = 2,457$ ). Following the  
488 BERT protocol<sup>110</sup>: 80 % of masked sites were properly masked, 10 % were substituted by a random nucle-  
489 otide, and 10 % were left unchanged to prevent the model from over-relying on the mask token distribution.  
490 To exploit the fact that double-stranded DNA is symmetric under reverse complementation, we applied re-  
491 verse-complement augmentation to each input sequence with probability 0.5, thereby encouraging the model  
492 to learn strand-invariant features, a strategy shown to substantially improve performance<sup>111,112</sup>.

493 We computed categorical cross-entropy over masked positions only and reweighted the loss by genomic  
494 region. In order to focus the model on regulatory sequences, we down-weighted exonic and repeat regions  
495 by a factor of 0.1, following the success of this strategy for repeats in prior work<sup>52</sup>. We trained the model  
496 using the Adam optimizer<sup>113</sup> (learning rate =  $1 \times 10^{-4}$ ;  $\beta_1 = 0.7$ ;  $\beta_2 = 0.9$ ), global clipnorm = 0.1, and a linear  
497 warmup over the first 20,000 steps.

498 Training proceeded for up to 10,000 epochs (minimum = 100; maximum = 10,000), each capped at 150 steps,  
499 with early stopping according to the validation loss (patience = 1,000 epochs). At the end of each epoch, we  
500 averaged validation loss over five independent mask-and-predict passes per example (`repeat_eval = 5`) to  
501 reduce sampling noise. We retained the checkpoint with the lowest validation loss as the final model.

#### 502 **Shorkie LM evaluation**

503 We evaluated held-out test windows using the same “mask → predict → tile” procedure. For each test se-  
504 quence, we iteratively masked 15 % of positions (2,457 positions) and predicted them until all positions had  
505 been covered. We specified the `--rc` flag to average predictions from both the forward sequence and its  
506 reverse complement. We then computed the per-position cross-entropy loss (Equation 1) for each 16k win-  
507 dow  $s$  as:

$$\mathcal{L}_{CE_s} = \frac{\sum_{j=1}^L \left[ -\sum_{n \in \{A, C, G, T\}} \mathbb{I}_{j,n}^{(s)} \ln p_{j,n}^{(s)} \right] w_j^{(s)}}{\sum_{j=1}^L w_j^{(s)}}$$

Equation 1

508 where  $L = 16,384$  is the window length,  $\mathbb{I}_{j,n}^{(s)} \in \{0,1\}$  is the one-hot indicator for base  $n$  at position  $j$  in win-  
 509 dow  $s$ ,  $p_{j,n}^{(s)}$  is the model's predicted probability for base  $n$ , and  $w_j^{(s)}$  is a position-specific weight (exon/re-  
 510 peat scaling) at that position. The global test loss is the average over all  $N$  test sequences:

$$\mathcal{L}_{CE} = \frac{1}{N} \sum_{s=1}^N \mathcal{L}_{CE_s}$$

Equation 2

513 Perplexity is then defined as

$$\text{Perplexity} = \exp(\mathcal{L}_{CE})$$

Equation 3

516 Finally, we calculated taxon-specific evaluation metrics by summing segment losses for each species and  
 517 dividing by the number of segments for that species, yielding a species-level loss. In addition to loss and  
 518 perplexity, we computed and stored a position probability matrix (PPM) for each sequence (`x_pred`),  
 519 alongside one-hot inputs (`x_true`), species labels, and scaling weights for each base pair in a `pred.npz` file.  
 520 Finally, we converted these PPMs into information content matrices (ICMs)<sup>114,115</sup> to facilitate the identifica-  
 521 tion of TF binding sites.

## 522 Computation of Shorkie LM's information content matrix

523 Let  $p_{j,n}^{(s)}$  be the PPM probability for nucleotide  $n \in \{A, C, G, T\}$  at position  $j$  in the 16,384 bp window  $s$ . We  
 524 transform this PPM into an ICM via the following steps. To avoid zeros, we first add a small pseudocount  
 525  $\varepsilon$ , yielding

$$\tilde{p}_{j,n}^{(s)} = p_{j,n}^{(s)} + \varepsilon$$

Equation 4

528 and then renormalize across nucleotides so that

$$\bar{p}_{j,n}^{(s)} = \frac{\tilde{p}_{j,n}^{(s)}}{\sum_{m \in \{A, C, G, T\}} \tilde{p}_{j,m}^{(s)}}, \quad \sum_n \bar{p}_{j,n}^{(s)} = 1$$

Equation 5

531 The Shannon entropy at position  $j$  is

$$H_j^{(s)} = - \sum_{n \in \{A, C, G, T\}} \bar{p}_{j,n}^{(s)} \log_2(\bar{p}_{j,n}^{(s)})$$

Equation 6

534 with larger  $H_j^{(s)}$  indicating greater nucleotide diversity (i.e. lower conservation). We then define the per-  
535 position information content as

$$536 \quad C_j^{(s)} = \log_2(4) - H_j^{(s)} = 2 - H_j^{(s)}, \quad C_j^{(s)} \in [0, 2]$$

537 Equation 7

538 For logo visualization, each nucleotide's column height is set to

$$539 \quad h_{j,n}^{(s)} = \bar{p}_{j,n}^{(s)} \times C_j^{(s)}, \quad \sum_n h_{j,n}^{(s)} = C_j^{(s)}$$

540 Equation 8

541 Stacking letters in ascending  $h_{j,n}^{(s)}$  order produces a DNA logo whose total column height reflects the infor-  
542 mation content,  $C_j^{(s)}$ , and whose individual letter heights encode the normalized nucleotide probabilities,  
543  $\bar{p}_{j,n}^{(s)}$ .

## 544 **Constructing a DNA logo of the SMT3 promoter region using SpeciesLM**

545 We adapted the Python notebook workflow of SpeciesLM<sup>52</sup> ([https://github.com/gagneurlab/dependencies\\_DNALM/blob/main/compute\\_and\\_visualize\\_dep\\_maps.ipynb](https://github.com/gagneurlab/dependencies_DNALM/blob/main/compute_and_visualize_dep_maps.ipynb)). First, we loaded a pretrained BERT-  
546 style masked-language model (BertForMaskedLM; “johahi/specieslm-fungi-upstream-k1”) and its Au-  
547 toTokenizer via Hugging Face Transformers<sup>116</sup>. Following the notebook, we applied this model to the 1  
548 kb upstream region of SMT3 in *S. cerevisiae*. For each sequence, we tokenized the nucleotides and pre-  
549 pended a proxy-species identifier. At inference, a softmax produced per-position base-probability distribu-  
550 tions. We then computed per-position information content (IC) from the reference (unmutated) probabilities  
551 against the yeast genomic background. Finally, we generated a sequence logo by setting the height of nucle-  
552 otide  $b$  at position  $j$  to its predicted probability multiplied by the IC (Equation 8), so total stack height re-  
553 flects conservation and relative letter heights encode base preferences (Figure 2B).

## 555 **Six genomic datasets and PPM construction for Shorkie LM evaluation**

556 We evaluated Shorkie LM across six diverse genomic datasets:

- 557 1. *S. cerevisiae* reference genome (R64).
- 558 2. Four randomly selected *S. cerevisiae* strains: *S. cerevisiae* YJM1202, YJM1400, YJM555, and YJM984.
- 559 3. Five genomes within the Saccharomycetales order: *Candida albicans*, *Eremothecium gossypii* FDAG1,  
560 *Kluyveromyces lactis* str. NRRL Y-1140, *Komagataella phaffii* CBS 7435, and *Candida glabrata*.
- 561 4. Four genomes from the broader Ascomycota phylum: *Aspergillus fumigatus*, *Neurospora crassa*, *Peni-*  
562 *cillium chrysogenum* str. P2niaD18, and *Tuber melanosporum*.
- 563 5. Four genomes from the Orbiliales order: *Arthrotrrys flagrans* str. CBS H-5679, *Arthrotrrys oligo-*  
564 *spora* ATCC 24927, *Dactylellina haptotyyla* CBS 200.50, and *Drechlerella stenobrocha* 248.

565 6. Four genomes from the Schizosaccharomycetales order: *Schizosaccharomyces cryophilus*, *Schizosac-*  
566 *charomyces japonicus*, *Schizosaccharomyces octosporus*, and *Schizosaccharomyces pombe*.

567 Except for *S. cerevisiae* and other Saccharomycetales genomes, all others were held out during training. We  
568 then applied the approach described in the “Shorkie LM Data Preprocessing” section to segment genomes  
569 into 16,384 bp windows, create ZLIB-compressed TFRecord shards, and predict with Shorkie LM.

## 570 **Motif discovery with TF-MoDISco-Lite**

### 571 **TF-Modisco run on the *S. cerevisiae* genome**

572 To identify salient sequence motifs from Shorkie LM’s predictions, we employed TF-MoDISco-Lite  
573 (<https://github.com/jmschrei/tfmodisco-lite>), a memory- and time-efficient reimplementa-  
574 tion of TF-MoDISco<sup>53,54</sup>. First, we loaded one-hot encoded inputs (`x_true`) and predicted probabilities (`x_pred`)  
575 from the `pred.npz` file generated during Shorkie LM evaluation. For each nucleotide channel  $i \in \{A, C, G, T\}$   
576 at every position  $j = 1, \dots, L$  within a 16,384 bp window, we added a pseudocount  $\varepsilon = 1 \times 10^{-4}$  to each  
577 probability  $p_{i,j}$ . We then computed the local background frequency at position  $j$  as

$$578 \quad \bar{p}_j = \frac{1}{4} \sum_{i \in \{A, C, G, T\}} (p_{i,j} + \varepsilon)$$

579 Equation 9

580 And transformed each adjusted probability into a log-odds score:

$$581 \quad \Delta_{i,j} = (p_{i,j} + \varepsilon) \log \left( \frac{p_{i,j} + \varepsilon}{\bar{p}_j} \right)$$

582 Equation 10

583 thereby accentuating deviations from the local background. We saved these log-odds matrices as “`x_true.npz`”  
584 and “`x_pred.npz`”, reshaped them to (samples  $\times$  positions  $\times$  channels), and then ran `modisco motifs -`  
585 `s x_true.npz -a x_pred.npz -n 1000000 -w 16384` to sample one million seqlets across  
586 all 16,384 bp windows. TF-MoDISco-Lite then clustered those seqlets into consolidated motifs by first com-  
587 puting cosine similarity on gapped k-mer representations and next refining clusters via fine-grained realign-  
588 ment. For each resulting cluster, it computes both a contribution weight matrix (CWM) and a PWM. The  
589 `modisco report` step then generated motif logos and an interactive HTML summary.

### 590 **Curating known TF-binding motifs from multiple yeast motif databases**

591 We curated known TF-binding motifs from six publicly available databases, covering literature-curated as-  
592 sociations, computational predictions, and high-throughput assay-derived specificity profiles. We included  
593 YEASTRACT, which provides 732 curated yeast motifs (average width 9.8 bp)<sup>59</sup> (<https://yeastract.com/>);  
594 SwissRegulon Yeast, with 158 genome-wide motifs<sup>58</sup> (<https://swissregulon.unibas.ch/pages/>); UniPROBE  
595 Yeast (GR09), containing 89 PBM-derived motifs<sup>57</sup> (<http://thebrain.bwh.harvard.edu/uniprobe/>); MacIsaac  
596 v1, offering 124 phylogenetically conserved motifs<sup>56</sup> ([https://fraenkel-nsf.csbi.mit.edu/improved\\_map/](https://fraenkel-nsf.csbi.mit.edu/improved_map/));  
597 SCPD, supplying 24 promoter-derived motifs ([https://esefinder.ahc.umn.edu/cgi-](https://esefinder.ahc.umn.edu/cgi-bin/tools/ESE3/esefinder.cgi)  
598 [bin/tools/ESE3/esefinder.cgi](https://esefinder.ahc.umn.edu/cgi-bin/tools/ESE3/esefinder.cgi)) we batch-normalize

599 )<sup>60</sup>; and YeTFaSCo, a comprehensive compendium of 1,709 motifs for 256 yeast proteins calculated and  
600 quality-assessed by multiple metrics (<http://yetfasco.cabr.utoronto.ca/>)<sup>55</sup>. All motifs were downloaded in  
601 MEME format from the respective websites and merged to obtain a comprehensive yeast TF-binding motif  
602 database for integration into our downstream analyses.

### 603 **Genome-wide enrichment analysis of TF-MoDISco motifs around TSS**

604 To map TF-MoDISco-derived motifs to the genome coordinates and assess their proximity to transcription  
605 start sites (TSS), we first parsed the HDF5 output from the TF-MoDISco-Lite run (`modisco_re-`  
606 `sults.h5`) to extract seqlet positions (start, end, example index, strand) for each pattern. We then con-  
607 verted these coordinates to genomic coordinates by adding the seqlet offsets to the corresponding entries in  
608 the original BED file of one-hot input windows, yielding a unified BED of motif hits annotated with species,  
609 chromosome, start, end, strand, and the best-matching known motif with its associated q-value.

610 Next, we generated a comprehensive TSS annotation by parsing GTF files. For each gene, we recorded the  
611 the 5' end of its transcripts (start coordinate for “+” strand, end for “-” strand) and sorted these TSS positions  
612 per chromosome. For each motif hit, we computed its midpoint and located the nearest TSS via binary search  
613 on the sorted list of transcripts, defining signed distance as negative for upstream (“-”) and positive for  
614 downstream (“+”) relative to the TSS strand.

615 To establish a background distribution, we sampled one random position per seqlet from the metaclusters on  
616 the same chromosome, drawing positions uniformly across chromosome lengths, and calculated each posi-  
617 tion’s closest-TSS distance using the same procedure. Finally, we compared the observed and background  
618 distance distributions by plotting histograms over a  $\pm 2.5$  kb window centered on the TSS to highlight motif  
619 enrichment patterns.

### 620 **Embedding t-SNE clustering using Shorkie LM**

621 To generate low-dimensional embeddings of genomic contexts, we first defined three interval classes across  
622 all sixteen *S. cerevisiae* chromosomes: (1) Promoters: 500 bp immediately upstream of each start codon; (2)  
623 gene bodies: the span between each gene’s start and end coordinates; (3) intergenic regions: all regions not  
624 annotated as gene, exon, or CDS.

625 Each interval was retrieved from the reference genome using pysam, centered, and end-padded to 16,384  
626 bp. Sequences on the “-” strand were reverse-complemented. We then one-hot encoded the four nucleotide  
627 channels and concatenated a 165-dimensional species one-hot vector, setting the *S. cerevisiae* channel (index  
628 114) to 1. Batches of eight intervals (shape 16,384 $\times$ 170) were fed into a pretrained Shorkie LM. To capture  
629 intermediate representations, we defined for each selected layer a sub-model that takes the original LM in-  
630 puts and returns that layer’s activations. Ten selected layers are `max_pooling1d_6`, `multihead_attention`,  
631 `dense`, `dense_1`, `multihead_attention_7`, `dense_14`, `dense_15`, `dense_16`, `dense_28`, `dense_29`.

632 For each interval, we mean-pooled the per-position outputs across the sequence axis, concatenated the re-  
633 sulting vectors to form an embedding of dimension  $D$ , and stored the matrix ( $N$  intervals  $\times$   $D$ ) in HDF5  
634 datasets. We also saved accompanying metadata arrays: chromosome, coordinates, strand, feature class, and  
635 `gene_id`. Next, we aggregated metadata across all chromosomes and parsed each gene interval’s biotype  
636 from the Ensembl GTF, categorizing intervals into five groups: “Protein-coding gene,” “Intergenic region,”

637 “tRNA,” “Transposable element,” and “Promoter”. For each layer’s embedding matrix, we applied t-SNE  
638 <sup>117</sup> (scikit-learn `TSNE` with `n_components = 2`) to project the D-dimensional embeddings into two dimen-  
639 sions. The resulting 2D coordinates were visualized to assess clustering patterns by genomic feature.

## 640 **Shorkie and Shorkie\_Random\_Init data pre-processing**

### 641 **RNA-Seq perturbation experiments**

642 The design of the inducible genetic perturbation experiments builds on methods developed for the Induction  
643 Dynamics gene Expression Atlas<sup>34</sup>, where hundreds of transcription factors were independently induced and  
644 resultant gene expression changes were profiled over time using RNA hybridization microarrays. New data  
645 used for Shorkie-supervised training were generated at Calico Life Sciences LLC using updated miniaturized  
646 chemostats, or ministats. The instrumentation and RNA-sequencing protocols are described here.

### 647 **Strain construction and selection**

648 Each time-course experiment uses a strain selected from the Yeast Estradiol strains with Titratable Induction  
649 (YETI) collection<sup>118</sup>, where the native promoter for a gene of interest was replaced with a synthetic Z3EV  
650 inducible promoter which drives transcription in the presence of estradiol. New data collected for this study  
651 can be split into three partitions: (1) many replicates of pre-induction cultures of the MSN4 inducible strain,  
652 (2) a set of 8 TF perturbations in replicate, matched to previously measured microarray data, and (3) a set of  
653 460 other genes, including kinases, phosphatases, and other transcriptional regulators prioritized as likely to  
654 have downstream transcriptional changes as assessed using evidence from YeastMine curated annotations<sup>119</sup>,  
655 Phenome<sup>120</sup>, and Fitness Clusters<sup>121</sup>.

