

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Conservation of function without conservation of amino acid sequence in intrinsically disordered transcriptional activation domains

### Permalink

<https://escholarship.org/uc/item/55g3r44b>

### Journal

bioRxiv, 5(12-13)

### ISSN

2692-8205

### Authors

LeBlanc, Claire

Stefani, Jordan

Soriano, Melvin

et al.

### Publication Date

2024-12-05

### DOI

10.1101/2024.12.03.626510

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial-NoDerivatives License, available at

<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Peer reviewed

1

1 **Conservation of function without conservation of amino acid sequence**  
2 **in intrinsically disordered transcriptional activation domains**

3

4

5 Claire LeBlanc<sup>1,2</sup>, Jordan Stefani<sup>1,2</sup>, Melvin Soriano<sup>1,2</sup>, Angelica Lam<sup>1,2,#</sup>, Marissa A.  
6 Zintel<sup>1</sup>, Sanjana R. Kotha<sup>1,2</sup>, Emily Chase<sup>1,2</sup>, Giovanni Pimentel-Solorio<sup>1,2,##</sup>, Aditya  
7 Vunnum<sup>1</sup>, Katherine Flug<sup>1</sup>, Aaron Fultineer<sup>3</sup>, Niklas Hummel<sup>4</sup>, Max V. Staller<sup>1,2,5,\*</sup>

8

9

10 **Affiliations:**

11 <sup>1</sup> Department of Molecular and Cell Biology, University of California Berkeley,  
12 Berkeley, 94720

13 <sup>2</sup> Center for Computational Biology, University of California Berkeley, Berkeley,  
14 94720

15 <sup>3</sup> Department of Physics, University of California Berkeley, Berkeley, 94720

16 <sup>4</sup> Department of Biology, Technische Universität Darmstadt, Darmstadt, Germany

17 <sup>5</sup> Chan Zuckerberg Biohub–San Francisco, San Francisco, CA 94158

18 # Present address: University of California San Francisco, San Francisco, CA  
19 94158

20 ## Present address: University of California Davis, Davis, CA

21

22 \*Corresponding author: 16 Barker Hall, Berkeley, CA 94720, USA.

23 [mstaller@berkeley.edu](mailto:mstaller@berkeley.edu)

## 24 Abstract:

25 Protein function is canonically believed to be more conserved than amino  
26 acid sequence, but this idea is only well supported in folded domains, where  
27 highly diverged sequences can fold into equivalent 3D structures. In contrast,  
28 intrinsically disordered protein regions (IDRs) do not fold into a stable 3D  
29 structure, thus it remains unknown when and how function is conserved for IDRs  
30 that experience rapid amino acid sequence divergence. As a model system for  
31 studying the evolution of IDRs, we examined transcriptional activation domains,  
32 the regions of transcription factors that bind to coactivator complexes. We  
33 systematically identified activation domains on 502 orthologs of the  
34 transcriptional activator Gcn4 spanning 600 MY of fungal evolution. We find that  
35 the central activation domain shows strong conservation of function without  
36 conservation of sequence. This conservation of function without conservation of  
37 sequence is facilitated by evolutionary turnover (gain and loss) of key acidic and  
38 aromatic residues, the positions most important for function. This high sequence  
39 flexibility of functional orthologs mirrors the physical flexibility of the activation  
40 domain coactivator interaction interface, suggesting that physical flexibility  
41 enables evolutionary plasticity. We propose that turnover of short functional  
42 elements, sometimes individual amino acids, is a general mechanism for  
43 conservation of function without conservation of sequence during IDR evolution.  
44

## 45 Key words

46 Intrinsically disordered proteins; transcription; transcription factor;  
47 activation domains; evolution; evolutionary turnover; high-throughput assays  
48

## 49 Introduction:

50 The evolution of eukaryotic transcription factor (TF) function contains a  
51 paradox: TF protein sequences diverge quickly but maintain function over long  
52 evolutionary distances. For example, the master regulator of eye development in  
53 mice, Pax6, induces ectopic eyes in fly, and fly Pax6 (*eye/less*) creates ectopic  
54 eye structures in frogs and mice<sup>1-3</sup>. While the DNA-binding domains (DBD) are  
55 96% identical, eye induction requires the intrinsically disordered regions (IDRs),  
56 which are only 35.5% identical. These IDRs must share a conserved function  
57 despite substantial sequence divergence. In contrast, small sequence changes in  
58 TFs can lead to large functional changes that drive the evolution of new traits<sup>4,5</sup>.  
59 Some TFs maintain function despite low conservation of sequence<sup>6</sup>, while other  
60 TFs drive evolutionary innovations with limited sequence changes.

61 For folded domains, function is more conserved than sequence because  
62 highly diverged sequences can fold into the same 3D structure and maintain  
63 function<sup>7-9</sup>. Here, we seek an analogous framework for understanding the  
64 evolution of and functional constraint on IDRs. Small-scale studies have found  
65 examples of diverged IDRs that conserve function<sup>10-12</sup> and diverged IDRs that do  
66 not conserve function<sup>13,14</sup>. Transcriptional activation domains provide an  
67 excellent model system for studying IDR evolution because they are one of the

68 oldest classes of functional IDRs<sup>15</sup>, they are required for TF function, and their  
69 activity can be measured in high throughput<sup>16</sup>. Our goal is to identify molecular  
70 mechanisms by which TF IDR function can be conserved in the face of rapid  
71 sequence divergence.

72 We hypothesized that TF IDRs can maintain function despite sequence  
73 divergence through evolutionary turnover of functional elements. Evolutionary  
74 turnover is repeated gain and loss of functional elements. Mutations create new  
75 functional elements and negative selection maintains a minimum number of  
76 elements, allowing ancestral elements to be lost. As a result, on long timescales,  
77 neutral drift will give the appearance of functional elements moving around the  
78 sequence. For TFs, it is unclear if the functional elements will be entire activation  
79 domains, short linear interaction motifs (SLiMs)<sup>17</sup>, or individual amino acids. Here,  
80 we aim to identify the functional units and test the hypothesis that evolutionary  
81 turnover can explain conservation of function without conservation of sequence.

82 Evolutionary studies of acidic activation domains in yeast benefit from  
83 high-throughput data that define sequence features controlling their function<sup>16,18-  
84 23</sup>. These data have trained neural network models for predicting activation  
85 domains from protein sequence<sup>18,21,23-26</sup>. Our acidic exposure model further  
86 provides a biophysical mechanism for the observed features: aromatic and  
87 leucine residues make key contacts with hydrophobic surfaces of coactivator  
88 complexes, but these residues can also interact with each other and drive  
89 collapse into an inactive state<sup>16,27-30</sup>. The acidic residues repel each other, expand  
90 the activation domain, and promote exposure of the hydrophobic residues. In  
91 many cases, the aromatic and leucine residues are arranged into short linear  
92 motifs. Large-scale mutagenesis showed the acidic exposure model applies to  
93 hundreds of human activation domains<sup>31</sup>.