### 656 **Growth conditions**

657 For all experiments, cells were grown under continuous culturing conditions. Cultures were maintained un-  
658 der phosphate limitation in minimal chemically defined media prepared by mixing 20mL of 1000x vitamin  
659 solution, 20mL of 1000x metals, 40mL of 10g/L KH<sub>2</sub>PO<sub>4</sub>, 1L 40% dextrose, and 2L 10X salts solution in  
660 16L of milliQ water to bring the total volume to 20L. (Stock solutions defined in Table S1).

### 661 **Instrumentation and perturbation experiments**

662 To increase throughput and coverage of gene expression changes, we developed a ministat array system  
663 capable of growing 24 30mL cultures in parallel to steady state, applying a chemical perturbation, and gener-  
664 ating samples for time-resolved omic measurements. The system is housed in a 30°C warm room and  
665 consists of four banks of six 100 mL round-bottom vessels, integrated via a 24-vessel manifold into a com-  
666 plete ministat array (Figure S11).

667 Each ministat array incorporates three commercial peristaltic pumps: (1) a media pump that delivers fresh  
668 media to each vessel at a constant rate (the dilution rate), (2) a sampling pump that extracts a bolus of cells  
669 from each culture, and (3) an input pump used initially to inoculate cultures and subsequently to deliver a  
670 chemical perturbation (e.g., estradiol). Each pump uses separate tubing for each of the 24 ministats. Each  
671 vessel is also equipped with a tube for the effluent, which is continuously weighed to monitor and control  
672 the dilution rate, and a tube for delivering air. Air is supplied via a Flow Master 2400 airflow regulator  
673 connected to a humidifier to minimize culture evaporation. The regulator also splits the airflow into 24 in-  
674 dividually adjustable lines, one for each culture.

675 For each culture, a pre-induction ( $T_0$ ) sample was collected prior to perturbation. The input pump then de-  
676 livered 100  $\mu$ L of 500  $\mu$ M b-estradiol in base media from a 96-well plate into each vessel to initiate the time  
677 course. Following induction, cultures were typically sampled at 8 timepoints using the sampling pump to  
678 collect each sample into 96 deep-well plates pre-filled with chilled lysis buffer containing RNase inhibitor  
679 (Takara). Samples from the 24 ministats were staggered across quadrants of each 96-well plate, enabling  
680 systematic re-arraying of the eight timepoint-specific plates into two consolidated 96-well plates for down-  
681 stream processing. Sample plates were flash frozen in liquid nitrogen immediately after collection.

## 682 **RNA-seq library preparation**

683 Samples from four timepoints were combined into a single 96-well PCR plate using the Bravo BenchCel  
684 system and processed using a miniaturized, high-throughput adaptation of standard library preparation pro-  
685 tocols. After addition of oligo-dT primer (sequence: AAGCAGTGGTATCAACGCAGAG-  
686 TACTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTV), samples were lysed by three freeze-thaw cycles and  
687 incubated at 42°C for 3 minutes. cDNA synthesis was performed using the SMARTscribe Reverse Tran-  
688 scriptase kit (Takara), with 2.4  $\mu$ M LNA template switch oligo (sequence: AAGCAGTGGTATCAAC-  
689 GCAGAGTACrGrG+G) and RNase inhibitor. For cDNA amplification, SeqAmp DNA polymerase (Takara)  
690 and ISO PCR primer (sequence: AAGCAGTGGTATCAACGCAGAGT) were added. The samples were  
691 placed in a thermocycler with the following conditions: 95°C for 1 min, 10 cycles of [98°C for 10 s, 65°C  
692 for 30 s, 68°C for 3 min], 72°C for 10 min. Cleanup was performed using RNAClean XP beads (Beckman  
693 Coulter) at 0.9x reaction volume.

694 Relative cDNA concentration was measured using Quant-iT PicoGreen (Fisher). For absolute quantification  
695 and quality control, 11 representative samples (chosen to span the range of PicoGreen values) were assessed  
696 using a High Sensitivity DNA kit (Agilent) on the BioAnalyzer. These values were used to generate a stand-  
697 ard curve, and all samples were diluted to 200 pg/ $\mu$ L using a Mantis dispenser.

698 Library preparation was performed using the Nextera XT kit on cDNA from two ministat batches at a time,  
699 re-arrayed into a 384-well PCR plate. Samples were pooled, and pooled libraries underwent a double  
700 RNAClean XP bead cleanup. Final library concentration and quality were confirmed using a High Sensitivity  
701 DNA kit. Libraries were sequenced by Genewiz on the Illumina NovaSeq.

## 702 ***S. cerevisiae* R64 reference genome pre-processing**

703 Contigs from *S. cerevisiae* R64 reference were first split at assembly gaps and at hypervariable regions  
704 identified in the Rossi et. al. study (these sites include the rDNA locus, tRNA genes and telomere regions  
705 and are available in 02\_References\_and\_Features\_Files at [https://github.com/CEGRcode/2021-Rossi\\_Na-  
706 ture](https://github.com/CEGRcode/2021-Rossi_Nature)) hereafter referred to as rossi\_mask.bed.

707 We then trimmed 1,024 bp from each contig end and discarded those shorter than 16,384 bp. Contigs longer  
708 than 786,432 bp were split in half. The genome was then segmented into overlapping 16,384 bp windows  
709 with a 6,165 bp stride. These windows were shuffled and partitioned into eight cross-validation folds by  
710 balancing total nucleotide counts. Unmappable positions (from hypervariable regions described in  
711 rossi\_mask.bed) were annotated, and any window with > 50% unmappability (--umap\_clip 0.5) was  
712 removed.

## 713 **ChIP-exo and ChIP-MNase samples preprocessing**

714 Paired-end sequencing data was obtained directly from Rossi et al<sup>65</sup> (GSE147927) and sequence alignment  
715 was performed using Bwa-0.7.17 mem algorithm<sup>122</sup> and multi-mappers removed using SAMtools<sup>123</sup>. For  
716 ChIP-exo data the position of the 5' end of Read 1 was used, while the full span for MNase was used to  
717 generate tracks. BAM files were additionally qc'd using PICARD (<https://github.com/broadinstitute/picard>)  
718 to mark and remove duplicates. Reads overlapping regions described in rossi\_mask.bed were removed and  
719 experiments with less than 10,000 remaining reads and greater than 75 percent duplication rate were dropped  
720 from the dataset. Lastly, BEDtools<sup>124</sup> and bedGraphToBigWig (<https://www.encodeproject.org/software/bedgraph-tobigwig/>)  
721 were used to transform the BAM files to BigWig tracks, yielding 1,128 ChIP-exo,  
722 and 20 ChIP-MNase tracks.

## 723 **RNA-seq samples preprocessing**

724 First adapter sequences were trimmed from the FASTQ files using bmap (<https://github.com/BioInfoTools/BBMap>). Then transcript alignment and quantification was performed using STAR<sup>125</sup>. A genomic  
725 index was created using S288c\_R63-3 and quantification was performed using GeneCounts with the follow-  
726 ing parameters (`--outFilterMultimapNmax 1, --bamRemoveDuplicatesType UniqueIdentical, --alignIntronMin 10, --alignIntronMax 2500, --alignMatesGapMax 2500`). BAM files were processed to remove PCR duplicates marked using PICARD<sup>126</sup>. Lastly, bam\_cov.py  
727 ([https://github.com/calico/basenji/blob/master/bin/bam\\_cov.py](https://github.com/calico/basenji/blob/master/bin/bam_cov.py)) was used to produce BigWig tracks. Data  
728 from 1000 strains (Caudal et al.) used unpaired RNA-seq<sup>66</sup>, and was filtered to retain samples with greater  
729 than 150,000 reads and less than 80 percent duplication rate while the in-house generated induction experi-  
730 ments used paired-end RNA-seq and was filtered to retain samples with greater than 150,000 reads and mean  
731 insert size greater than 250 bp. This produced 3,053 induction RNA-Seq tracks and 1,014 1,000-strain RNA-  
732 Seq tracks.  
733  
734  
735

## 736 **Validation**

737 To ensure that gene expression dynamics measured using the new ministat array system are reliable and  
738 accurate, we compared data collected on the high-throughput ministat array to previously collected micro-  
739 array data from Hackett et al.<sup>34</sup>. To assess the correspondence between transcriptional profiles in each system,  
740 we calculated the Pearson correlation of gene expression fold-changes across all genes for a given matched  
741 timepoint post induction (Figure S12A), using both raw and log<sub>2</sub> shrunken fold-changes. Next, to evaluate  
742 the new system's sensitivity in detecting previously identified differentially expressed genes, we conducted  
743 an ROC analysis. To select a positive set of differentially expressed genes, we selected genes with an abso-  
744 lute log<sub>2</sub> fold-change greater than one in the microarray data as a heuristic (Figure S12B). We then tested for  
745 differential expression in the new system by separately fitting a standard ordinary least squares linear regres-  
746 sion for each gene's expression time course following perturbation. In this model, gene expression relative  
747 to the mean  $T_0$  expres-sion is described as a function of the time point, coded as a categorical variable for  
748 each timepoint, and a co-variate for the culture vessel corresponding to experiment replicate. Then we fit an  
749 ANOVA for each gene expression time course regression and used the F-statistic, describing the variance  
750 between timepoints to the variance among replicates within a timepoint. By varying the cutoff F-statistic,  
751 we calculated the AUROC (Figure S12C-D).

## 752 **Track data transformation**

753 We processed 5,215 BigWig tracks: 3,053 induction timepoint RNA-Seq datasets<sup>34</sup>, 1,014 RNA-Seq datasets  
754 from various yeast strains<sup>66</sup>, 1,128 ChIP-exo tracks, and 20 ChIP-MNase tracks<sup>65</sup>. For each 16,384-bp win-  
755 dow, we extracted per-base coverage and imputed missing values with the window's median. To reduce  
756 edge effects, we cropped 1,024 bp from each end, retaining a 14,336-bp interior. We then summed coverage  
757 across consecutive 16-bp intervals to form 896 non-overlapping bins, yielding a 896-element vector per  
758 window. We saved these vectors as float16 in HDF5. Finally, we serialized the one-hot DNA sequence (via  
759 pysam), the corresponding binned coverage vector, and the unmappability mask into ZLIB-compressed  
760 TFRecord files (256 windows per file) organized by cross-validation fold.

## 761 **Shorkie model architecture and hyperparameters**

762 Shorkie fine-tunes the pretrained Shorkie LM backbone (13.7 M parameters) to predict RNA-Seq, ChIP-exo,  
763 and ChIP-MNase signals across 16,384 bp windows. Its trunk replicates Shorkie LM exactly:

- 764 • Initial Conv1D projection: 11 bp kernel  $\times$  96 filters (linear activation)
- 765 • Seven residual down-sampling blocks: each block is BatchNorm $\rightarrow$ GELU $\rightarrow$ Conv1D(5 bp) with filter  
766 counts increasing from 96 to 384 in 32-filter steps, followed by 5% dropout, a skip connection, and  
767 MaxPool1D (pool size = 2)
- 768 • Transformer bottleneck: eight layers operating over 128 positions. Each layer contains LayerNorm $\rightarrow$ 4-  
769 head self-attention (model dimension = 384, key dimension = 64) with 20% dropout, followed by a feed-  
770 forward (LayerNorm $\rightarrow$ Dense $\rightarrow$ ReLU $\rightarrow$ Dropout $\rightarrow$ Dense $\rightarrow$ Dropout) and residual connections.

771 The decoder uses the same U-Net up-sampling scheme as Shorkie LM but with only three up-sampling  
772 stages to restore the 16 bp resolution. At each stage, feature maps are batch-normalized, passed through a  
773 GELU nonlinearity, projected via a Dense layer (384  $\rightarrow$  384), and doubled in length with UpSampling1D,  
774 then merged with the corresponding encoder output via a U-Net skip connection. At the final 16 bp resolution,  
775 a Cropping1D layer (cropping = 64) removes convolutional padding artifacts and is followed by GELU. A  
776 single Dense layer then projects each position into 5,215 channels: TF-perturbed RNA-Seq ( $n = 3,053$ ),  
777 1,000-strain RNA-Seq ( $n = 1,014$ ), ChIP-exo ( $n = 1,128$ ), and ChIP-MNase histone marks ( $n = 20$ ), and a  
778 Softplus activation ensures all outputs remain positive.

779 We fine-tuned Shorkie with Adam ( $\beta_1 = 0.7$ ,  $\beta_2 = 0.9$ ; global clip-norm = 0.1), learning rate =  $2 \times 10^{-5}$   
780 with 20,000 warm-up steps, batch = 8, using a Poisson + Multinomial loss (5x scaling of multinomial loss  
781 component) and early stopping with a patience of 150 epochs.

## 782 **Shorkie\_Random\_Init model architecture and hyperparameters**

783 The Shorkie\_Random\_Init model uses an identical architecture as Shorkie but with all weights initialized  
784 from scratch. We trained Shorkie\_Random\_Init in a supervised manner using Adam ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ;  
785 global clip-norm = 0.1), learning rate =  $1 \times 10^{-4}$  with 5,000 warm-up steps, and otherwise identical hyperpa-  
786 rameters. Shorkie\_Random\_Init thus serves as a baseline for quantifying the performance improvements  
787 attributable to transfer learning.

## 788 **Shorkie and Shorkie\_Random\_Init bin-level, gene-level and track-level evaluation**

### 789 **Bin-level evaluation**

790 We generated per-base predictions for each 16,384 bp input window by first averaging the model's outputs  
 791 across both the forward and reverse-complement strands (strand-ensemble) and then applying a shift-ensem-  
 792 ble over offsets  $S = \{0, 1\}$  bp to smooth boundary artifacts. Specifically, the ensemble prediction at genomic  
 793 position  $i$  and track  $j$  is defined as:

$$794 \hat{y}_{i,j}^{(\text{ens})} = \frac{1}{|S|} \sum_{s \in \{0,1\}} \frac{y_{i+s,j}^{\text{fwd}} + y_{i+s,j}^{\text{rc}}}{2}$$

795 Equation 11

796 where  $y_{i,j}^{\text{fwd}}$  and  $y_{i,j}^{\text{rc}}$  are the strand-specific model outputs. We compared  $\hat{y}_{i,j}^{(\text{ens})}$  to measured coverage  $y_{i,j}$   
 797 and reported bin-level Pearson correlation and coefficient of determination ( $R^2$ ) across all tracks and cross-  
 798 validation folds.

### 799 **Raw (non-normalized) gene-level evaluation**

800 Let  $G$  denote all annotated genes from the *S. cerevisiae* GTF. We assigned a given 16,384 bp bin  $i$  to gene  
 801  $g$  if at least 50% of the bin's width ( $p$ ) overlaps and exon of gene  $g$ :  $\mathcal{B}(g) = \{b: \text{overlap}(b, g) \geq 0.5 p\}$ .

802 For each gene  $g$  and track  $j$ , we aggregated the strand- and shift-ensemble predictions  $\hat{y}_{i,j}^{(\text{ens})}$  and the true  
 803 coverages  $y_{i,j}$  over  $\mathcal{B}(g)$  and stabilized variance via a  $\log_2$ -transform with a pseudocount of 1:

$$804 \hat{Y}_{g,j} = \log_2 \left( \sum_{i \in \mathcal{B}(g)} \hat{y}_{i,j}^{(\text{ens})} + 1 \right), \quad Y_{g,j} = \log_2 \left( \sum_{i \in \mathcal{B}(g)} y_{i,j} + 1 \right)$$

805 Equation 12

806 For each track  $j$ , we computed gene-level Pearson  $r$  and  $R^2$  between  $\{Y_{g,j}\}_{g \in G}$  and  $\{\hat{Y}_{g,j}\}_{g \in G}$ .

### 807 **Assay-normalized gene-level evaluation**

808 To account for the differing dynamic ranges across assays, we first partition tracks into assay-specific groups  
 809  $T_k$ , including TF-perturbed RNA-seq, 1,000-strain RNA-seq, ChIP-exo, and ChIP-MNase. Within each  
 810 group  $T_k$ , we independently applied quantile normalization (QN) to the  $\log_2$ -transformed aggregated cover-  
 811 ages to equalize their distributions:

$$812 \hat{Z}_{g,j} = \text{QN}_k(\hat{Y}_{g,j}), \quad Z_{g,j} = \text{QN}_k(Y_{g,j}), \quad \forall j \in T_k$$

813 Equation 13

814 We then mean-centered each gene within its assay group:

$$815 \Delta \hat{Z}_{g,j} = \hat{Z}_{g,j} - \frac{1}{|T_k|} \sum_{j' \in T_k} \hat{Z}_{g,j'}, \quad \Delta Z_{g,j} = Z_{g,j} - \frac{1}{|T_k|} \sum_{j' \in T_k} Z_{g,j'}$$

816 Equation 14

817 We summarized the gene-level concordance with Pearson  $r$  and  $R^2$  between  $\{\Delta Z_{g,j}\}_{g \in G}$  and  $\{\Delta \hat{Z}_{g,j}\}_{g \in G}$  for  
 818 each  $j$ . Normalizing by assay-specific variability yields metrics that isolate predictive skill from absolute  
 819 signal scale and enable evaluation of track specificity—that is, whether the model reproduces the differences  
 820 among tracks within each assay.

### 821 **Within-gene bin-level consistency**

822 To evaluate how precisely Shorkie reconstructs fine-scale positional coverage within individual genes, we  
 823 analyzed consistency at the bin-level within each gene. First, for each gene  $g$  and track  $j$ , we derived  $\log_2$ -  
 824 transformed coverage vectors from the predictions and measured coverages:

$$825 \quad \hat{v}_{g,j}^{\text{raw}} = \left\{ \log_2 \left( \hat{y}_{i,j}^{(\text{ens})} + 1 \right) \right\}_{i \in B(g)}, \quad v_{g,j}^{\text{raw}} = \left\{ \log_2 (y_{i,j} + 1) \right\}_{i \in B(g)}$$

826 Equation 15

827 We computed Pearson  $r$  for gene-track pairs with non-trivial variation (both variances  $> 10^{-6}$ ). To de-em-  
 828 phasize nearly flat profiles, we further required  $\text{Var}(\hat{v}_{g,j}^{\text{raw}}) > \tau_j$ , where  $\tau_j$  is the 80th percentile of predicted  
 829 variances across genes for track  $j$ . For each gene we summarized:

$$830 \quad r_g^{\text{within, raw}} = \frac{1}{|\{j : \text{Var}(\hat{v}_{g,j}^{\text{raw}}) > \tau_j\}|} \sum_{j: \text{Var}(\hat{v}_{g,j}^{\text{raw}}) > \tau_j} r_{g,j}^{\text{bin, raw}}$$

831 An identical procedure was conducted using quantile-normalized and mean-centered vectors to compute  
 832  $r_g^{\text{within, norm}}$ .

### 833 **Track-level evaluation**

834 To evaluate Shorkie’s ability to accurately capture each gene’s pattern across RNA-Seq assays specifically,  
 835 we performed a track-level analysis focusing exclusively on RNA-Seq tracks. For each gene  $g$ , we con-  
 836 structed vectors of coverage across all RNA-Seq tracks, both in their raw and normalized (mean-centered  
 837 and quantile-normalized) forms:

$$838 \quad Y_g = \{Y_{g,j}\}_{j \in \mathbb{J}_{\text{RNA-Seq}}}, \quad \hat{Y}_g = \{\hat{Y}_{g,j}\}_{j \in \mathbb{J}_{\text{RNA-Seq}}}, \quad \Delta Z_g = \{\Delta Z_{g,j}\}_{j \in \mathbb{J}_{\text{RNA-Seq}}}, \quad \Delta \hat{Z}_g = \{\Delta \hat{Z}_{g,j}\}_{j \in \mathbb{J}_{\text{RNA-Seq}}}$$

839 Equation 16

840 We reported Pearson  $r$  and  $R^2$  between predicted and observed vectors for both raw ( $Y_g$  and  $\hat{Y}_g$ ) and normal-  
 841 ized ( $\Delta Z_g$  and  $\Delta \hat{Z}_g$ ) settings, summarizing per gene.

### 842 **Attention weight matrix visualization from selected Shorkie embeddings**

843 To inspect how Shorkie’s transformer blocks focus on different parts of a gene, we extracted 16,384 bp  
 844 windows centered on two example genes (EFM5 at chr VII: 489,391–505,775; RPL7A at chr VII: 356,973–  
 845 373,357) from the *S. cerevisiae* R64 reference using pysam. We center-trimmed or padded each window  
 846 with “N”s to exactly 16,384 bp and parsed the Ensembl GTF with PyRanges<sup>127</sup> to build an annotation table  
 847 of gene bodies, exons, and 5’/3’ UTRs for downstream overlay.

848 We evaluated three models: Shorkie LM (pretrained LM), Shorkie (fine-tuned), and Shorkie\_Random\_Init  
849 (baseline). For each  $16,384 \times 4$  input tensor, we averaged forward and reverse-complement predictions and  
850 concatenated results across replicates to yield a  $1 \times N_{\text{reps}} \times 16,384$  coverage array. To capture self-  
851 attention, we computed dot-products between learned query and key vectors at every position pair  $(i, j)$ ,  
852 applied softmax to obtain attention weights, and collected these across all eight cross-validation folds for (i)  
853 the first transformer block and (ii) the final two blocks. This yielded a single  $128 \times 128$  attention map per  
854 block set.

855 For visualization, we clipped each fold-averaged attention map below  $10^{-4}$  and above 0.05 to enhance con-  
856 trast, then displayed it as a heatmap across the full window. We plotted predicted coverage as filled curves  
857 along the top and right margins. We converted genomic features into “attention-bin” coordinates  
858 (bin index =  $\lfloor (\text{position} - \text{window\_start})/128 \rfloor$ ) and overlaid as colored boxes for gene bodies or lines  
859 for UTRs and exons, with strand-specific coloring.

## 860 ***In silico* mutagenesis analysis of Shorkie and Shorkie\_Random\_Init**

### 861 **Sequence extraction and formatting**

862 We applied *in silico* mutagenesis (ISM) to quantify the effect of every single-nucleotide variant (SNV)  
863 within yeast promoter regions based on Shorkie and Shorkie\_Random\_Init predictions. We defined input  
864 regions as 16,384-bp windows centered on promoter segments (450 bp upstream and 50 bp downstream of  
865 the TSS). The promoter set consisted of 137 ribosomal protein genes, 64 ribosome/rRNA biosynthesis (RRB)  
866 genes, and 3,258 additional protein-coding genes from the *S. cerevisiae* R64 reference genome. We extracted  
867 the corresponding sequences from the reference FASTA while preserving strand orientation.

868 For each 16 kb window, we computed ISM maps across the 500 bp promoter segment. For each reference  
869 sequence  $s$  and position  $p = 1, \dots, 500$ , we generated three mutant sequences by substituting the reference  
870 base with each of the three alternative nucleotides  $n \in \{A, C, G, T\} \setminus \{\text{ref}\}$ , yielding  $3 \times 500$  mutant se-  
871 quences per window.

### 872 **ISM importance score matrix construction**

873 We one-hot encoded all reference and mutant sequences and ran them through Shorkie and Shorkie\_Ran-  
874 dom\_Init with strand-ensemble averaging. For each sequence  $s$ , genomic bin  $i \in \{1, \dots, 896\}$ , and track  $j$ ,  
875 let  $\hat{y}_{\text{ref}, i, j}^{(s)}$  be the model prediction for the reference sequence, and  $\hat{y}_{\text{alt}, i, j}^{(s, p, n)}$  be the prediction for the variant at  
876 position  $p$  with nucleotide  $n$ . We computed the  $\log_2$  fold-change score to quantify the variant effect:

$$877 \quad \log_2 \text{FC}_j^{(s, p, n)} = \log_2 \left( \sum_{i \in \mathcal{B}} \hat{y}_{\text{alt}, i, j}^{(s, p, n)} + 1 \right) - \log_2 \left( \sum_{i \in \mathcal{B}} \hat{y}_{\text{ref}, i, j}^{(s)} + 1 \right)$$

878 Equation 17

879 where  $\mathcal{B}$  indexes the 896 output bins. We saved the resulting  $\log_2$  fold-change scores, together with metadata  
880 (chromosome, start, end, strand, reference, and alternate alleles), to an HDF5 file (scores.h5) via h5py  
881 (<https://www.h5py.org/>).

## 882 ISM importance score matrix normalization and ISM map visualization

883 To derive per-position importance profiles, we defined the set  $J$  of  $T_0$  RNA-Seq tracks. For each sequence  
884  $s$ , mutation  $(p, n)$ , and position  $i$ , we averaged the  $\log_2$  fold-change scores across  $T_0$  tracks:

$$885 \quad M_{i,n}^{(s)} = \frac{1}{|J|} \sum_{j \in J} \log_2 \text{FC}_{i,n,j}^{(s)}$$

886 Equation 18

887 where  $i$  indexes positions and  $n$  denotes nucleotides. Next, we zero-mean normalized at each position  $i$   
888 across the four nucleotides:

$$889 \quad \tilde{M}_{i,n}^{(s)} = M_{i,n}^{(s)} - \frac{1}{4} \sum_{m \in \{A, C, G, T\}} M_{i,m}^{(s)}$$

890 Equation 19

891 Finally, to focus on reference-base contributions, we computed:

$$892 \quad \text{Logo}_i^{(s)} = \sum_{n \in \{A, C, G, T\}} \tilde{M}_{i,n}^{(s)} \times H_{i,n}^{(s)}$$

893 Equation 20

894 where  $H_{i,n}^{(s)} \in \{0, 1\}$  is the one-hot indicator for the reference base at position  $i$ . Finally, we visualized per-  
895 position scores  $\text{Logo}_i^{(s)}$  as DNA sequence logos (ISM maps), highlighting the magnitude and direction of  
896 variant effects across the promoter window.

## 897 Motif discovery using TF-MoDISco-Lite

898 To identify TF-binding motifs, we first constructed two arrays: `ref.npz` of shape  $(N_{\text{seq}}, 4, M)$ , containing  
899 one-hot encoded reference sequences, and `pred.npz` of shape  $(N_{\text{seq}}, 4, M, |J|)$ , containing per-replicate var-  
900 iant scores for  $T_0$  RNA-Seq tracks.

901 Here,  $N_{\text{seq}}$  is the number of promoter windows,  $M = 500$  is the promoter window length, and  $|J|$  is the num-  
902 ber of  $T_0$  RNA-Seq replicates. We reshaped `pred.npz` to  $(N_{\text{seq}} \times |J|, 4, M)$  by concatenating across repli-  
903 cates, then ran TF-MoDISco-Lite (see “TF-MoDISco run on the *S. cerevisiae* genome” section) on these  
904 matrices to uncover recurring sequence motifs.

## 905 Time-series motif analysis of TF-induction with Shorkie

906 To capture dynamic motif changes during TF induction (e.g., MSN2, MSN4, MET4), we filtered the RNA-  
907 Seq metadata to select only the tracks corresponding to each TF time course and grouped biological repli-  
908 cates by sampling time. We parsed time-point annotations (e.g., “T\_0”, “T\_15”, “T\_30”) directly from the  
909 track identifiers, ordered them chronologically, and for each time point  $T_t$  defined  $J_t$  as the set of associated  
910 RNA-Seq track indices.

911 For each promoter window  $s$ , base position  $i \in \{1, 2, \dots, 500\}$ , nucleotide channel  $n \in \{A, C, G, T\}$ , and time  
 912 point  $T_t$ , we computed an ISM importance score matrix by averaging  $\log_2$  fold-change scores across all  
 913 replicate tracks:

$$914 \quad M_{i,n}^{(s,T_t)} = \frac{1}{|\mathbb{J}_t|} \sum_{j \in \mathbb{J}_t} \log_2 \text{FC}_{i,n,j}^{(s,T_t)}$$

915 Equation 21

916 To focus on differential saliency relative to baseline expression, we baseline-corrected each time point by  
 917 subtracting the  $T_0$  map:

$$918 \quad \Delta M_{i,n}^{(s,T_t)} = M_{i,n}^{(s,T_t)} - M_{i,n}^{(s,T_0)}$$

919 Equation 22

920 We then applied zero-mean normalization across nucleotides at each position to  $\Delta M$  and visualized the re-  
 921 sulting  $\Delta$ ISM maps for each  $T_t$ . Finally, we reshaped per-timepoint  $\Delta$ ISM matrices to  $(N_{\text{seq}} \times |\mathbb{J}_t|, 4, M)$  and  
 922 applied TF-MoDISco-Lite to identify motifs whose importance trajectories changed over the induction time  
 923 course.

#### 924 **Gene-level coverage calculation for experimental measurements and Shorkie predictions**

925 To compare experimental and predicted gene-level coverages at Shorkie's native 16 bp resolution, we  
 926 summed bin values over each gene's span. For gene  $g$  and track  $j$ , let  $\mathcal{B}(g)$  be the set of overlapping 16 bp  
 927 bins. We define:

$$928 \quad \text{cov}_j^{\text{exp}}(g) = \sum_{i \in \mathcal{B}(g)} y_{i,j}^{\text{exp}}, \quad \text{cov}_j^{\text{pred}}(g) = \sum_{i \in \mathcal{B}(g)} y_{i,j}^{\text{pred}}$$

929 Equation 23

930 where  $\mathcal{B}(g)$  indexes the bins overlapping gene  $g$ . Next, we normalized to reads-per-million (RPM) using  
 931 each sample's library size  $\text{libsize}_j$ ,

$$932 \quad \text{cov}_j^{\text{exp,RPM}}(g) = \frac{\text{cov}_j^{\text{exp}}(g)}{\text{libsize}_j} \times 10^6, \quad \text{cov}_j^{\text{pred,RPM}}(g) = \frac{\text{cov}_j^{\text{pred}}(g)}{\text{libsize}_j} \times 10^6$$

933 Equation 24

935 To characterize time-dependent gene expression trajectories in TF-induction experiments, we averaged  
 936 RPM-normalized coverages across replicate tracks at each time point  $T_t$ :

$$937 \quad \overline{\text{cov}}^{(T_t, \text{exp,RPM})}(g) = \frac{1}{|\mathbb{J}_t|} \sum_{j \in \mathbb{J}_t} \text{cov}_j^{\text{exp,RPM}}(g), \quad \overline{\text{cov}}^{(T_t, \text{pred,RPM})}(g) = \frac{1}{|\mathbb{J}_t|} \sum_{j \in \mathbb{J}_t} \text{cov}_j^{\text{pred,RPM}}(g)$$

938 Equation 25

939 This yields matched, per-gene coverage trajectories for both experimental measurements and Shorkie pre-  
940 dictions.

#### 941 **Euclidean distance calculation of ISM maps**

942 To quantify motif importance changes during induction, we extracted the normalized importance score ma-  
943 trix  $\mathbf{P}_{s,t} \in \mathbb{R}^{L \times 4}$  for each promoter window  $s$  and timepoint  $t$ . We flattened this matrix into a vector  $\mathbf{v}_{s,t} \in$   
944  $\mathbb{R}^{4L}$ , and computed pairwise Euclidean distances between timepoints:

$$945 \quad D_s(t_1, t_2) = \|\mathbf{v}_{s,t_1} - \mathbf{v}_{s,t_2}\|_2, \quad \forall t_1, t_2 \in \{1, \dots, T\}$$

946 Equation 26

947 To summarize across all promoter windows, we conducted element-wise averaging:

$$948 \quad \bar{D}(t_1, t_2) = \frac{1}{S} \sum_{s=1}^S D_s(t_1, t_2)$$

949 Equation 27

950 where  $S$  is the total number of promoter windows. We visualized the resulting mean Euclidean-distance  
951 matrix  $\bar{D}$  as a heatmap, effectively capturing motif shifts over the TF induction time course.

#### 952 ***cis*-eQTL analysis with Shorkie and DREAM challenge models**

953 To evaluate model performance on *cis*-eQTLs, we benchmarked Shorkie against the DREAM challenge  
954 models<sup>93</sup> using two independent eQTL datasets.

#### 955 **Caudal et al. eQTLs from the pan-transcriptome of ~1,000 yeast natural isolates**

956 We imported the GWAS summary statistics<sup>66</sup> (file: GWAS\_combined\_lgcCorr\_ldPruned\_noBonfer-  
957 roni\_20221207.tab; downloaded from The 1002 Yeast Genome website: [http://1002genomes.u-  
958 strasbg.fr/files/RNAseq](http://1002genomes.u-<br/>958 strasbg.fr/files/RNAseq)) into a DataFrame and removed variants labeled as masked (ld\_mask = “masked”)  
959 were. We normalized Phenotype (Pheno\_pos) and SNP (ChrPos) positions by chromosome length, then  
960 classified a variant as *cis* if  $|\text{ChrPos} - \text{Pheno\_pos}| \leq 8,000$  bp on the same chromosome and as *trans*  
961 otherwise. We further stratified variants by subtype (SNP vs. CNV) (Figure S20A).

962 We parsed the gVCF containing 1,011 yeast isolates (1011Matrix.gvcf; downloaded from [http://1002ge-  
963 nomes.u-strasbg.fr/files/](http://1002ge-<br/>963 nomes.u-strasbg.fr/files/))<sup>95</sup> with pysam to extract reference and alternate alleles, chromosome, position, and  
964 quality. We then merged the *cis* and *trans* eQTL tables with the gVCF DataFrame on chromosome and  
965 position to identify intersecting and unique variant sets.

966 We removed variants with missing or non-positive P-values and computed significance as  $-\log_{10}(\text{PValue})$ .  
967 Finally, we generated a Manhattan plot by plotting  $-\log_{10}(\text{PValue})$  against cumulative genomic position  
968 with alternating colors per chromosome, overlaid a genome-wide significance line at  $P = 5 \times 10^{-8}$ , and  
969 placed chromosome ticks at median cumulative positions (Figure S20B).