94 We investigated the molecular mechanisms by which full-length TFs can  
95 maintain activator function over long evolutionary distances despite divergence  
96 of their amino acid sequences. As a model system, we used 502 diverse  
97 orthologs of Gcn4, a nutrient stress TF, and screened for activation domains with  
98 a high-throughput functional assay in *Saccharomyces cerevisiae*<sup>16</sup>. All orthologs  
99 contain at least one 40 AA region that functions as an activation domain, and we  
100 see widespread conservation of function without conservation of sequence. We  
101 demonstrate evolutionary turnover of entire activation domains and turnover of  
102 key residues within an activation domain. The N-terminal activation domains are  
103 repeatedly gained and lost. In contrast, the central activation domain is  
104 functionally conserved because of turnover of key acidic and hydrophobic  
105 residues. This work illustrates how functional screening can unravel the complex  
106 evolution of activation domains and IDRs.

107

## 108 Results:

### 109 **Characterization of a tiling-library of Gcn4 orthologs**

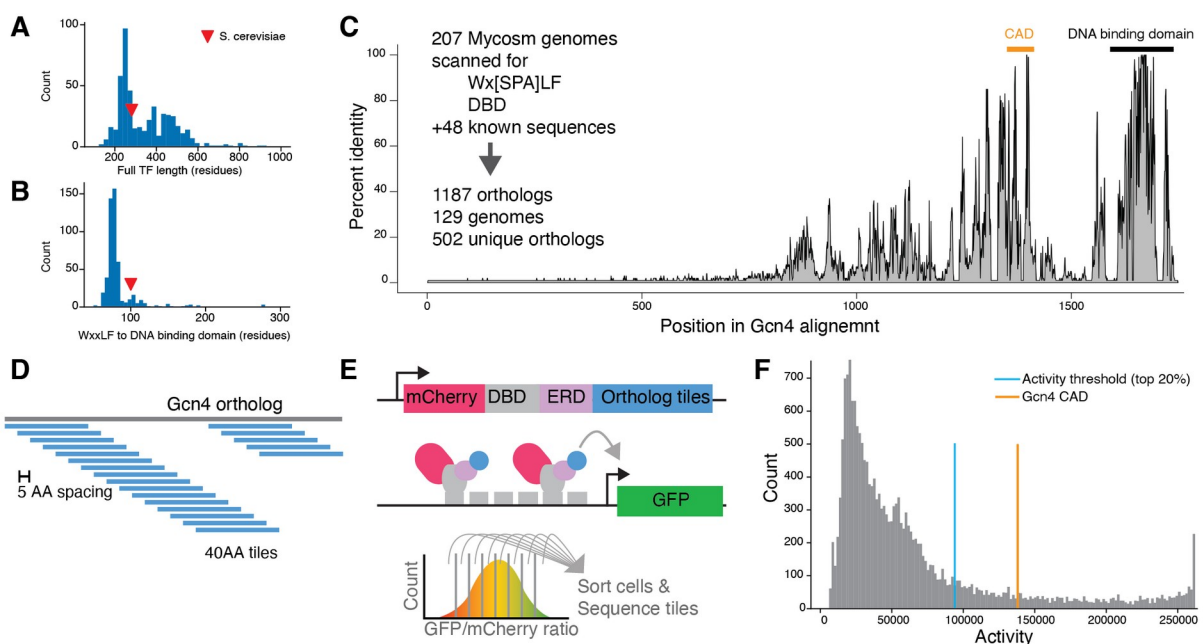
110 To study the evolutionary dynamics of activator function, we sought to  
111 experimentally map activation domains across a diverse collection of  
112 orthologous TFs. We and others have shown that protein fusion libraries,  
113 designed to tile across protein sequences with short, 30-60 amino acid peptides,

114 can faithfully measure activation domain activity<sup>16,18,19,21,22,32</sup>. Furthermore,  
 115 because activation domain function in yeast is a reliable measure of endogenous  
 116 function in humans<sup>33</sup>, viruses<sup>34</sup>, *Drosophila*<sup>35,36</sup>, plants<sup>23,37,38</sup>, and other yeast  
 117 species<sup>39</sup>, we reasoned that the activity of fungal orthologues in our assay would  
 118 serve as a reliable measure of activity in their native context. In all subsequent  
 119 analysis, we assume that tile activity measured in *S. cerevisiae* is a good proxy  
 120 for TF function in their native species.

121 As a null hypothesis, we assumed the TF function is conserved and that  
 122 the observed diversity of sequence is the result of neutral drift. Absent strong  
 123 evidence to the contrary, neutral drift is a strong null hypothesis<sup>40</sup>. Mutation  
 124 processes introduce changes, and selection acts at the level of the full protein.  
 125 Purifying (negative) selection will tolerate all changes that do not reduce function  
 126 below a minimum level. The neutral space for IDRs is potentially much larger  
 127 than that of folded proteins because there are no structural constraints.  
 128 Supporting this assumption, we found evidence for weak negative selection on  
 129 the full-length TF using a high-quality set of thirty-six true Gcn4 homologs from  
 130 the yeast gene order browser (**Figure S1E**)<sup>41</sup>. It follows that most of the  
 131 sequence differences we see in extant species are neutral. We aim to find the  
 132 (potentially rare or diffuse) sequence features that are functional and conserved.

133 We chose a diverse set of orthologous Gcn4 protein sequences for  
 134 functional characterization in *S. cerevisiae*. We found 502 unique Gcn4 ortholog  
 135 sequences from 129 genomes that span the Ascomycota, the largest phylum of  
 136 Fungi, representing >600 million years of evolution<sup>42</sup> (**Figure S1, S2**). While the  
 137 Gcn4 orthologs vary in length (**Figure 1A**), 500 have the DBD at the C-terminus,  
 138 and the distance between the WxxLF motif and the DBD is very consistent  
 139 (**Figure 1B**).