## 970 **Kita et al. eQTLs from 85 diverse *S. cerevisiae* isolates**

971 We imported the summary statistics (pnas.1717421114.sd01; downloaded from  
972 [https://www.pnas.org/doi/suppl/10.1073/pnas.1717421114/suppl\\_file/pnas.1717421114.sd01.txt](https://www.pnas.org/doi/suppl/10.1073/pnas.1717421114/suppl_file/pnas.1717421114.sd01.txt)), yielding  
973 1,640 eQTLs. From these, we selected 683 variants in four genomic contexts: Promoter, UTR5, UTR3, and  
974 ORF. For each variant, we computed the absolute distance between its genomic coordinate (ChrPos) and the  
975 target gene's TSS, labeling it *cis* if  $|\text{ChrPos} - \text{TSS}| \leq 8,000$  bp and *trans* otherwise. Finally, we retrieved  
976 the corresponding allele sequences from 1011Matrix.gvcf using pysam.

## 977 **Negative-eQTL sampling**

978 We generated four independent negative-eQTL sets by sampling common, non-coding variants that matched  
979 the positive eQTLs in allele composition and distance to the gene TSS. We extracted *S. cerevisiae* TSS  
980 coordinates from the Ensembl GTF. From the gVCF, we retained variants located outside CDS/exon inter-  
981 vals with allele frequency (AF)  $\geq 0.05$  and excluded any variant that matched a positive eQTL.

982 For each positive eQTL, we shuffled its matching candidate list of negative variants on the same chromo-  
983 some and attempted to select one negative variant whose distance to a randomly chosen gene's TSS matched  
984 the positive's distance within  $\pm 100$  bp (with a fallback window of  $\pm 200$  bp if no match emerged). We en-  
985 forced that each negative was used only once per iteration, yielding one negative per positive. We repeated  
986 this sampling four times to produce four distinct negative-eQTL sets.

## 987 **Variant effect prediction with Shorkie**

988 We predicted how individual variants alter gene-level coverage profiles with Shorkie. For each variant, we  
989 (1) extracted a 16,384 bp window centered on the SNP; (2) verified that the reference allele matched the  
990 extracted sequence; (3) generated one-hot encodings for both the reference and alternate alleles; and (4)  
991 averaged predictions over forward and reverse-complement strands. Within the window, we summed the  
992 predicted coverage values over all bins overlapping the annotated gene  $g$ 's exons:

$$993 \quad \text{Cov}_{\text{ref}} = \sum_{i \in \mathcal{B}(g)} y_{\text{ref}}(i), \quad \text{Cov}_{\text{alt}} = \sum_{i \in \mathcal{B}(g)} y_{\text{alt}}(i)$$

994 Equation 28

995 where  $\mathcal{B}(g)$  indexes the bins overlapping gene  $g$ . We then defined the  $\log_2$  fold-change score as

$$996 \quad \log_2 \text{FC}_{\text{Shorkie}} = \log_2(\text{Cov}_{\text{alt}} + 1) - \log_2(\text{Cov}_{\text{ref}} + 1)$$

997 Equation 29

## 998 **eQTL ISM analysis with Shorkie**

999 To probe local sequence drivers, we applied the ISM pipeline to 80-bp windows centered on each SNP (to  
1000 match DREAM input length), generated an  $80 \times 3$   $\Delta$  matrix per variant, formed the reference-average ma-  
1001 trix, and rendered sequence logos.

## 1002 **Predicting eQTL SNPs with DREAM challenge models**

1003 To benchmark Shorkie against established DREAM challenge models, we evaluated three pretrained models:  
1004 convolutional (DREAM-CNN), recurrent (DREAM-RNN), and attention-based (DREAM-Atten), on both  
1005 positive eQTLs and the four independently sampled negative sets.

## 1006 **Sequence extraction and formatting**

1007 For each variant (positive or negative), we extracted an 80bp window centered on the SNP from the *S.*  
1008 *cerevisiae* genome (40 bp upstream and 39 bp downstream), then prepended a fixed 17bp upstream flank  
1009 (TGCATTTTTTTCACATC) and appended a 13 bp downstream flank (GGTTACGGCTGTT), following  
1010 the DREAM models' expected input format, yielding final input sequences of 110 bp.

## 1011 **Model inference and log fold-change score calculation**

1012 We one-hot encoded the 110-bp reference and alternate sequences and scored them on a GPU to obtain scalar  
1013 predictions, then converted them to  $\log_2$  FC using the shared definition.

## 1014 **ISM analysis of DREAM-RNN**

1015 We performed single-nucleotide ISM on the 110-bp inputs by iterating over each position  $p = 1, \dots, 110$  and  
1016 substituting the native nucleotide with each of the three alternatives ( $n \in \{A, C, G, T\} \setminus \{\text{ref}\}$ ); we then com-  
1017 puted  $\Delta(p, n)$ , constructed the reference-average matrix, and visualized sequence logos.

## 1018 **Massively parallel reporter assay (MPRA) evaluation of promoter variants**

1019 We evaluated Shorkie's predictive performance for regulatory variants measured by a publicly available  
1020 MPRA experiment<sup>93</sup>. The MPRA included five single-sequence sub-libraries: high\_exp, low\_exp, yeast\_exp,  
1021 random\_exp, challenging\_exp, each containing unique 110 bp promoter sequences, and three dual-sequence  
1022 sub-libraries: SNVs\_exp, motif\_perturbation\_exp, motif\_tiling\_exp, consisting of paired reference and al-  
1023 ternate sequences. We excluded constructs present in any public leaderboard and randomly drew a uniform  
1024 sample of 1,000 sequences without replacement from each category.

## 1025 **Gene selection and library sampling**

1026 We stratified  $T_0$  RNA-Seq expression into three quantiles (5–25%, 25–75%, 75–95%) separately by strand.  
1027 From each quantile, we randomly chose three to four genes, resulting in 22 genes:

- 1028 • Forward strand genes (10 total):
  - 1029 ○ 5–25%: GPM3/SLI1/VPS52
  - 1030 ○ 25–75%: YMR160W/MRPS28/YCT1
  - 1031 ○ 75–95% RDL2/PHS1/RTC3/MSN4
- 1032 • Reverse strand genes (12 total):
  - 1033 ○ 5–25%: COA4/ERI1/RSM25/AIM11
  - 1034 ○ 25–75% ERD1/MRM2/SNT2/MRPL1
  - 1035 ○ 75–95% CSI2/RPE1/PKC1/MAE1

## 1036 Promoter insertion design

1037 To examine how Shorkie’s predictions vary relative to TSS, we first defined a promoter insertion window  
1038 of 110 bp ( $\pm 55$  bp) and enforced a minimum 100-bp offset from the TSS to avoid direct overlap. We then  
1039 chose eleven distinct insertion offsets (ranging from 100 to 200 bp, incremented by 10 bp). For each selected  
1040 gene, we calculated the insertion midpoints as follows:

- 1041 • Forward strand: midpoint = TSS – offset
- 1042 • Reverse strand: midpoint = TSS + offset

1043 Each midpoint defined a 110 bp replacement window from (midpoint – 55 bp) to (midpoint + 55 bp). We  
1044 then replaced these native genomic windows with each MPRA-derived sequence and used Shorkie to predict  
1045 the regulatory impact. In contrast, DREAM challenge models directly utilized the 110 bp MPRA windows  
1046 as input to predict scalar expression scores.

## 1047 Promoter variants effect quantification

1048 To quantify gene-level effects from Shorkie’s predictions, we identified model-predicted coverage bins  
1049  $B(g)$  overlapping the exonic regions of each target gene ( $g$ ).

1050 For single-sequence libraries, we computed  $\text{Cov}_{\text{native}} = \sum_{i \in B(g)} Y_{\text{native}}(i)$  and  $\text{Cov}_{\text{MPRA}} = \sum_{i \in B(g)} Y_{\text{MPRA}}(i)$ .  
1051 We reported the regulatory effect as  $\log_2 \text{FC}_{\text{single}} = \log_2(\text{Cov}_{\text{MPRA}} + 1) - \log_2(\text{Cov}_{\text{native}} + 1)$ .  $\log_2 \text{FC}_{\text{single}}$   
1052 reflects changes in predicted coverage due to the MPRA insert relative to the native sequence.

1053 In dual-sequence libraries, each construct includes both reference and alternate promoter variants. Shorkie  
1054 predicted REF and ALT coverage separately, and we computed  $\log_2 \text{FC}_{\text{ref}}$  and  $\log_2 \text{FC}_{\text{alt}}$  relative to the na-  
1055 tive promoter, and summarized the predicted differential effect as  $\Delta \log_2 \text{FC}_{\text{dual}} = \log_2 \text{FC}_{\text{alt}} - \log_2 \text{FC}_{\text{ref}}$ .

1056 Experimental reporter assays yielded corresponding expression scores for the reference ( $S_{\text{ref}}$ ) and alternate  
1057 ( $S_{\text{alt}}$ ) sequences. We quantified their difference as  $\Delta S = S_{\text{alt}} - S_{\text{ref}}$ . Finally, we visualized scatterplots to as-  
1058 sess the concordance between the Shorkie-predicted  $\Delta \log_2 \text{FC}_{\text{dual}}$  and the experimentally derived  $\Delta S$  scores.

## 1059 Data and Code Availability

- 1060 • The new Induction Dynamics Gene Expression Atlas RNA-seq datasets generated by Calico Life Sci-  
1061 ences LLC are hosted on Google Cloud Storage (GCS). Coverage tracks (BigWig) are available at  
1062 <gs://shorkie-paper/data/supervised/bigwigs/>, and processed TFRecords are at [gs://shorkie-paper/data/su-  
1063 pervised/processed](gs://shorkie-paper/data/supervised/processed/).
- 1064 • The genomes used for self-supervised language-model pretraining are at [gs://shorkie-paper/data/unsu-  
1065 pervised/genome/](gs://shorkie-paper/data/unsupervised/genome/), with corresponding TFRecords at <gs://shorkie-paper/data/unsupervised/processed/>.
- 1066 • The Shorkie LM and Shorkie models are implemented in TensorFlow. Shorkie LM is available at:  
1067 [gs://seqnn-share/shorkie\\_lm/](gs://seqnn-share/shorkie_lm/). Shorkie models are available at: <gs://seqnn-share/shorkie/>.
- 1068 • Parameters, training code and evaluation scripts for both the Shorkie LM and Shorkie models are avail-  
1069 able under the Apache-2.0 license at <https://github.com/calico/baskerville-yeast> and  
1070 <https://github.com/calico/shorkie-paper> under license Apache-2.0.

## 1071 **Author Contribution**

1072 J.L. and D.R.K. conceived the project.

1073 K-H.C., M.M.M., J.L., and D.R.K. designed the research.

1074 K-H.C., J.L., and D.R.K. developed the baskerville-yeast repository.

1075 S.H. led design of the yeast induction experiments and instrumentation.

1076 E.S. optimized library preparation miniaturization, processed samples and generated sequencing libraries for  
1077 the induction experiments.

1078 M.M.M. processed the RNA-Seq, ChIP-exo, and ChIP-MNase data for model training.

1079 K-H.C. and J.L. trained the Shorkie LM models.

1080 K-H.C. and M.M.M. trained the Shorkie\_Random\_Init models.

1081 K-H.C. trained the Shorkie models.

1082 K-H.C. and J.L. conducted analyses on Shorkie LM transcription factor motif inference and Shorkie model  
1083 interpretability.

1084 K-H.C. and J.L. conducted the MPRA and eQTL analyses.

1085 K-H.C., M.M.M., E.S., S.H., J.L., and D.R.K. wrote the manuscript.

## 1086 **Funding**

1087 This research was supported in part by the U.S. National Institutes of Health (NIH) under grants R01-  
1088 HG006677 and R35-GM156470, and by the U.S. National Science Foundation (NSF) under grant DBI-  
1089 2412449. Computational analyses were performed using resources provided by the Advanced Research  
1090 Computing at Hopkins (ARCH) core facility, supported in part by NSF grant OAC-1920103. Additional  
1091 funding was provided by Calico Life Sciences LLC.

## 1092 **Acknowledgements**

1093 We gratefully acknowledge all members of the Kelley Lab at Calico Life Sciences LLC for their insightful  
1094 discussions. We thank David Botstein for guidance and feedback planning the induction time course exper-  
1095 iments. We also are grateful for contributions from Griffin Kim for designing the ministat arrays, Rebecca  
1096 Wang for adapting the library prep and automation of the higher throughput miniaturized vessels, and  
1097 Thomas Li for assistance generating the experimental data. We also thank Steven Salzberg and Mihaela  
1098 Pertea for their valuable ideas and thoughtful feedback. The Shorkie logo was generated with the help of  
1099 OpenAI in the style of Borzoi ([https://github.com/calico/borzoi/blob/main/borzoi\\_logo.png](https://github.com/calico/borzoi/blob/main/borzoi_logo.png)).

## 1100 **Reference**

1101 1. Harbison, C. T. *et al.* Transcriptional regulatory code of a eukaryotic genome. *Nature* **431**, 99–104  
1102 (2004).

- 1103 2. He, Q., Johnston, J. & Zeitlinger, J. ChIP-nexus enables improved detection of in vivo transcription  
1104 factor binding footprints. *Nat Biotechnol* **33**, 395–401 (2015).
- 1105 3. Lee, T. I. *et al.* Transcriptional Regulatory Networks in *Saccharomyces cerevisiae*. *Science* **298**, 799–  
1106 804 (2002).
- 1107 4. Rhee, H. S. & Pugh, B. F. Comprehensive Genome-wide Protein-DNA Interactions Detected at Single-  
1108 Nucleotide Resolution. *Cell* **147**, 1408–1419 (2011).
- 1109 5. Skene, P. J. & Henikoff, S. An efficient targeted nuclease strategy for high-resolution mapping of  
1110 DNA binding sites. *eLife* **6**, e21856 (2017).
- 1111 6. Struhl, K. Molecular mechanisms of transcriptional regulation in yeast. *Annual review of biochemistry*  
1112 **58**, 1051–1077 (1989).
- 1113 7. Venters, B. J. *et al.* A Comprehensive Genomic Binding Map of Gene and Chromatin Regulatory Pro-  
1114 teins in *Saccharomyces*. *Molecular Cell* **41**, 480–492 (2011).
- 1115 8. Venters, B. J. & Pugh, B. F. A canonical promoter organization of the transcription machinery and its  
1116 regulators in the *Saccharomyces* genome. *Genome Res.* **19**, 360–371 (2009).
- 1117 9. Weiner, A. *et al.* High-Resolution Chromatin Dynamics during a Yeast Stress Response. *Molecular*  
1118 *Cell* **58**, 371–386 (2015).
- 1119 10. Zentner, G. E., Kasinathan, S., Xin, B., Rohs, R. & Henikoff, S. ChEC-seq kinetics discriminates tran-  
1120 scription factor binding sites by DNA sequence and shape in vivo. *Nat Commun* **6**, 8733 (2015).
- 1121 11. Zhu, J. *et al.* Integrating large-scale functional genomic data to dissect the complexity of yeast regula-  
1122 tory networks. *Nat Genet* **40**, 854–861 (2008).
- 1123 12. Beer, M. A. & Tavazoie, S. Predicting Gene Expression from Sequence. *Cell* **117**, 185–198 (2004).
- 1124 13. Sharon, E. *et al.* Inferring gene regulatory logic from high-throughput measurements of thousands of  
1125 systematically designed promoters. *Nat Biotechnol* **30**, 521–530 (2012).
- 1126 14. De Boer, C. G. *et al.* Deciphering eukaryotic gene-regulatory logic with 100 million random promot-  
1127 ers. *Nat Biotechnol* **38**, 56–65 (2020).

- 1128 15. Vaishnav, E. D. *et al.* The evolution, evolvability and engineering of gene regulatory DNA. *Nature*  
1129 **603**, 455–463 (2022).
- 1130 16. Zhou, Z. *et al.* DNABERT-2: Efficient Foundation Model and Benchmark For Multi-Species Genome.  
1131 Preprint at <https://doi.org/10.48550/ARXIV.2306.15006> (2023).
- 1132 17. Ji, Y., Zhou, Z., Liu, H. & Davuluri, R. V. DNABERT: pre-trained Bidirectional Encoder Representa-  
1133 tions from Transformers model for DNA-language in genome. *Bioinformatics* **37**, 2112–2120 (2021).
- 1134 18. Nguyen, E. *et al.* Sequence modeling and design from molecular to genome scale with Evo. *Science*  
1135 **386**, eado9336 (2024).
- 1136 19. Zhai, J. *et al.* Cross-species modeling of plant genomes at single nucleotide resolution using a pre-  
1137 trained DNA language model. Preprint at <https://doi.org/10.1101/2024.06.04.596709> (2024).
- 1138 20. Benegas, G., Batra, S. S. & Song, Y. S. DNA language models are powerful predictors of genome-  
1139 wide variant effects. *Proc. Natl. Acad. Sci. U.S.A.* **120**, e2311219120 (2023).
- 1140 21. Sanabria, M., Hirsch, J., Joubert, P. M. & Poetsch, A. R. DNA language model GROVER learns se-  
1141 quence context in the human genome. *Nat Mach Intell* **6**, 911–923 (2024).
- 1142 22. Karollus, A. *et al.* Species-aware DNA language models capture regulatory elements and their evolu-  
1143 tion. *Genome Biol* **25**, 83 (2024).
- 1144 23. Dalla-Torre, H. *et al.* Nucleotide Transformer: building and evaluating robust foundation models for  
1145 human genomics. *Nat Methods* **22**, 287–297 (2025).
- 1146 24. Dao, T. & Gu, A. Transformers are SSMS: Generalized Models and Efficient Algorithms Through  
1147 Structured State Space Duality. Preprint at <https://doi.org/10.48550/ARXIV.2405.21060> (2024).
- 1148 25. Brix, G. *et al.* Genome modeling and design across all domains of life with Evo 2. Preprint at  
1149 <https://doi.org/10.1101/2025.02.18.638918> (2025).
- 1150 26. Schiff, Y. *et al.* Caduceus: Bi-Directional Equivariant Long-Range DNA Sequence Modeling. Preprint  
1151 at <https://doi.org/10.48550/ARXIV.2403.03234> (2024).

- 1152 27. Nguyen, E. *et al.* HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide Resolu-  
1153 tion. Preprint at <https://doi.org/10.48550/ARXIV.2306.15794> (2023).
- 1154 28. Hayes, T. *et al.* Simulating 500 million years of evolution with a language model. *Science* **387**, 850–  
1155 858 (2025).
- 1156 29. Madani, A. *et al.* Large language models generate functional protein sequences across diverse families.  
1157 *Nat Biotechnol* **41**, 1099–1106 (2023).
- 1158 30. Nijkamp, E., Ruffolo, J. A., Weinstein, E. N., Naik, N. & Madani, A. ProGen2: Exploring the bounda-  
1159 ries of protein language models. *Cell Systems* **14**, 968-978.e3 (2023).
- 1160 31. Bhatnagar, A. *et al.* Scaling unlocks broader generation and deeper functional understanding of pro-  
1161 teins. Preprint at <https://doi.org/10.1101/2025.04.15.649055> (2025).
- 1162 32. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model.  
1163 *Science* **379**, 1123–1130 (2023).
- 1164 33. Shen, X.-X. *et al.* Tempo and Mode of Genome Evolution in the Budding Yeast Subphylum. *Cell* **175**,  
1165 1533-1545.e20 (2018).
- 1166 34. Hackett, S. R. *et al.* Learning causal networks using inducible transcription factors and transcriptome-  
1167 wide time series. *Molecular Systems Biology* **16**, e9174 (2020).
- 1168 35. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interac-  
1169 tions. *Nat Methods* **18**, 1196–1203 (2021).
- 1170 36. Linder, J., Srivastava, D., Yuan, H., Agarwal, V. & Kelley, D. R. Predicting RNA-seq coverage from  
1171 DNA sequence as a unifying model of gene regulation. *Nat Genet* (2025) doi:10.1038/s41588-024-  
1172 02053-6.
- 1173 37. Ronneberger, O., Fischer, P. & Brox, T. U-Net: Convolutional Networks for Biomedical Image Seg-  
1174 mentation. in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015* (eds.  
1175 Navab, N., Hornegger, J., Wells, W. M. & Frangi, A. F.) vol. 9351 234–241 (Springer International  
1176 Publishing, Cham, 2015).

- 1177 38. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of Phylo-  
1178 genomic Data. *Mol Biol Evol* **33**, 1635–1638 (2016).
- 1179 39. Letunic, I. & Bork, P. Interactive Tree of Life (iTOL) v6: recent updates to the phylogenetic tree dis-  
1180 play and annotation tool. *Nucleic Acids Research* **52**, W78–W82 (2024).
- 1181 40. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and  
1182 annotation. *Bioinformatics* **23**, 127–128 (2007).
- 1183 41. Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLoS Comput Biol* **14**,  
1184 e1005944 (2018).
- 1185 42. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using MinHash. *Genome*  
1186 *Biol* **17**, 132 (2016).
- 1187 43. Bailey, T. L. *et al.* MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Research* **37**,  
1188 W202–W208 (2009).
- 1189 44. Flynn, J. M. *et al.* RepeatModeler2 for automated genomic discovery of transposable element families.  
1190 *Proc. Natl. Acad. Sci. U.S.A.* **117**, 9451–9457 (2020).
- 1191 45. Smit, A. F. A., Hubley, R. & Green, P. RepeatMasker Open-4.0. (2013).
- 1192 46. Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community resource of trans-  
1193 posable element families, sequence models, and genome annotations. *Mobile DNA* **12**, 2 (2021).
- 1194 47. Ou, S. & Jiang, N. LTR\_retriever: A Highly Accurate and Sensitive Program for Identification of Long  
1195 Terminal Repeat Retrotransposons. *Plant Physiol.* **176**, 1410–1422 (2018).
- 1196 48. Price, A. L., Jones, N. C. & Pevzner, P. A. De novo identification of repeat families in large genomes.  
1197 *Bioinformatics* **21**, i351–i358 (2005).
- 1198 49. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
- 1199 50. Rafi, A. M., Kiyota, B., Yachie, N. & De Boer, C. Detecting and avoiding homology-based data leak-  
1200 age in genome-trained sequence models. Preprint at <https://doi.org/10.1101/2025.01.22.634321> (2025).

- 1201 51. Sahu, B. *et al.* Sequence determinants of human gene regulatory elements. *Nat Genet* **54**, 283–294  
1202 (2022).
- 1203 52. Tomaz Da Silva, P. *et al.* Nucleotide dependency analysis of DNA language models reveals genomic  
1204 functional elements. Preprint at <https://doi.org/10.1101/2024.07.27.605418> (2024).
- 1205 53. Shrikumar, A. *et al.* Technical Note on Transcription Factor Motif Discovery from Importance Scores  
1206 (TF-MoDISco) version 0.5.6.5. Preprint at <https://doi.org/10.48550/ARXIV.1811.00416> (2018).
- 1207 54. Shrikumar, A. *et al.* Technical Note on Transcription Factor Motif Discovery from Importance Scores  
1208 (TF-MoDISco) version 0.5.6.5. Preprint at <https://doi.org/10.48550/arXiv.1811.00416> (2020).
- 1209 55. De Boer, C. G. & Hughes, T. R. YeTFaSCo: a database of evaluated yeast transcription factor se-  
1210 quence specificities. *Nucleic Acids Research* **40**, D169–D179 (2012).
- 1211 56. MacIsaac, K. D. *et al.* An improved map of conserved regulatory sites for *Saccharomyces cerevisiae*.  
1212 *BMC Bioinformatics* **7**, (2006).
- 1213 57. Newburger, D. E. & Bulyk, M. L. UniPROBE: an online database of protein binding microarray data  
1214 on protein-DNA interactions. *Nucleic Acids Research* **37**, D77–D82 (2009).
- 1215 58. Pachkov, M., Erb, I., Molina, N. & Van Nimwegen, E. SwissRegulon: a database of genome-wide an-  
1216 notations of regulatory sites. *Nucleic Acids Research* **35**, D127–D131 (2007).
- 1217 59. Teixeira, M. C. *et al.* The YEASTRACT database: an upgraded information system for the analysis of  
1218 gene and genomic transcription regulation in *Saccharomyces cerevisiae*. *Nucl. Acids Res.* **42**, D161–  
1219 D166 (2014).
- 1220 60. Zhu, J. & Zhang, M. Q. SCPD: a promoter database of the yeast *Saccharomyces cerevisiae*. *Bioinform-*  
1221 *atics* **15**, 607–611 (1999).
- 1222 61. Yabana, N. & Yamamoto, M. *Schizosaccharomyces pombe map1<sup>+</sup>* Encodes a MADS-Box-Family Pro-  
1223 tein Required for Cell-Type-Specific Gene Expression. *Molecular and Cellular Biology* **16**, 3420–  
1224 3428 (1996).

- 1225 62. Nielsen, O., Friis, T. & Kjærulff, S. The Schizosaccharomyces pombe map1 gene encodes an SRF /  
1226 MCM1-related protein required for P-cell specific gene expression. *Mol Gen Genet* **253**, 387–392  
1227 (1996).
- 1228 63. Casselton, L. A. Mate recognition in fungi. *Heredity* **88**, 142–147 (2002).
- 1229 64. Oliva, A. *et al.* The Cell Cycle–Regulated Genes of Schizosaccharomyces pombe. *PLoS Biol* **3**, e225  
1230 (2005).
- 1231 65. Rossi, M. J. *et al.* A high-resolution protein architecture of the budding yeast genome. *Nature* **592**,  
1232 309–314 (2021).
- 1233 66. Caudal, É. *et al.* Pan-transcriptome reveals a large accessory genome contribution to gene expression  
1234 variation in yeast. *Nat Genet* **56**, 1278–1287 (2024).
- 1235 67. Martin, D. E., Soulard, A. & Hall, M. N. TOR Regulates Ribosomal Protein Gene Expression via PKA  
1236 and the Forkhead Transcription Factor FHL1. *Cell* **119**, 969–979 (2004).
- 1237 68. Schawalder, S. B. *et al.* Growth-regulated recruitment of the essential yeast ribosomal protein gene ac-  
1238 tivator Ifh1. *Nature* **432**, 1058–1061 (2004).
- 1239 69. Rudra, D., Zhao, Y. & Warner, J. R. Central role of Ifh1p–Fhl1p interaction in the synthesis of yeast  
1240 ribosomal proteins. *EMBO J* **24**, 533–542 (2005).
- 1241 70. Reja, R., Vinayachandran, V., Ghosh, S. & Pugh, B. F. Molecular mechanisms of ribosomal protein  
1242 gene coregulation. *Genes Dev.* **29**, 1942–1954 (2015).
- 1243 71. Shore, D. RAP1: a protean regulator in yeast. *Trends in Genetics* **10**, 408–412 (1994).
- 1244 72. Wade, J. T., Hall, D. B. & Struhl, K. The transcription factor Ifh1 is a key regulator of yeast ribosomal  
1245 protein genes. *Nature* **432**, 1054–1058 (2004).
- 1246 73. Warner, J. R. The economics of ribosome biosynthesis in yeast. *Trends in Biochemical Sciences* **24**,  
1247 437–440 (1999).
- 1248 74. Jorgensen, P. *et al.* A dynamic transcriptional network communicates growth potential to ribosome  
1249 synthesis and critical cell size. *Genes Dev.* **18**, 2491–2505 (2004).

- 1250 75. Wade, C. H., Umbarger, M. A. & McAlear, M. A. The budding yeast rRNA and ribosome biosynthesis  
1251 (RRB) regulon contains over 200 genes. *Yeast* **23**, 293–306 (2006).
- 1252 76. Arnone, J. T. & McAlear, M. A. Adjacent Gene Pairing Plays a Role in the Coordinated Expression of  
1253 Ribosome Biogenesis Genes *MPP10* and *YJR003C* in *Saccharomyces cerevisiae*. *Eukaryot Cell* **10**,  
1254 43–53 (2011).
- 1255 77. Hughes, J. D., Estep, P. W., Tavazoie, S. & Church, G. M. Computational identification of Cis -regula-  
1256 tory elements associated with groups of functionally related genes in *Saccharomyces cerevisiae* 1 1Ed-  
1257 ited by F. E. Cohen. *Journal of Molecular Biology* **296**, 1205–1214 (2000).
- 1258 78. Jorgensen, P., Nishikawa, J. L., Breikreutz, B.-J. & Tyers, M. Systematic Identification of Pathways  
1259 That Couple Cell Growth and Division in Yeast. *Science* **297**, 395–400 (2002).
- 1260 79. Brown, S. J., Cole, M. D. & Erives, A. J. Evolution of the holozoan ribosome biogenesis regulon.  
1261 *BMC Genomics* **9**, 442 (2008).
- 1262 80. Robinson, J. T., Thorvaldsdottir, H., Turner, D. & Mesirov, J. P. igv.js: an embeddable JavaScript im-  
1263 plementation of the Integrative Genomics Viewer (IGV). *Bioinformatics* **39**, btac830 (2023).
- 1264 81. Schirman, D., Yakhini, Z., Pilpel, Y. & Dahan, O. A broad analysis of splicing regulation in yeast us-  
1265 ing a large library of synthetic introns. *PLoS Genet* **17**, e1009805 (2021).
- 1266 82. Moore, M. J., Query, C. C., Sharp, P. A., & others. Splicing of precursors to mRNAs by the spliceo-  
1267 some. *Cold Spring Harbor Monograph Series* **24**, 303–303 (1993).
- 1268 83. Parker, R., Siliciano, P. G. & Guthrie, C. Recognition of the TACTAAC box during mRNA splicing in  
1269 yeast involves base pairing to the U2-like snRNA. *Cell* **49**, 229–239 (1987).
- 1270 84. Zavanelli, M. I. & Ares, M. Efficient association of U2 snRNPs with pre-mRNA requires an essential  
1271 U2 RNA structural element. *Genes & Development* **5**, 2521–2533 (1991).
- 1272 85. Engel, S. R. *et al.* *Saccharomyces* Genome Database: advances in genome annotation, expanded bio-  
1273 chemical pathways, and other key enhancements. *GENETICS* **229**, iyae185 (2025).