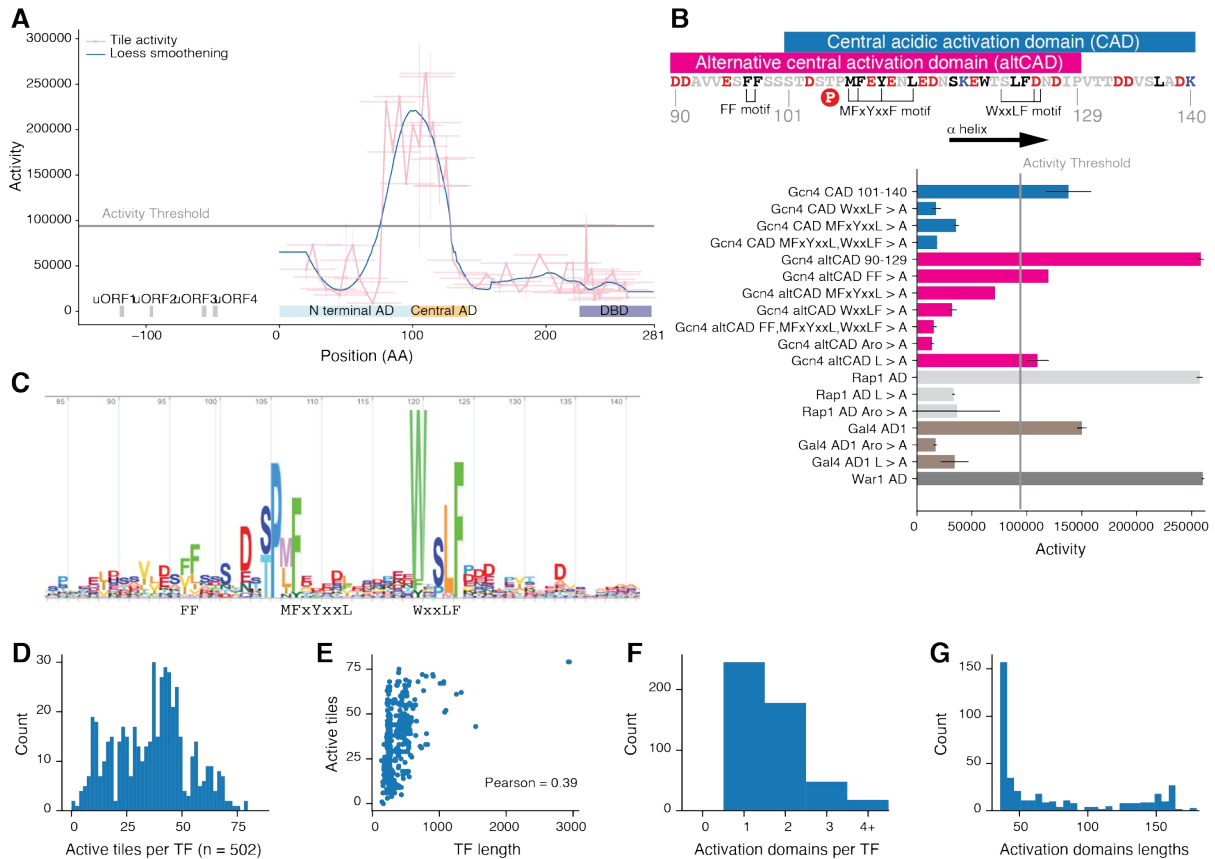
140  
 141



142  
 143 **Figure 1: Screening fragments of Gcn4 orthologs for activation domain**  
 144 **activity in *S. cerevisiae*.**

145 A) Gcn4 ortholog lengths. Red arrow, *S. cerevisiae*. B) The distance between the  
 146 WxxLF motif and the start of the DBD is conserved. C) The MSA of 500 orthologs  
 147 shows the DBD binding domain is highly conserved, and the Central Activation  
 148 Domain around the WxxLF motif is moderately conserved. D) The tiling strategy  
 149 for oligo design and the high-throughput activation domain assay. E) The high-  
 150 throughput assay for measuring activation domain function uses a synthetic TF  
 151 with mCherry for quantification of abundance, the Zif268 DNA binding domain  
 152 (DBD), an estrogen response domain (ERD) for inducible activation, and a C-  
 153 terminally fused tile. Tile activity was calculated based on barcode abundance in  
 154 eight equally sized bins of a FACS sorting experiment. Bins were set based on  
 155 GFP/mCherry ratios. F) The distribution of measured tile activities with our  
 156 activity threshold (top 20%). *S. cerevisiae* Gcn4 CAD activity is shown in orange.

159 The Gcn4 multiple sequence alignment (MSA) typifies eukaryotic TF  
 160 evolution, with a highly conserved DBD and lower conservation in the rest of the  
 161 protein (**Figure 1C**). The central activation domain (CAD) shows intermediate  
 162 levels of conservation, driven in part by the WxxLF motif (**Figure 2B, S3**).  
 163 Sequence divergence is driven by insertions: 54% of columns in the MSA contain  
 164 fewer than 1% of sequences (**Figure 1C, S4**). Distant pairs of sequences do not  
 165 align outside of the DBD.  
 166



167

168 **Figure 2: In the *S. cerevisiae* central activation domain, residues that**  
 169 **are critical for activity are poorly conserved.**

170 **A)** Schematic of *S. cerevisiae* Gcn4 with the upstream open reading frames  
 171 (uORFs) that regulate translation, the NAD and the CAD. Individual measured  
 172 tiles are indicated as pink lines with a pink point at the center, and the standard  
 173 deviation of the two replicates is shown vertically. We imputed activity at each  
 174 position with a Loess smoothing (blue). **B)** Schematic of the CAD and altCAD  
 175 (most active tile) with key motifs,  $\alpha$ -helix, and phosphosites indicated. Mutating  
 176 motifs, aromatic residues, or leucine residues reduced activity in all cases. **C)**  
 177 The sequence logo from the 4th iteration of a search for Gcn4 orthologs in fungal  
 178 genomes with HMMER. This independent analysis confirmed the WxxLF motif is  
 179 more conserved than the FF and MFxYxxL motifs. **D)** The number of active tiles  
 180 found on each full-length TF (tiles that map to multiple orthologs can count  
 181 multiple times in this analysis). **E)** There is a weak correlation between TF length  
 182 and the number of active tiles. **F-G)** Combining overlapping active tiles shows  
 183 that most TFs have 2 or more activation domains with a wide distribution of  
 184 lengths.  
 185

186 **High-throughput measurement of orthologs for activation domain**  
 187 **function**

188 To study the evolution of TF function, we measured the activation domain  
 189 activity of all the orthologs. For each of the 502 Gcn4 orthologs, we tiled across  
 190 the full-length protein with 40 AA tiles spaced every 5 AA, and measured  
 191 activities of all tiles in *S. cerevisiae* using our established high-throughput  
 192 assay<sup>16</sup> (**Figure 1D, 1E**). We recovered 18947 of 20731 designed tiles (91.4%),  
 193 and these data were of high quality (Methods, **Figure S5, S6**). The tiles had a  
 194 range of activities (**Figure 1F**), and mutations in control activation domains  
 195 behaved as expected (**Figure 2A, 2B, S7**). As a threshold for highly-active tiles,  
 196 we used the top 20% of sequences, but other thresholds led to similar results  
 197 (Methods, **Figure S8**). Many more tiles are active than datasets that naively tile  
 198 all TFs in a proteome, as we would expect if most Gcn4 orthologs are activators.  
 199 Due to the divergence of the orthologs, the sequences of the active tiles are very  
 200 diverse, allowing us to study sequence-to-function relationships controlling  
 201 activation domain function. To our knowledge, this dataset is the largest  
 202 functional study of TF evolution to date.

## 203 **Activator function is conserved across the Gcn4 orthologs**

204 All the Gcn4 orthologs had at least one tile that functioned as an activation  
 205 domain in our assay, indicating that activator function is conserved across 600  
 206 million years of evolution (**Figure 2D**, Supplemental note 1). *A priori*, it was not  
 207 a given that all the Gcn4 orthologs would be activators, because on long  
 208 evolutionary timescales, a family of TFs that share a conserved DBD will include  
 209 both activators and repressors<sup>23,31,32,38</sup>. Gcn4 activator function is highly  
 210 conserved despite divergence of the sequence.

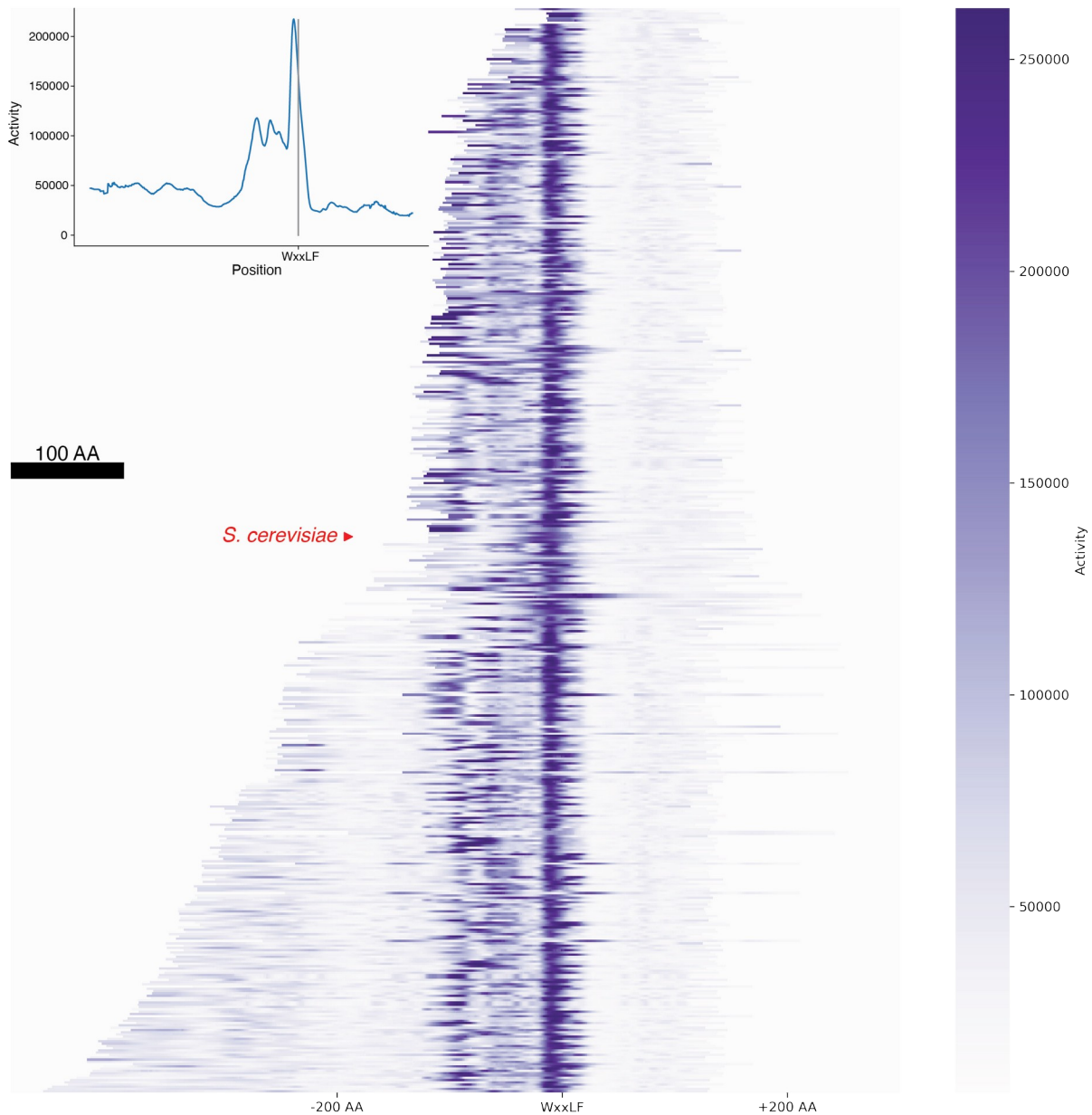
## 211 **The central acidic activation domain shows strong functional** 212 **conservation.**

213 Our finding that all the orthologs are activators combined with the  
 214 sequence divergence in the MSA indicates there is conservation of function  
 215 without conservation of primary amino acid sequence. We examined three  
 216 hypotheses for this conservation of function without conservation of sequence:  
 217 1) turnover of entire activation domains, 2) turnover of motifs within activation  
 218 domains, and 3) turnover of key residues within activation domains. We found  
 219 turnover of entire N-terminal activation domains and turnover of key residues  
 220 within the central activation domain.

221 The central activation domain is functionally conserved across the  
 222 orthologs. An advantage of our tiling strategy is the ability to infer the activity of  
 223 each position in each full-length protein (**Figure 3**, Methods). We found that all  
 224 orthologs had high activity in the central region (Supplemental note 1). The peak  
 225 of activity is ten AA residues upstream of the WxxLF motif (**Figure 3, inset**).  
 226 Aligning on the WxxLF motif or the DBD led to similar results (**Figure S9-S12**).  
 227 Projecting the activity heatmap onto the local species tree or gene tree  
 228 illustrates how the central activation domain can drift side-to-side but stays near  
 229 the WxxLF motif (**Figure S13, S14**). Intriguingly, the integral of activity across  
 230 each ortholog was highly consistent, suggesting conservation of total activity  
 231 (**Figure S15C**).

232 The second major result is that N-terminal activation domains come and  
 233 go, providing evidence for turnover of entire activation domains (**Figure 3**). After  
 234 combining overlapping active tiles, the majority of orthologs have more than one  
 235 activation domain (**Figure 2F**). Projecting activity onto the MSA or sorting the  
 236 heatmap by activity at the WxxLF motif emphasizes how the N-terminal  
 237 activation domains come and go (**Figure S11, S12**). Using our stringent  
 238 threshold for activity (top 20%), thirteen orthologs lost activity at the WxxLF  
 239 motif, but all of these have gained additional upstream activation domains. The  
 240 N-terminal activation domains show intermediate conservation in the MSA  
 241 (**Figure S15**) and their sequences are very diverse, ruling out the possibility that  
 242 one ancestral activation domain is recurrently lost (**Figure S16**). Together, these  
 243 data demonstrate turnover of entire activation domains.  
 244





245

246 **Figure 3: The central acidic activation domain of Gcn4 is functionally**  
 247 **conserved.**