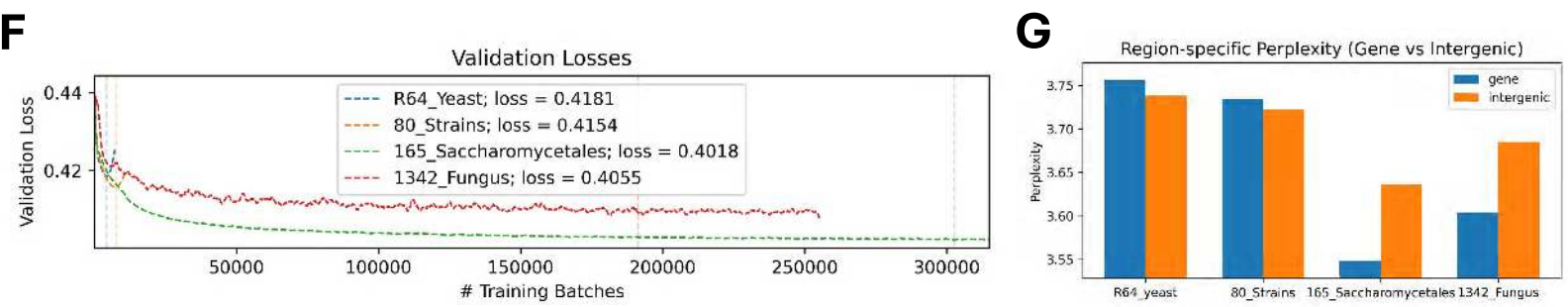
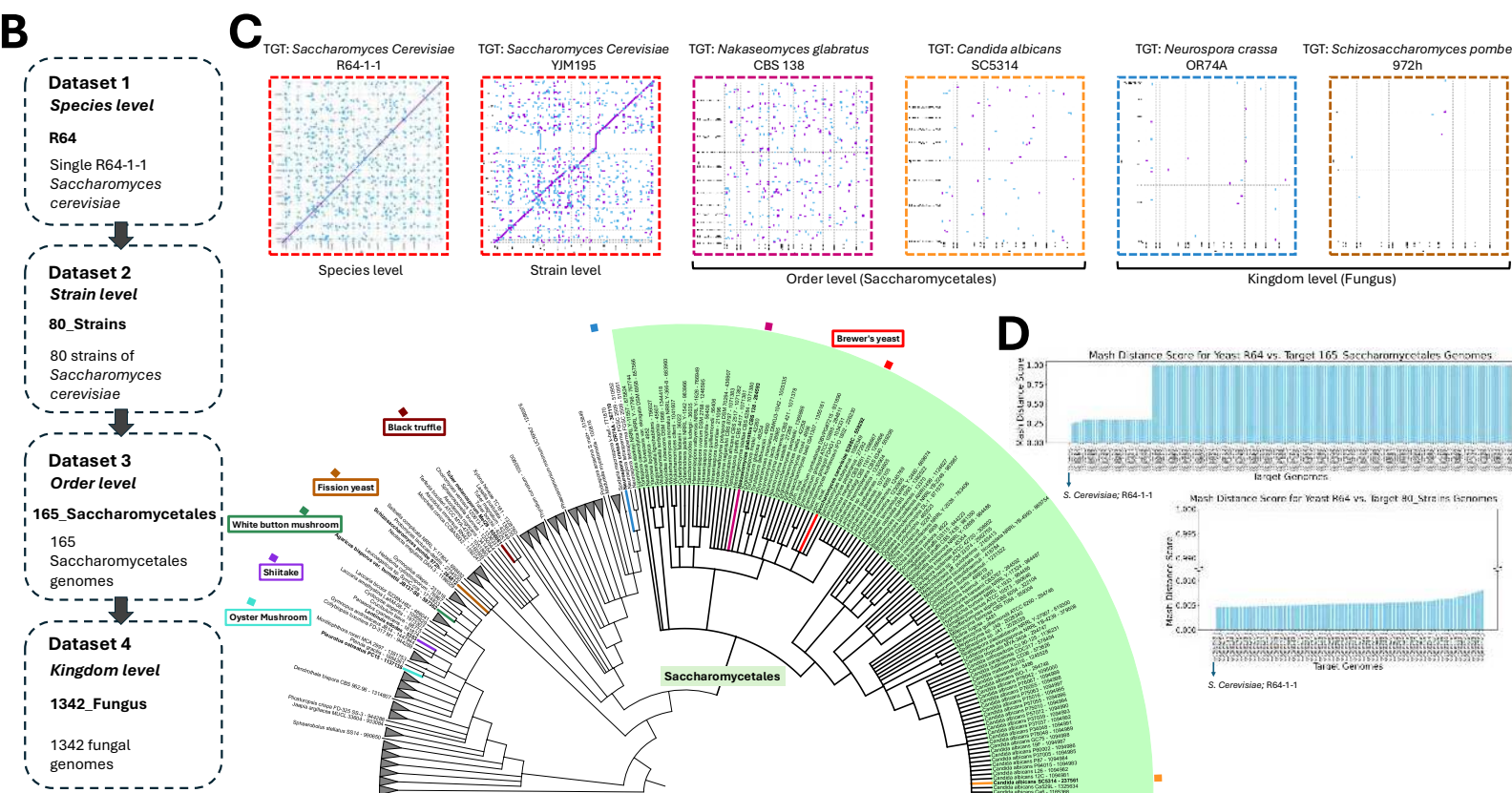
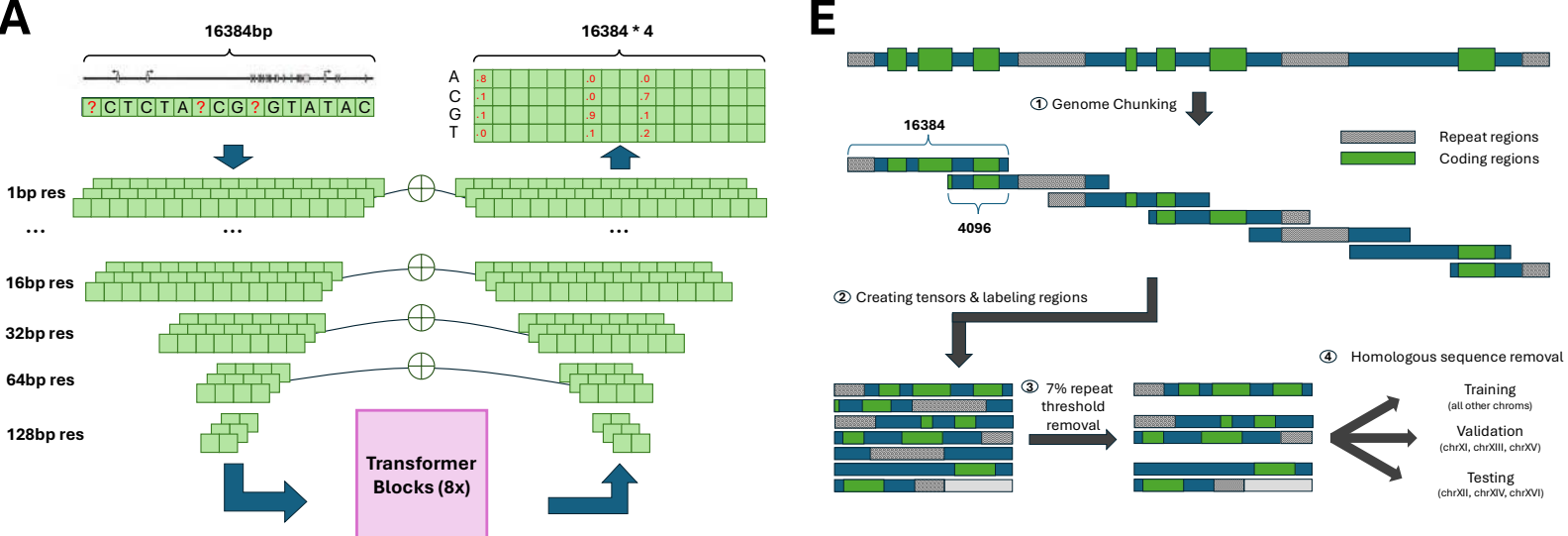
- 1274 86. Garmendia-Torres, C., Goldbeter, A. & Jacquet, M. Nucleocytoplasmic Oscillations of the Yeast Tran-  
1275 scription Factor Msn2: Evidence for Periodic PKA Activation. *Current Biology* **17**, 1044–1049 (2007).
- 1276 87. Ni, L. *et al.* Dynamic and complex transcription factor binding during an inducible response in yeast.  
1277 *Genes Dev.* **23**, 1351–1363 (2009).
- 1278 88. Martínez-Pastor, M. T. *et al.* The *Saccharomyces cerevisiae* zinc finger proteins Msn2p and Msn4p are  
1279 required for transcriptional induction through the stress response element (STRE). *The EMBO Journal*  
1280 **15**, 2227–2235 (1996).
- 1281 89. Gasch, A. P. *et al.* Genomic Expression Programs in the Response of Yeast Cells to Environmental  
1282 Changes. *MBoC* **11**, 4241–4257 (2000).
- 1283 90. Causton, H. C. *et al.* Remodeling of Yeast Genome Expression in Response to Environmental  
1284 Changes. *MBoC* **12**, 323–337 (2001).
- 1285 91. Kuras, L., Barbey, R. & Thomas, D. Assembly of a bZIP-bHLH transcription activation complex: for-  
1286 mation of the yeast Cbfl-Met4-Met28 complex is regulated through Met28 stimulation of Cbfl DNA  
1287 binding. *The EMBO Journal* **16**, 2441–2451 (1997).
- 1288 92. Lee, T. A. *et al.* Dissection of Combinatorial Control by the Met4 Transcriptional Complex. *MBoC* **21**,  
1289 456–469 (2010).
- 1290 93. Rafi, A. M. *et al.* A community effort to optimize sequence-based deep learning models of gene regu-  
1291 lation. *Nat Biotechnol* (2024) doi:10.1038/s41587-024-02414-w.
- 1292 94. Kita, R., Venkataram, S., Zhou, Y. & Fraser, H. B. High-resolution mapping of *cis* -regulatory varia-  
1293 tion in budding yeast. *Proc. Natl. Acad. Sci. U.S.A.* **114**, (2017).
- 1294 95. Peter, J. *et al.* Genome evolution across 1,011 *Saccharomyces cerevisiae* isolates. *Nature* **556**, 339–344  
1295 (2018).
- 1296 96. Kuras, L. *et al.* Dual Regulation of the Met4 Transcription Factor by Ubiquitin-Dependent Degradation  
1297 and Inhibition of Promoter Recruitment. *Molecular Cell* **10**, 69–80 (2002).

- 1298 97. Leroy, C., Cormier, L. & Kuras, L. Independent Recruitment of Mediator and SAGA by the Activator  
1299 Met4. *Molecular and Cellular Biology* **26**, 3149–3163 (2006).
- 1300 98. Lin, L., Chamberlain, L., Zhu, L. J. & Green, M. R. Analysis of Gal4-directed transcription activation  
1301 using Tra1 mutants selectively defective for interaction with Gal4. *Proc. Natl. Acad. Sci. U.S.A.* **109**,  
1302 1997–2002 (2012).
- 1303 99. Chandrasekaran, S. & Skowyra, D. The emerging regulatory potential of SCFMet30 -mediated  
1304 polyubiquitination and proteolysis of the Met4 transcriptional activator. *Cell Div* **3**, 11 (2008).
- 1305 100. Hendrycks, D. & Gimpel, K. Gaussian Error Linear Units (GELUs). Preprint at  
1306 <https://doi.org/10.48550/ARXIV.1606.08415> (2016).
- 1307 101. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol* **5**, R12  
1308 (2004).
- 1309 102. Baker, D. N. & Langmead, B. Genomic sketching with multiplicities and locality-sensitive hashing  
1310 using Dashing 2. *Genome Res.* gr.277655.123 (2023) doi:10.1101/gr.277655.123.
- 1311 103. Ertl, O. SetSketch: Filling the Gap between MinHash and HyperLogLog. (2021)  
1312 doi:10.48550/ARXIV.2101.00314.
- 1313 104. Ertl, O. ProbMinHash – A Class of Locality-Sensitive Hash Algorithms for the (Probability) Jac-  
1314 card Similarity. *IEEE Trans. Knowl. Data Eng.* 1–1 (2020) doi:10.1109/TKDE.2020.3021176.
- 1315 105. Broder, A. Z. On the resemblance and containment of documents. in *Proceedings. Compression*  
1316 *and Complexity of SEQUENCES 1997 (Cat. No.97TB100171)* 21–29 (IEEE Comput. Soc, Salerno, It-  
1317 aly, 1998). doi:10.1109/SEQUEN.1997.666900.
- 1318 106. Jaccard, P. Étude comparative de la distribution florale dans une portion des Alpes et des Jura. *Bull*  
1319 *Soc Vaudoise Sci Nat* **37**, 547–579 (1901).
- 1320 107. Wheeler, T. J. *et al.* Dfam: a database of repetitive DNA based on profile hidden Markov models.  
1321 *Nucleic Acids Research* **41**, D70–D82 (2012).
- 1322 108. Pertea, G. & Pertea, M. GFF Utilities: GffRead and GffCompare. *F1000Res* **9**, 304 (2020).

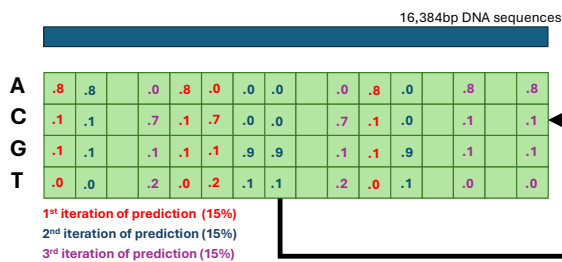
- 1323 109. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: as-  
1324 sessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**,  
1325 3210–3212 (2015).
- 1326 110. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional  
1327 Transformers for Language Understanding. Preprint at <https://doi.org/10.48550/ARXIV.1810.04805>  
1328 (2018).
- 1329 111. Mallet, V. & Vert, J.-P. Reverse-complement equivariant networks for DNA sequences. *Advances*  
1330 *in neural information processing systems* **34**, 13511–13523 (2021).
- 1331 112. Zhou, H., Shrikumar, A. & Kundaje, A. Towards a better understanding of reverse-complement  
1332 equivariance for deep learning models in genomics. in *Machine Learning in Computational Biology* 1–  
1333 33 (PMLR, 2022).
- 1334 113. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. Preprint at  
1335 <https://doi.org/10.48550/ARXIV.1412.6980> (2014).
- 1336 114. Schneider, T. D., Stormo, G. D., Gold, L. & Ehrenfeucht, A. Information content of binding sites  
1337 on nucleotide sequences. *Journal of Molecular Biology* **188**, 415–431 (1986).
- 1338 115. Schneider, T. D. & Stephens, R. M. Sequence logos: a new way to display consensus sequences.  
1339 *Nucl Acids Res* **18**, 6097–6100 (1990).
- 1340 116. Wolf, T. *et al.* HuggingFace’s Transformers: State-of-the-art Natural Language Processing. Pre-  
1341 print at <https://doi.org/10.48550/ARXIV.1910.03771> (2019).
- 1342 117. Hinton, G. E. & Roweis, S. Stochastic neighbor embedding. *Advances in neural information pro-*  
1343 *cessing systems* **15**, (2002).
- 1344 118. Arita, Y. *et al.* A genome-scale yeast library with inducible expression of individual genes. *Molec-*  
1345 *ular Systems Biology* **17**, e10207 (2021).
- 1346 119. Balakrishnan, R. *et al.* YeastMine—an integrated data warehouse for *Saccharomyces cerevisiae*  
1347 data as a multipurpose tool-kit. *Database* **2012**, (2012).

- 1348 120. Turco, G. *et al.* Global analysis of the yeast knockout phenome. *Sci. Adv.* **9**, eadg5702 (2023).
- 1349 121. Hou, J. *et al.* The Hidden Complexity of Mendelian Traits across Natural Yeast Populations. *Cell*  
1350 *Reports* **16**, 1106–1114 (2016).
- 1351 122. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint  
1352 at <https://doi.org/10.48550/ARXIV.1303.3997> (2013).
- 1353 123. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079  
1354 (2009).
- 1355 124. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.  
1356 *Bioinformatics* **26**, 841–842 (2010).
- 1357 125. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 1358 126. Picard toolkit. *Broad Institute, GitHub repository* (2019).
- 1359 127. Stovner, E. B. & Sætrom, P. PyRanges: efficient comparison of genomic intervals in Python. *Bioin-*  
1360 *formatics* **36**, 918–919 (2020).

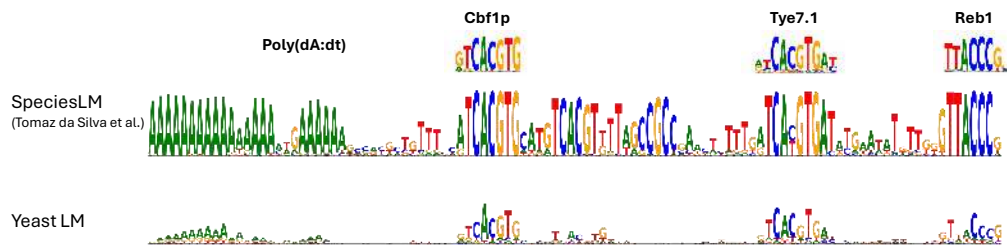
1361



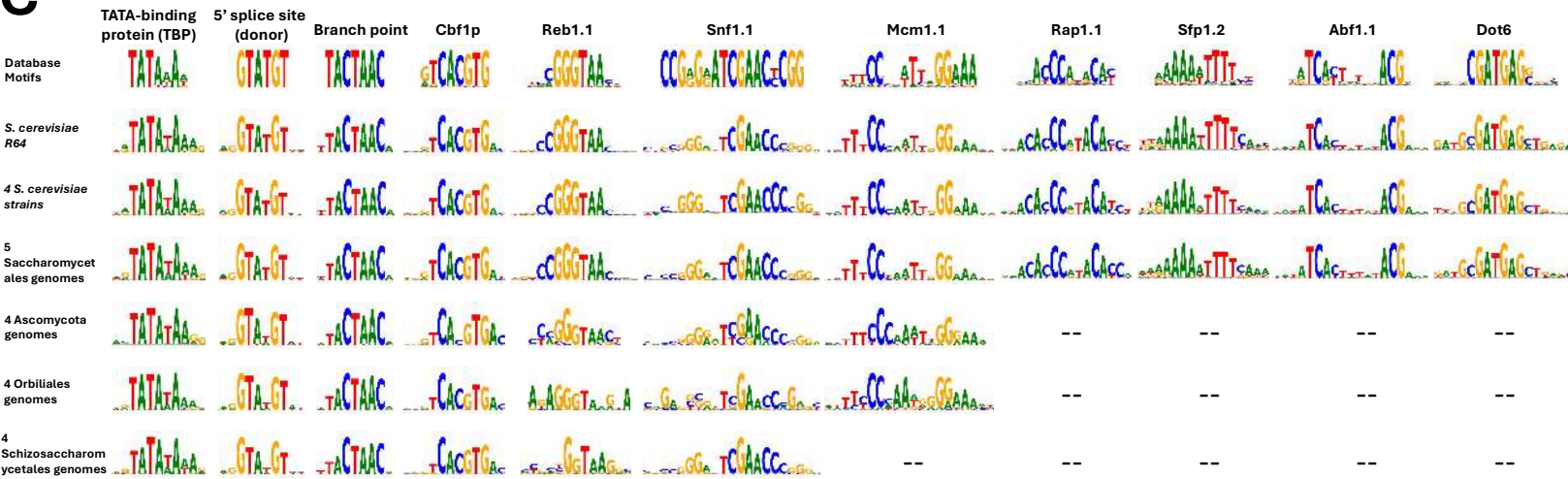
# A Predicting 15% masked regions for each iteration