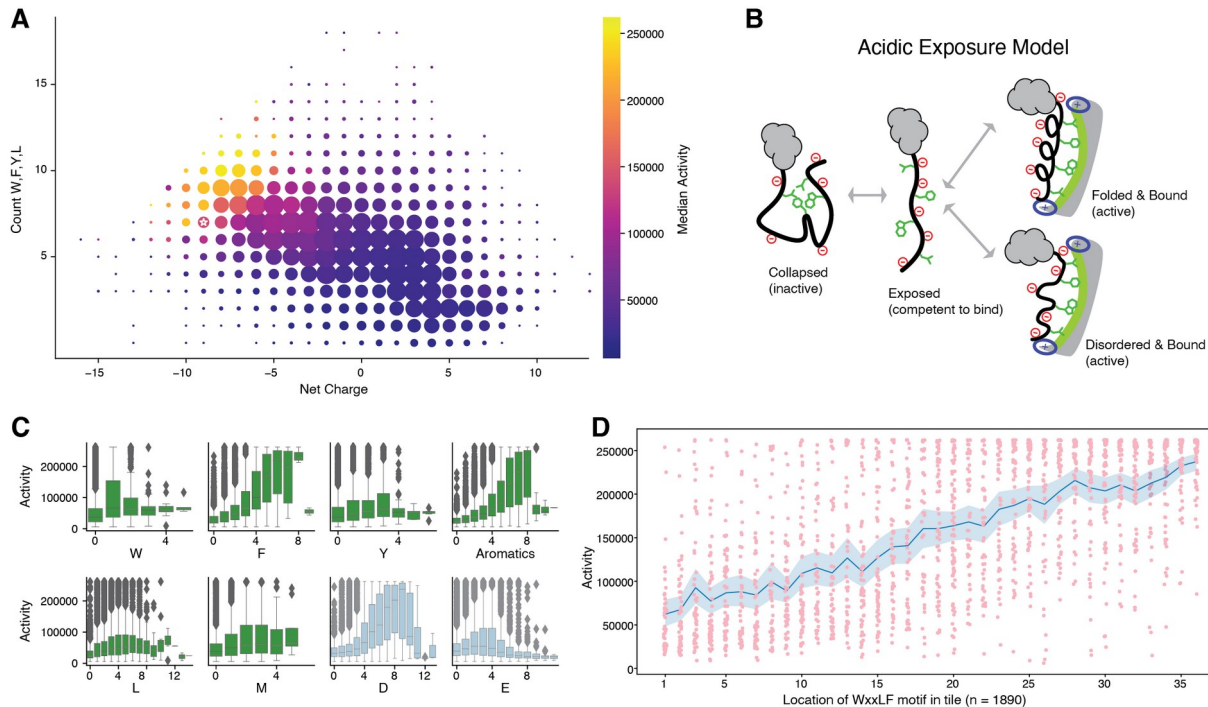
248 We used the tile activity data to impute the activity of each residue in all the  
 249 orthologs. These activities are visualized as a heatmap, with color representing  
 250 imputed activity. The 476 shortest orthologs are sorted by length and aligned on  
 251 the WxxLF motif. Inset, vertically averaging the heatmap. Activity is consistently  
 252 high around the WxxLF motif, indicating deep functional conservation. Upstream,  
 253 N-terminal activity is more salt and pepper, indicating recurrent gain and loss of  
 254 activation domains. Aligning on the DBD or including the longer sequences yields  
 255 similar results (**Figure S9-11**). Red arrow, *S. cerevisiae*. Black scale bar, 100 AA.  
 256

257 **Conservation of function without conservation of sequence in the**  
 258 **Central Acidic Activation Domain of Gcn4**

259 The central activation domain region with high-functional conservation  
 260 shows intermediate conservation in the multiple-sequence alignment (**Figure**

261 **1C, S3A**). We conclude that there is conservation of activation domain function  
 262 without conservation of the sequence. To understand the sequence features  
 263 underlying this conservation of function without conservation of sequence, we  
 264 first describe the amino acid sequence features controlling activity of individual  
 265 tiles and then apply these lessons to the orthologs.

266  
 267



268  
 269 **Figure 4: Highly active tiles contain many acidic, aromatic, and leucine**  
 270 **residues, supporting the acid exposure model of acidic activation**  
 271 **domain function.**

272 **A)** For each tile, we compute net charge and count the number of WFYL residues.  
 273 The size of the point indicates the number of tiles with the combination of  
 274 properties. The color is the median activity of tiles with each combination. White  
 275 star, *S. cerevisiae* Gcn4. **B)** The acidic exposure model of acidic activation  
 276 domain function. **C)** Boxplots for the residues that make the largest contributions  
 277 to activity. **D)** For each tile with the WxxLF motif, activity is plotted against the  
 278 location of the W. Blue, mean and 95% confidence interval. The location of the  
 279 motif is correlated with activity.

280

## 281 **The sequence features of active tiles support the acidic-exposure** 282 **model**

283 The Gcn4 ortholog dataset contains all previously observed relationships  
 284 between sequence and function, but many relationships are stronger and more  
 285 visible than previously reported (Supplemental Note 1). As predicted by the  
 286 acidic exposure model, many active tiles contain both acidic residues and WFYL  
 287 residues (**Figure 4A, 4B**). These key residues make quantitatively different  
 288 contributions to activity (**Figure 4C, S17, S18**). Aspartic acid (D) makes  
 289 stronger contributions to activity than glutamic acid (E), likely because the

charge is slower to the backbone and better promotes exposure<sup>43</sup> (**Figure 4C, S18**). In the control activation domains, all published motifs of aromatic and leucine residues made large contributions to activity, but no individual motif was sufficient for full activity (**Figure S7**). These sequence features of active tiles with or without the WxxLF motif are highly similar, suggesting the N-terminal activation domains function similarly to the central activation domain, as has been shown in *S. cerevisiae*<sup>44</sup> (**Figure 19**). Tiling orthologs reveals sequence rules more efficiently than tiling genomes (**Figure S20**).

Amino acid composition strongly contributes to activation domain function. Ordinary least squares (OLS) regression on single amino acids explains 49.9% of variance in activity (**Table 1**, AUC = 0.9346, PRC = 0.7620, **Table S9**). Regression on dipeptides<sup>21</sup> led to 69 significant parameters that explain 60.2% of the variance in activity (**Table 1**, AUC = 0.9472, PRC = 0.8190). More complex sequence motifs did not improve the regression models: published motifs explained 33.1%, and 40 *de novo* motifs explained 50.5% of the variance in activity (**Table 1**). Combining the *de novo* motifs with single amino acids performed similarly to dipeptides. This result implies that complex motifs capture very little additional information beyond adjacent pairwise amino acid relationships in dipeptides.