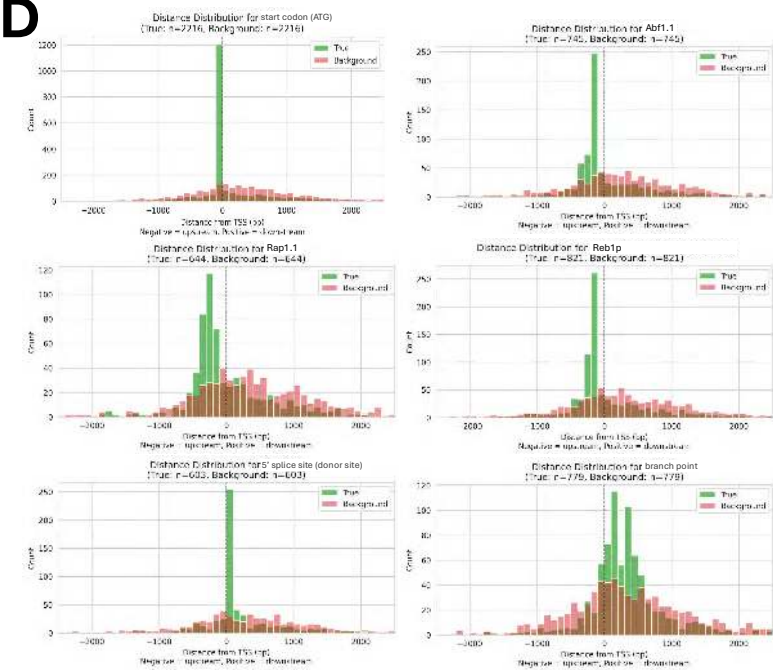
# B



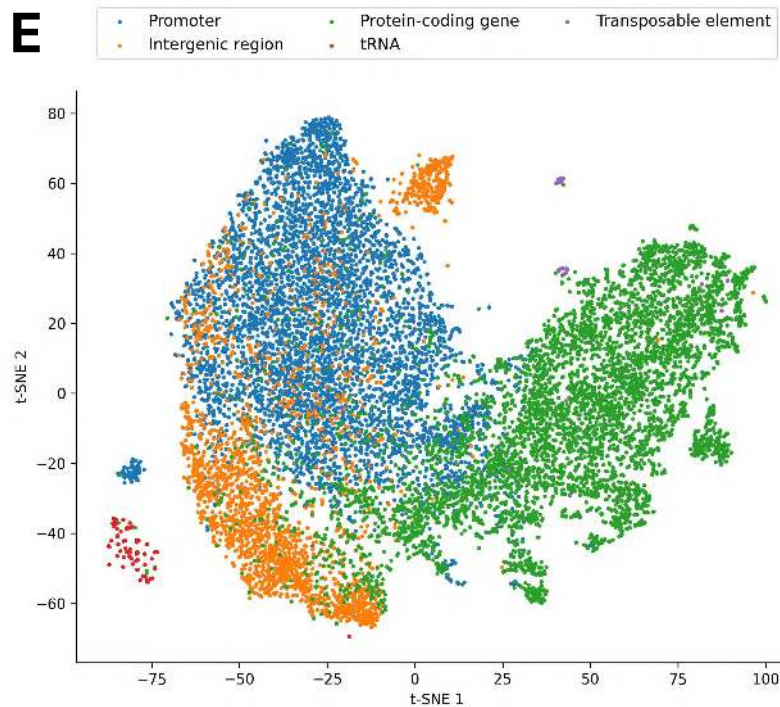
# C

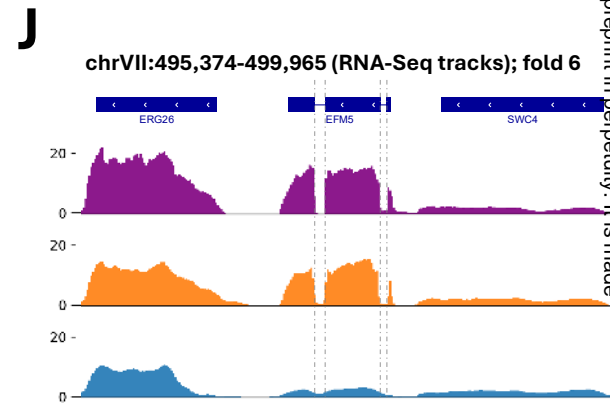
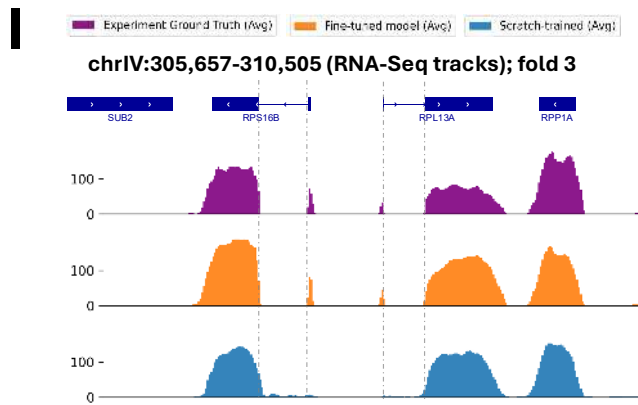
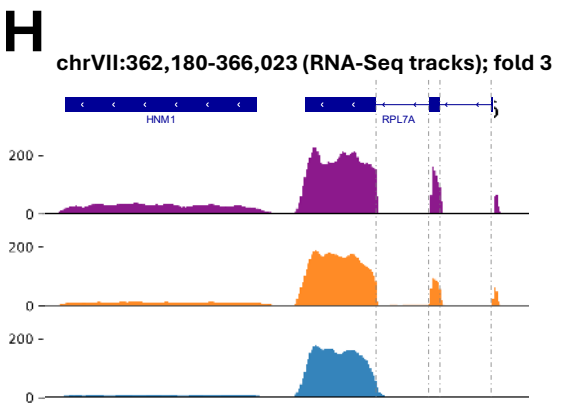
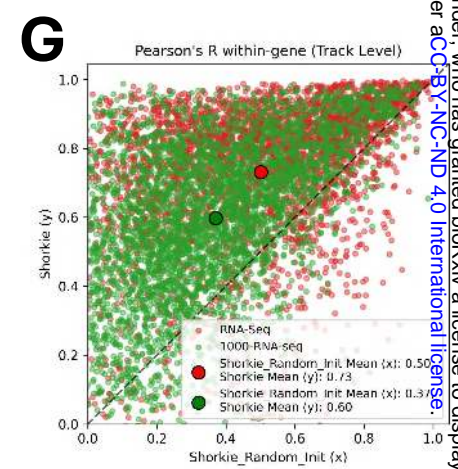
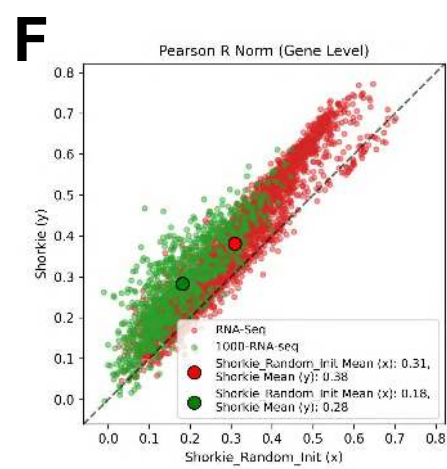
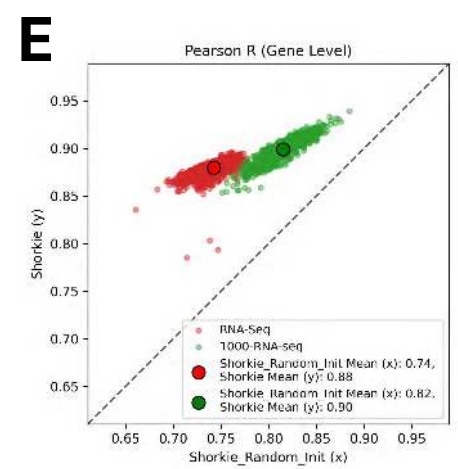
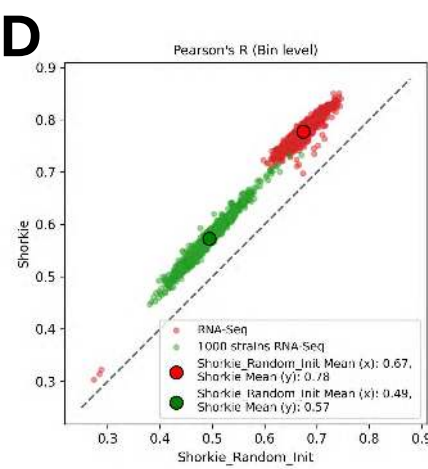
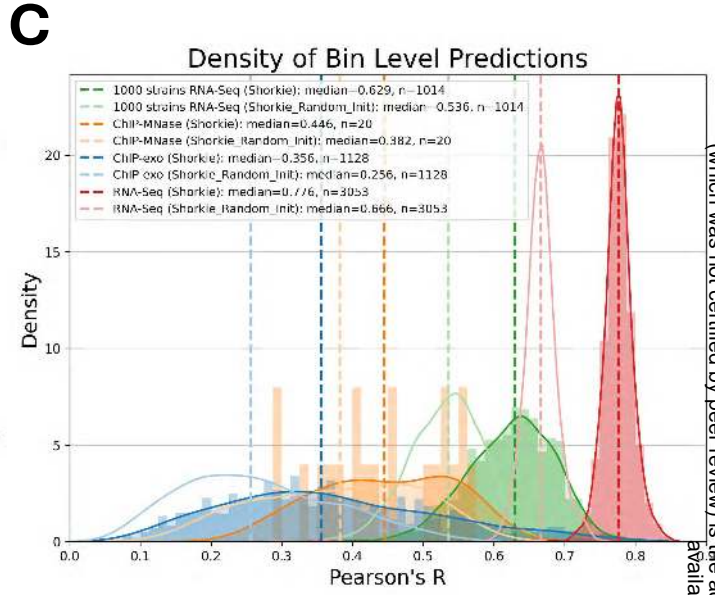
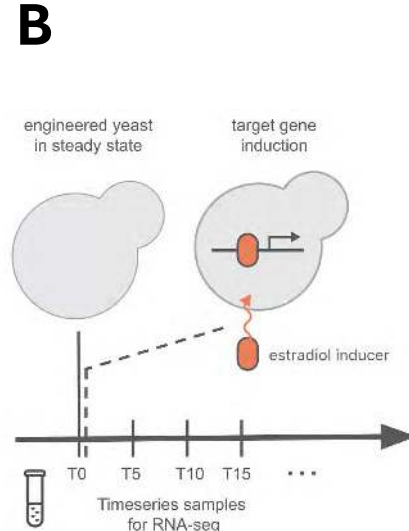
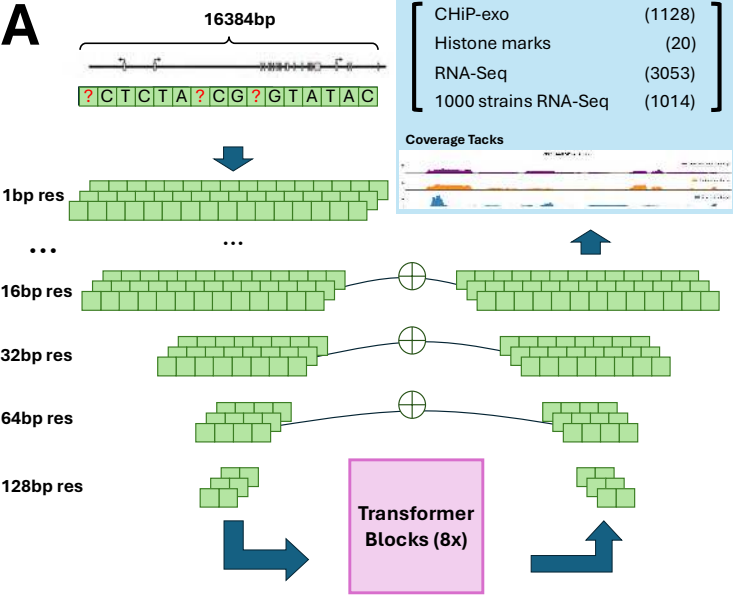


# D



# E





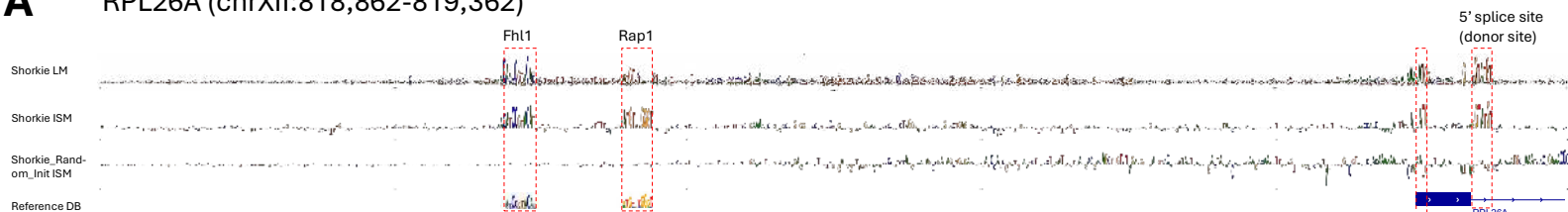
is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

450 nt

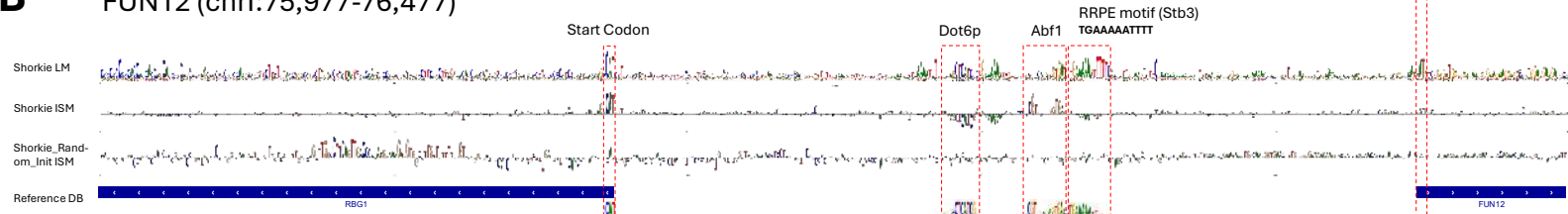
50 nt

**A**

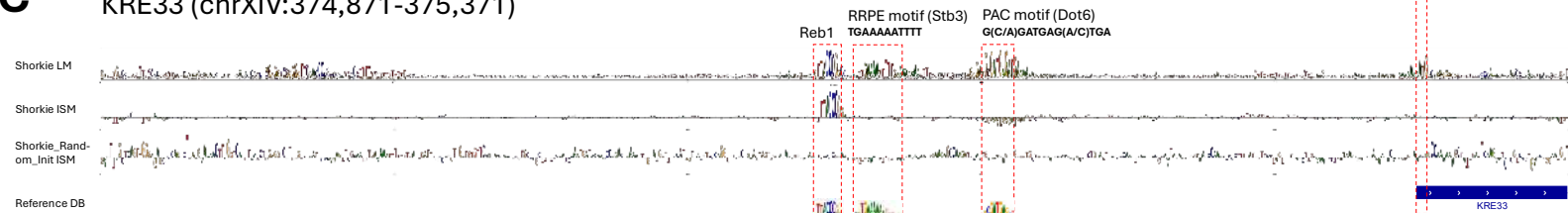
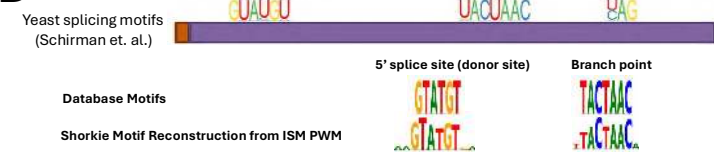
RPL26A (chrXII:818,862-819,362)

**B**

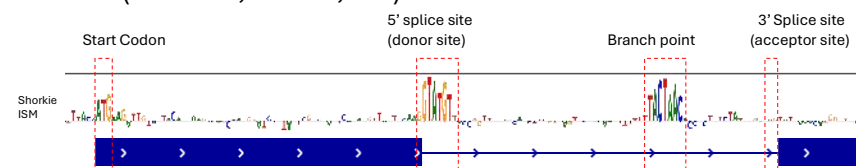
FUN12 (chrI:75,977-76,477)

**C**

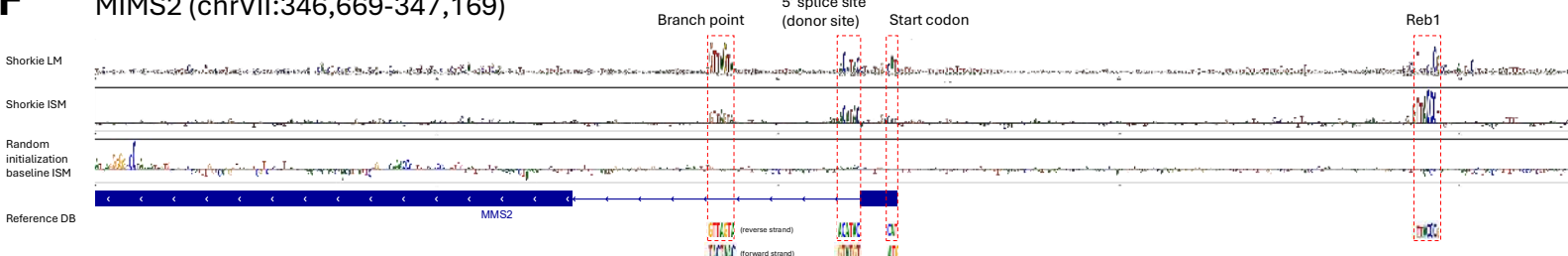
KRE33 (chrXIV:374,871-375,371)

**D****E**

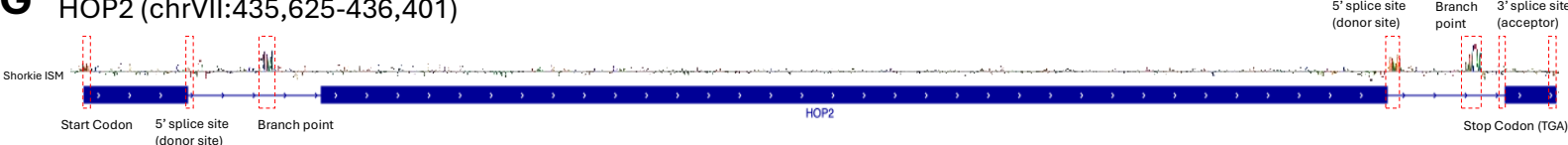
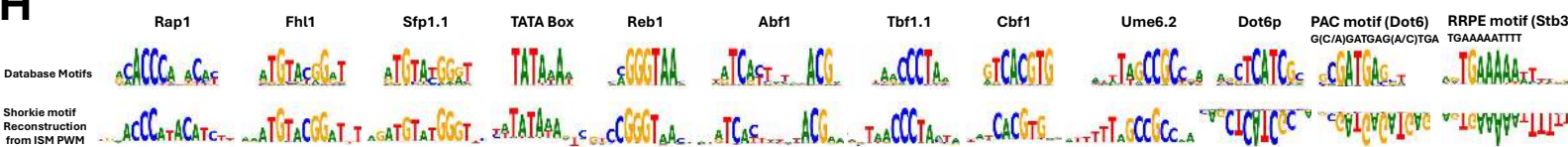
DTD1 (chrIV:65,235-65,431)

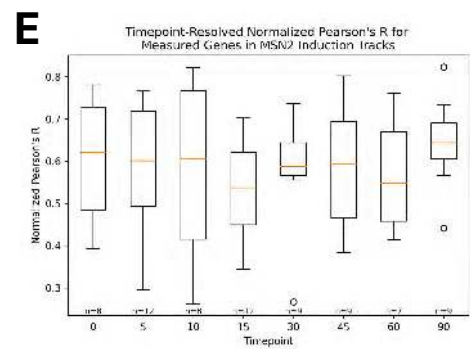
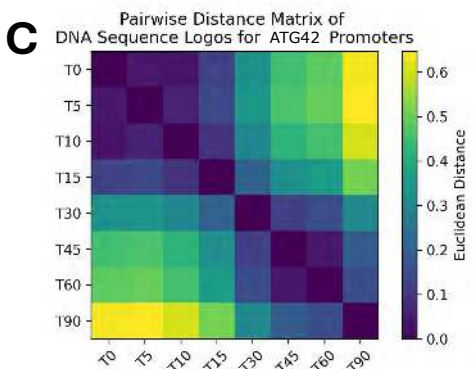
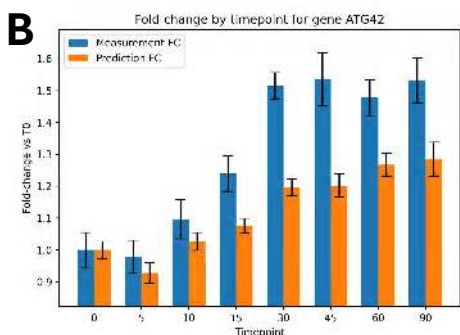
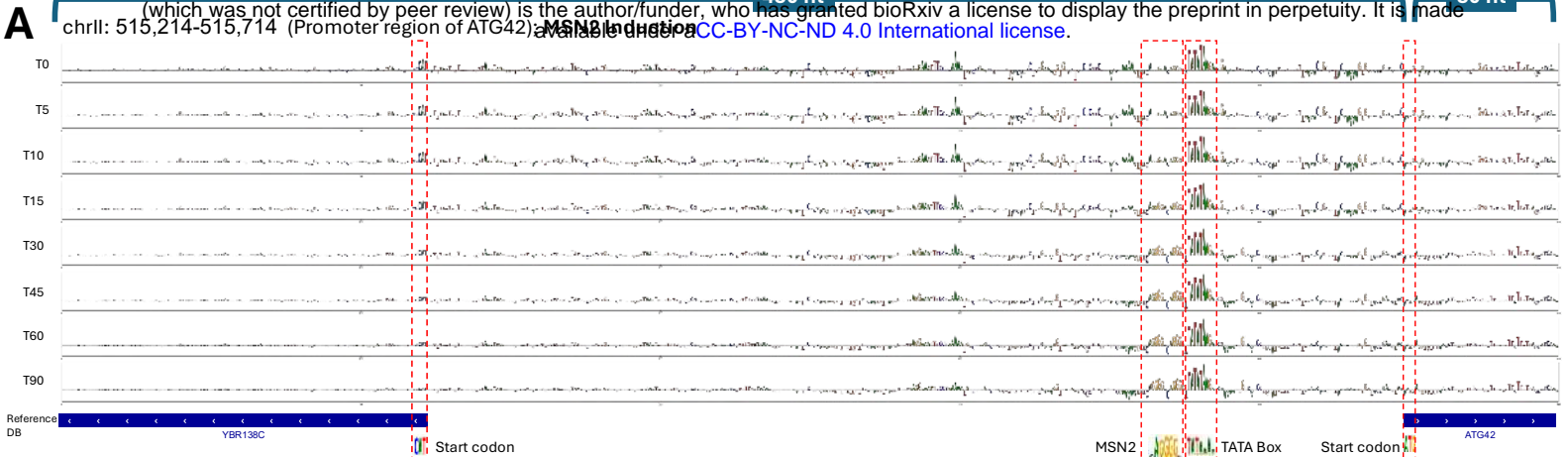
**F**

MIMS2 (chrVII:346,669-347,169)

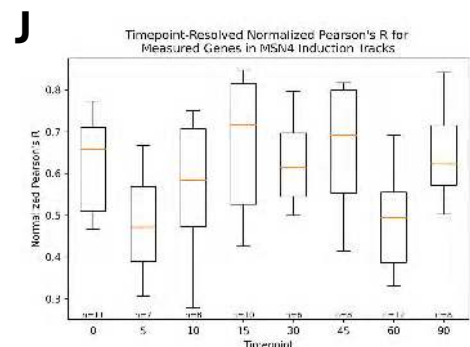
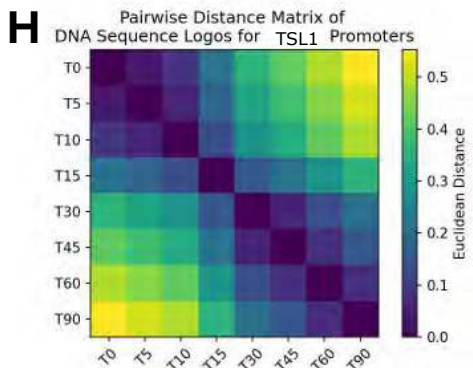
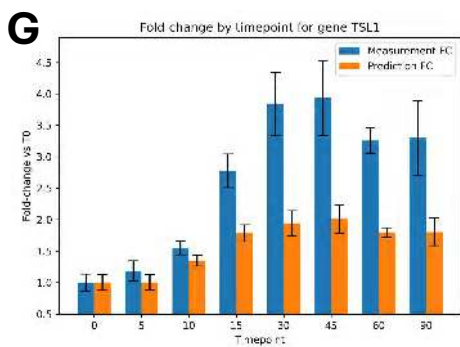
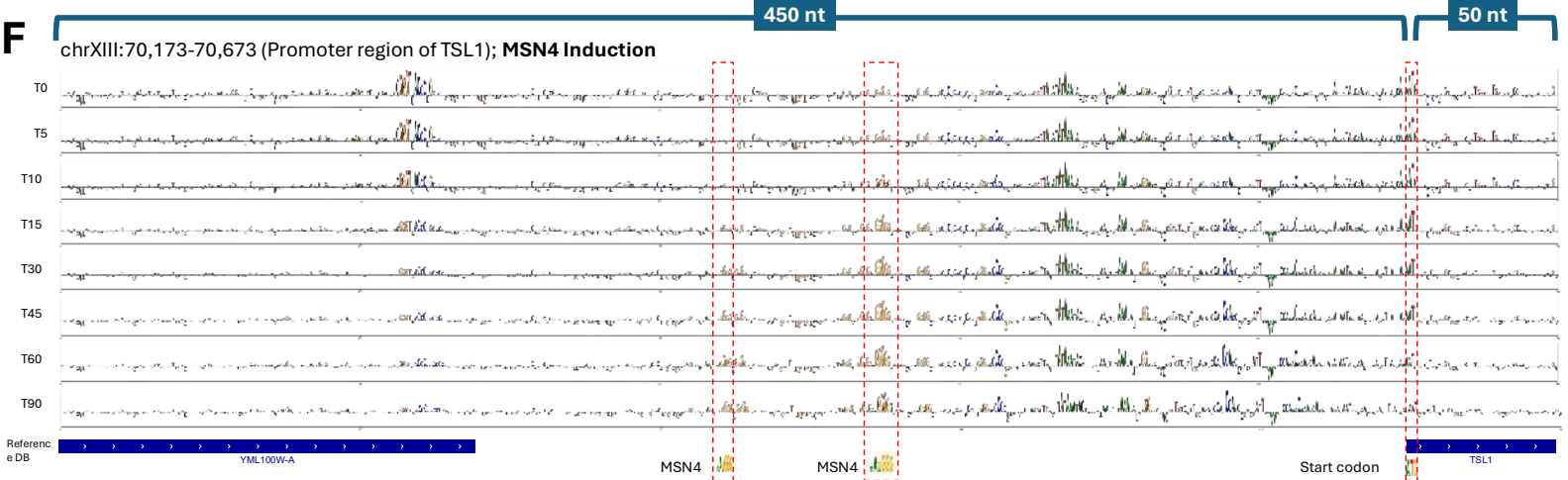
**G**

HOP2 (chrVII:435,625-436,401)

**H**

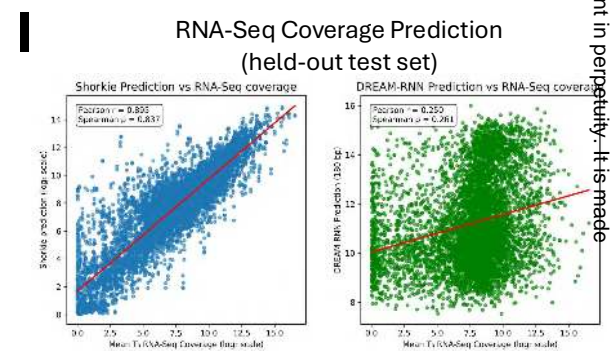
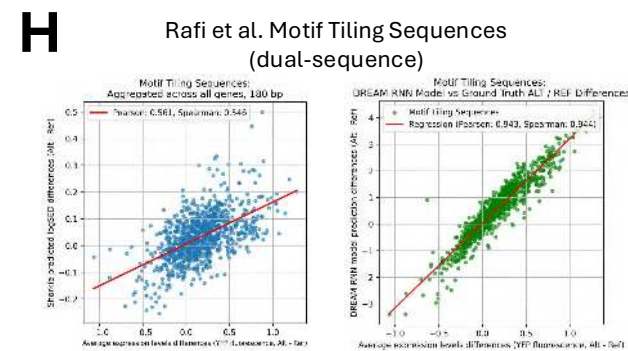
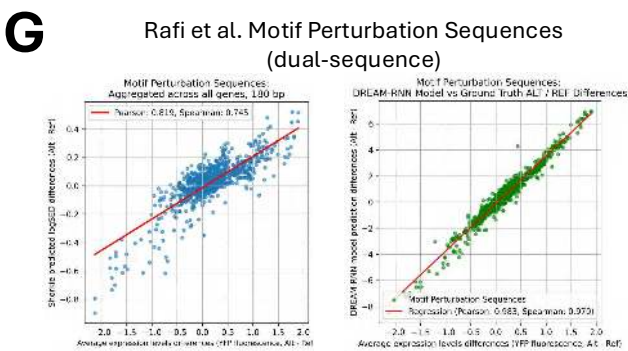
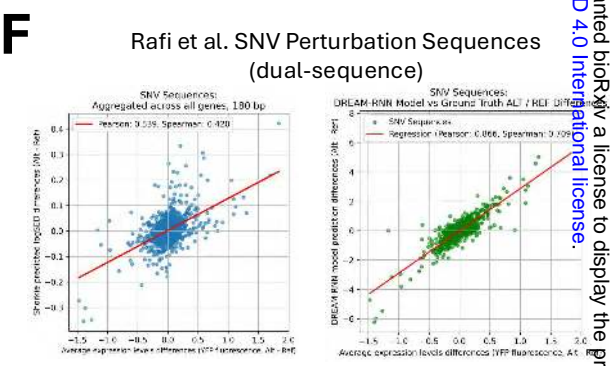
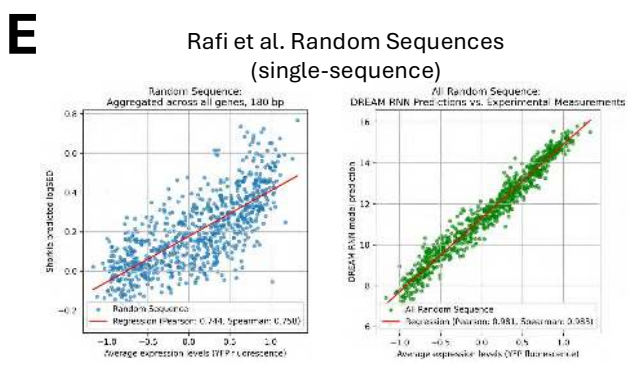
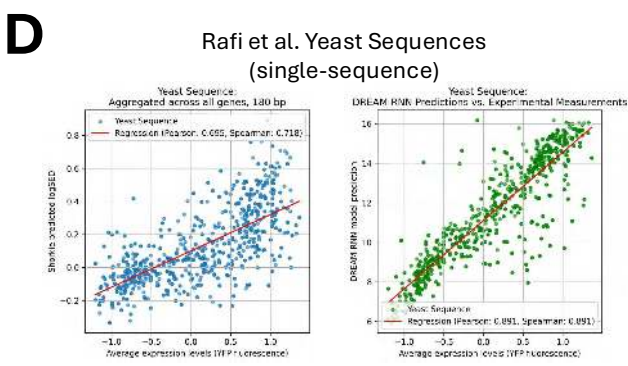
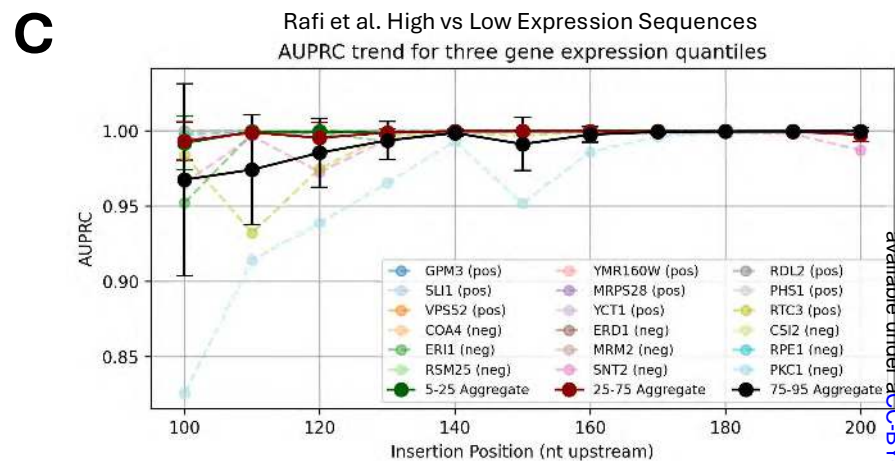
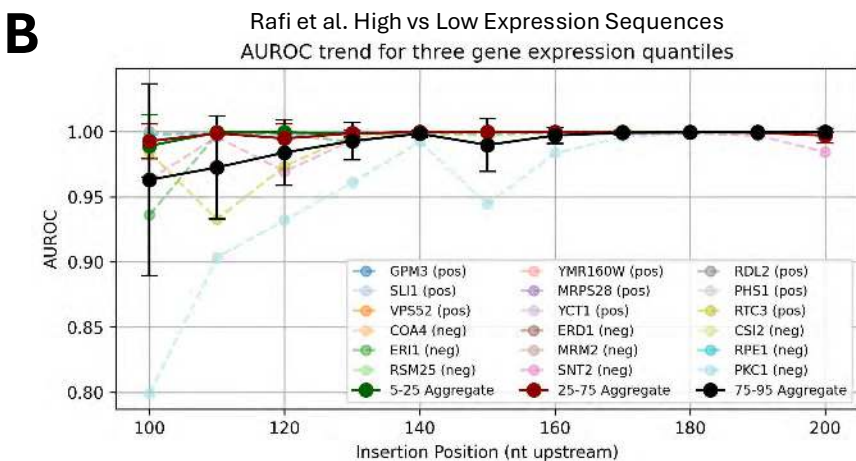
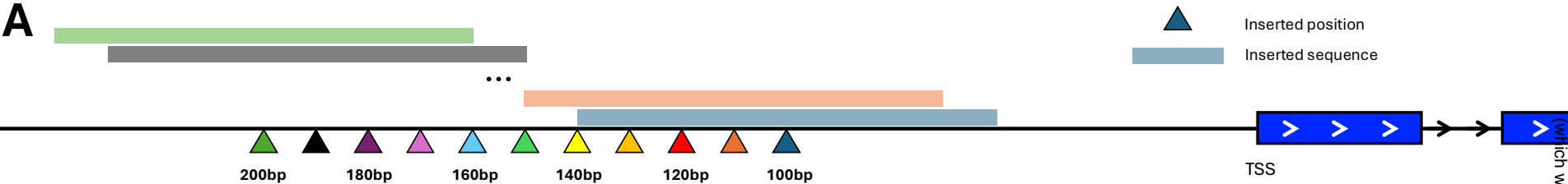


**D** MSN2 Induction



**I** MSN4 Induction





available under aCC-BY-NC-ND 4.0 International license.