309

Model	Number parameters	Number of statistically significant parameters	Adjusted R <sup>2</sup>
Single AAs	20	16	.498
Single AAs - reduced	16		.498
Dipeptides	400	69	.651
Dipeptides - reduced	69		.608
Published Motifs	7	5	.334
<i>de novo</i> motifs	40	27	.502
<i>de novo</i> motifs - reduced	27		.500
<i>de novo</i> motifs + single AAs	60	37	.606
<i>de novo</i> motifs + single AAs	37		.604

**Table 1:** Ordinary Least Squares regression on tile composition explains a large fraction of the variance in measured activation domain activity

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

### The WxxLF motif requires acidic context and supporting hydrophobic residues.

The absence of clear motifs raises the question of how the arrangement of amino acids, the sequence grammar, controls activation domain function. As an anchor point, we used the WxxLF motif, which makes large contributions to activity in the CAD but not all tiles with this motif are active (**Figure 2B, S8, S20**). We compared tiles with the WxxLF motif that had high or low activity: highly active tiles were more acidic and had more WFYLM residues (**Figure S21C,D**). The first grammar signal we found is that tiles with more evenly intermixed acidic and W,F,Y,L residues are more active, supporting the acidic exposure model (**Figure S21E**). The strongest grammar signal is that tiles with the WxxLF motif near the C-terminus are active (**Figure 4D, S22**). The additional negative charge of the C-terminus may increase exposure of the motif. Weak C-terminal effects have been seen for aromatic residues<sup>20,37</sup>. This result

327 emphasizes how even a conserved short linear motif requires an acidic context  
 328 and supporting hydrophobic residues to create an activation domain. Together,  
 329 our analysis indicates that yeast activation domains are nucleated by a cluster of  
 330 aromatic residues surrounded by acidic residues and supported by leucine and  
 331 methionine residues.

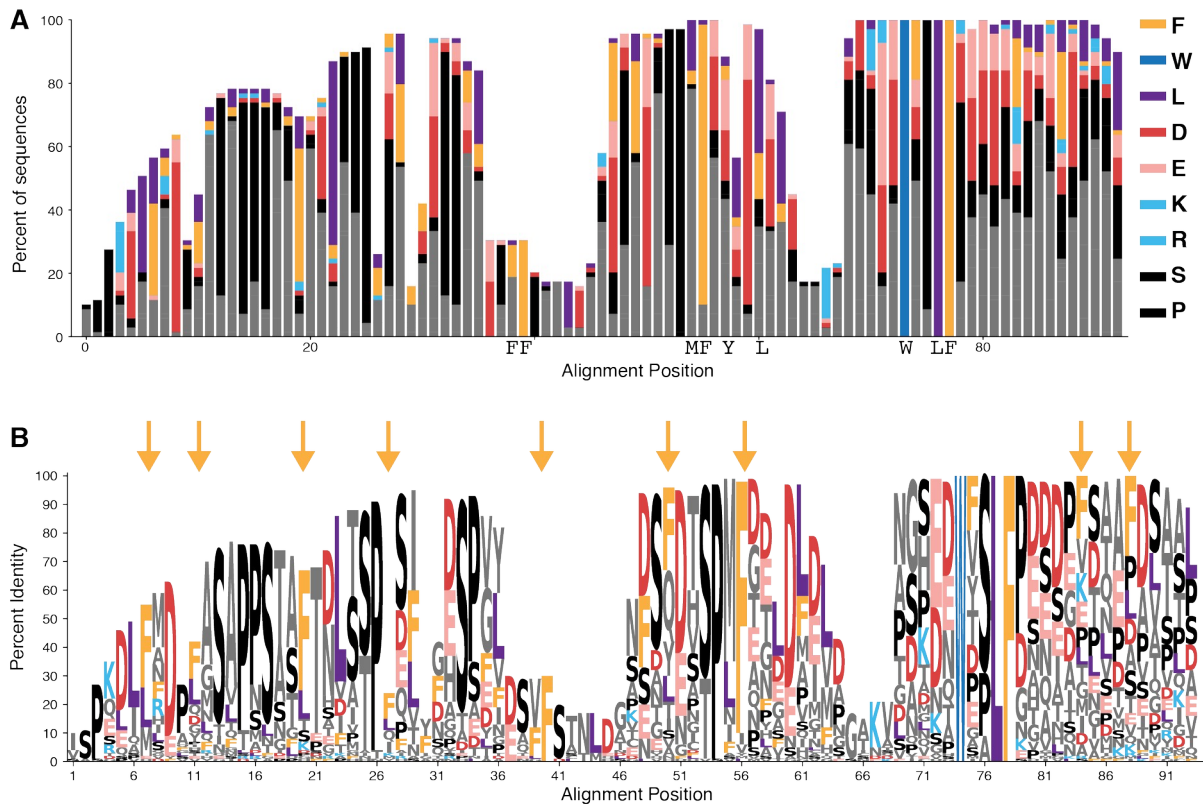
### 332 **The alpha helix from *S. cerevisiae* is dispensable for full activity.**

333 The sequence diversity of strongly active tiles with the WxxLF motif  
 334 strongly suggests that coupled folding and binding is not necessary for activity.  
 335 In *S. cerevisiae*, the disordered CAD folds into a short alpha helix upon binding  
 336 the Gal11/Med15 coactivator<sup>45,46</sup>. Inserting a proline into this helix has little effect  
 337 on activity<sup>16,45</sup>. The immediate vicinity of the helix in *S. cerevisiae* has 115 unique  
 338 sequences: 23 contain 3 prolines (20%), and 3 contain 4 prolines, e.g.  
 339 **GPSDPWYPLFPSDTA**. Using a 70 residue region, we predicted alpha helix  
 340 propensity, but only 38/138 (28%) are predicted to form a helix (**Figure S23**,  
 341 methods<sup>47</sup>). These sequences may still fold into a helix when binding to the  
 342 cognate partner. Amphipathic helices are enriched in activation domains<sup>18,32</sup>  
 343 because they are a convenient way to present hydrophobic residues to a partner,  
 344 but they are not the only way to create a strong activation domain. CAD function  
 345 is more conserved than alpha helix formation. This analysis suggests the alpha  
 346 helix is not the relevant functional unit for evolutionary turnover.

347

348

349



350

351 **Figure 5:** Evolutionary turnover of aromatic and acidic residues explains the

352 conservation of function without conservation of sequence in the central

353 activation domain of Gcn4. **A)** For the 69 most active unique regions around the  
 354 WxxLF motif, a bar plot showing the relative amino acid frequencies from the  
 355 MSA. MSA positions with >90% gaps have been removed. The acidic residues, D  
 356 and E, interchange. **B)** A sequence logo for the MSA. Arrows indicate the 9  
 357 positions where F is the most abundant residue. There is some interchange  
 358 between F and L. Black, SP motifs. See **Figure S25** for the MSA.

359

### 360 **Conservation of function without conservation of sequence in the** 361 **Central Acidic Activation Domain of Gcn4.**

362 The central activation domain region of the Gcn4 orthologs showed strong  
 363 conservation of function without conservation of sequence. Our two hypotheses  
 364 for this phenomenon were evolutionary turnover of motifs or evolutionary  
 365 turnover of key residues.

366 We found no evidence for turnover of motifs. Each of the published motifs  
 367 contributed to activity (**Figure 2B**) and was enriched in active tiles (**Figure**  
 368 **S21A**), but only the WxxLF motif was conserved (**Figure 2C, S3**). We did not  
 369 detect the emergence of new instances of these motifs, so we can reject the  
 370 motif turnover hypothesis.

371 The most conserved sequence feature of the CAD region besides the  
 372 WxxLF motif is an SP motif, which is not typically associated with activation  
 373 domain function. The full-length orthologs contain up to 4 SP motifs upstream of  
 374 the WxxLF motif. In *S. cerevisiae*, this SP is a TP (T105), which is phosphorylated  
 375 to create a phosphodegron that shuts down the Gcn4 program during the  
 376 recovery from starvation<sup>48,49</sup>. The majority of tiles (10752) contain an SP motif, so  
 377 it makes little contribution to activity on its own. We believe these motifs are  
 378 conserved due to regulated degradation. However, it remains possible that  
 379 multisite phosphorylation can increase activation domain activity<sup>50,51</sup>. For 32 tiles  
 380 we performed a followup mutagenesis of the SP motifs to test the hypothesis  
 381 that phosphorylation can control activity (**Figure S24**). These data support the  
 382 possibility that some of the orthologs utilize phosphorylation to modulate  
 383 activation domain function.

384 We observe evolutionary turnover of acidic and F residues within the  
 385 central activation domain. We focused on a 70 residue region around the WxxLF  
 386 motif (W-50 : W+19) that contained many active tiles and the peak of inferred  
 387 activity (**Figure S25**). The top half of sequences contain many acidic residues,  
 388 but individual acidic positions (D and E) are not well conserved because they  
 389 interconvert (**Figure 5A,B, S26**). When pooled together, D+E conservation  
 390 matches or exceeds the conservation level of the aromatic residues. In addition,  
 391 the F residues that are critical for high activity exhibit evolutionary turnover.  
 392 There are only 2 positions where an F is present in the majority of sequences,  
 393 but an additional 7 positions where F is the most common residue. The critical F  
 394 residues experience evolutionary turnover, giving the appearance of moving  
 395 around the activation domain.

396 To this point, all of our analysis has used only the MSA, so we next  
 397 leveraged the additional information present in the species tree. We tested the  
 398 hypothesis that the gains of F residues precede the loss of F residues. In most  
 399 cases, there is too much evolutionary distance between the species to answer

400 this question. However, in the high quality YGOB alignment<sup>41</sup>, we see the gain of  
 401 an F precedes the loss of an ancestral F (**Figure S27**). This example of gains  
 402 preceding loss bolster then evidence for evolutionary turnover of key residues.

403 The turnover and conservation patterns we observed in the Gcn4  
 404 orthologs generalized to other systems. We reanalyzed a set of orthologs of Pdr1  
 405 (**Figure S28**)<sup>18</sup>. For four TFs, we searched for orthologs in the Y1000+ collection  
 406 and made alignments of their activation domains (**Figure S28**). In these MSAs,  
 407 aromatic residues were highly conserved and acidic residues interchange at  
 408 many positions. Some positions also showed interchange between aromatic  
 409 residues. Other regions showed turnover of aromatic and leucine residues. We  
 410 propose that evolutionary turnover of key aromatic, leucine, and acidic residues  
 411 is a general feature of eukaryotic acidic activation domains.

## 412 **Machine learning insights into activation domain function**

413 The Gcn4 orthologs provide a large, unique dataset to evaluate deep  
 414 learning models that predict activation domains from amino acid sequence. We  
 415 compared two first-generation neural networks<sup>18,21</sup> with a second-generation  
 416 model that we developed<sup>23,24</sup>. All the models can approximate the locations of  
 417 activation domains in full-length TFs, but the new model, TADA, is substantially  
 418 more accurate at predicting the activities of individual tiles and identifying  
 419 activation domain boundaries (**Figure S29**). TADA was intentionally built to  
 420 ignore sequence grammar by blurring the raw sequence with sliding windows  
 421 and its high performance supporting the idea that there is very weak or very  
 422 little grammar in these orthologs. The machine learning models cannot detect  
 423 'missing' grammar, supporting the weak grammar hypothesis. The high accuracy  
 424 of these models suggests they may be ready to enable evolutionary studies.

425 We used TADA to predict the contributions of F residues to activity in the  
 426 central activation domain. The model predicts that all the F residues contribute  
 427 to activity (**Figure S30**). The contributions of the most conserved F positions are  
 428 indistinguishable from recently evolved F positions (**Figure S30D**). This analysis  
 429 further supports evolutionary turnover of key F residues.

430

## 431 **Discussion:**

432 By functionally screening protein fragments from a family of orthologous  
 433 sequences, we demonstrate how activation domains show strong conservation of  
 434 function without conservation of sequence through turnover of critical acidic and  
 435 phenylalanine residues. Conservation of function without conservation of  
 436 sequence was established for full-length TFs, but here we demonstrate how this  
 437 phenomenon emerges from turnover of entire activation domains and turnover  
 438 of key residues within activation domains. Our results emphasize how IDR  
 439 function can be highly conserved and constrained yet invisible in traditional  
 440 comparative genomics.

441 The observed turnover of critical residues supports our acidic exposure  
 442 model for activation domain function and explains why it is so difficult to identify  
 443 motifs in activation domains. Multiple screens for activation domains have found  
 444 only one recurrent motif, LxxLL motif with an acidic context, which can be  
 445 important for binding the Kix domain<sup>17,18,21,31,52,53</sup>. These screens have also shown

446 that the 9aaTAD is not enriched in active sequences<sup>54</sup>. We argue that activation  
447 domains are nucleated by Clusters of W and F residues surrounded by acidic  
448 residues and boosted by Y, L, and M residues. Under this weak molecular  
449 grammar, individual residues are easily replaced, facilitating turnover. The  
450 WxxLF motif is one solution among many. When only a few sequences are  
451 examined, clusters look like motifs. Each TF family has a different conserved  
452 cluster of hydrophobic residues that represents a very good solution to binding  
453 the preferred coactivator. Each TF family will appear to have a conserved,  
454 essential motif, but convergent evolution of motifs is rare (Supplemental note 1).

455 We propose that the physical flexibility of the protein interaction interface  
456 between Gcn4 and Med15 allows for evolutionary plasticity. The Gcn4 CAD  
457 undergoes coupled folding and binding with the Med15 activation domain  
458 binding domains, but this interaction is a physically flexible, fuzzy  
459 interaction<sup>44,45,55</sup>. The short helix presents the WxxLF motif in many orientations  
460 to a shallow hydrophobic canyon on Med15. Molecular dynamics simulations  
461 suggest that these orientations interconvert<sup>46</sup>. This binding interaction imposes  
462 few structural constraints on the Gcn4 CAD.

463 The turnover of hydrophobic residues is possible because of this physical  
464 flexibility of the Gcn4-Med15 protein-protein interaction. The weak structural  
465 constraint of this interaction enables evolutionary plasticity. Binding one  
466 sequence in multiple orientations is a step towards binding diverse orthologs,  
467 which in turn is a step towards binding to many activation domains<sup>18,56,57</sup>. This  
468 flexibility likely requires at least one disordered partner<sup>58</sup>. Coactivators that  
469 impose weak structural constraints on activation domains can become engines  
470 for evolutionary diversification of activation domains through neutral drift,  
471 creating an enormous sequence reservoir for later selection. Although we favor  
472 the hypothesis that the observed sequence divergence in Gcn4 orthologs is  
473 neutral, stabilizing selection, it remains possible that there is selection to  
474 diversify. Acidic activation domains are highly evolutionarily successful,  
475 representing more than half of all known examples<sup>27</sup>. Our observation that acidic  
476 activation domains can easily diversify without compromising function suggests  
477 they are highly evolvable. This evolvability creates a diverse sequence reservoir  
478 that allows for rapid selection on standing variation. We speculate this  
479 evolvability allowed for acidic activation domains to bind new coactivators as  
480 they emerged with multicellularity<sup>59</sup>.

481 Activation domain evolution exemplifies how protein-protein interactions  
482 mediated by IDRs can drive evolutionary plasticity and sequence diversity.  
483 Another example of an IDR engaged in flexible PPIs enabling evolutionary  
484 plasticity is the human TRIM5 antiviral caging system, wherein short disordered  
485 loops make multivalent contacts with the viral capsid<sup>60</sup>. Physically flexible  
486 binding and avidity provide the emergent specificity to keep up in evolutionary  
487 arms races with fast-evolving viruses<sup>61</sup>.

488 Our results fit well with findings that at long evolutionary distances,  
489 transcriptional regulatory networks rewire, substituting individual TFs but  
490 maintaining circuit logic<sup>39,62,63</sup>. Here, we examined longer evolutionary distances  
491 and found that all the Gcn4 orthologs are activators. This consistency of TF  
492 function shows that the sign of TF connections in regulatory networks are more

493 conserved than individual connections. Changes in TF function are pleiotropic,  
494 affecting many targets. Slow or rare changes in TF function likely make it easier  
495 to substitute TFs at individual regulatory elements.

496 Our deep dive into the evolution of one IDR family complements other  
497 studies of IDR evolution. Using small numbers of sequences, conservation of IDR  
498 function across orthologs has been observed, but often the essential residues are  
499 unknown<sup>10</sup>. In other systems, there is functional conservation of diverged IDRs,  
500 but the key residues are conserved<sup>12</sup> or motifs are conserved<sup>64</sup>. In other cases,  
501 functional conservation results from the composition, but not the arrangement,  
502 of residues through emergent properties like net charge<sup>11,65-69</sup>. The closest  
503 parallel to our turnover of key residues is *de novo* evolution of phosphorylation  
504 motifs<sup>70</sup>. TF IDRs are not always functionally conserved, for example in Abf1<sup>13</sup>  
505 and the Msn2/4 IDRs have two overlapping functions, only one of which is  
506 conserved<sup>14</sup>. Sox family members from Chianoflagelites can substitute for Sox2  
507 in mouse iPSC reprogramming<sup>6</sup>. Cases where function emerges from physical  
508 properties may allow for even more turnover than we observe in Gcn4. There  
509 remains a need for better IDR-alignment algorithms or alignment-free methods  
510 to group functionally related IDRs.

511 The turnover of key hydrophobic residues in activation domain evolution  
512 bears strong parallels to the turnover of TF binding sites in enhancer evolution.  
513 In metazoans, enhancers are regulatory DNA that contain clusters of TF binding  
514 sites (TFBS). The DNA sequence of enhancers diverges rapidly as individual TFBS  
515 are gained and lost while maintaining function<sup>71-73</sup>. Orthologous enhancers can  
516 be impossible to detect in sequence alignments but are readily identified by  
517 searching for clusters of TFBS<sup>74,75</sup>. Two mechanistic insights led to this predictive  
518 power: 1) understanding that the key functional subunit is the TFBS and 2)  
519 understanding that individual TFBS can turnover. This conservation of total  
520 binding site content enables complex of regulatory DNA to identify conserved  
521 enhancers<sup>74,76</sup>. We find strong parallels in the evolution of TF protein sequence.  
522 TF protein sequence changes rapidly and is hard to align, but activation domain  
523 function is conserved. Analogous to the TFBS in enhancers, the functional units  
524 of activation domains are individual aromatic residues. In both cases, the  
525 grammar is extremely flexible<sup>77</sup>. Given that TFs function by binding to enhancers,  
526 it is striking that both the protein and the DNA are evolving in the same way.  
527 Turnover of TF binding sites endows enhancers with robustness to genetic  
528 variation, robustness to environmental stress, and evolutionary plasticity.  
529 Turnover of key residues in activation domains may similarly endow TFs with  
530 plasticity and robustness. If TFs and enhancers are evolving in the same way, it  
531 increases the potential for compensatory mutations, expanding the neutral  
532 space and creating diverse sequence reservoirs that can be selected in new  
533 environments.

534 The primary limitation of this work is that we measured the activities of  
535 short fragments in one species. Measuring short uniform fragments makes the  
536 experiments possible but can miss longer 'emergent' activation domains<sup>55,78</sup>. If,  
537 in some species, an activation domain and cognate coactivator together  
538 experience many compensatory mutations, the assay will miss these sequences.  
539 Our analysis of Med15 coactivator conservation shows that the four activation



540 domain binding domains are conserved (**Figure S31**). Activity of our reporter is  
541 well correlated with Med15 binding affinity *in vitro*<sup>18</sup>. The most active tiles are  
542 computationally predicted to bind Med15<sup>79</sup> (**Figure S32**). In the future, limited  
543 screening in additional species or screening tiles of multiple tile lengths would  
544 enrich this work. A secondary limitation is that we measured activity in just one  
545 condition. A future direction is to explore activity in other conditions and on other  
546 promoters.

## 547 Materials and Methods

### 548 Identification of ortholog sequences

549 We computationally screened for Gcn4 orthologs of *S. cerevisiae*. We  
 550 started with a hand-collected set of 49 orthologs, 48 of which contained the  
 551 WxxLF motif<sup>16,55</sup>. To find new orthologs, we used two criteria: the bZIP DNA  
 552 binding domain (IPR004827) and the regular expression Wx[SPA]LF for the  
 553 WxxLF motif. These criteria distinguished Gcn4 orthologs from other leucine  
 554 zipper DNA binding domain TFs. We scanned 207 diverse and representative  
 555 proteomes from the MycoCosm database (mycocosm.jgi.doe.gov). This  
 556 computational screen yielded 1188 gene models from 129 genomes. These 1188  
 557 gene models combine to yield 502 unique proteins (**Table S1, Figure S1, S2**).  
 558 Of these, >99% were reciprocal Blast best hits with *S. cerevisiae* Gcn4. This initial  
 559 analysis was performed in 2020 by Sumanth Mutte of MyGen Informatics. 84 of  
 560 the genomes were from MycoCosm, while the original ortholog collection  
 561 contributed 45 species. Genomes contained 1-32 gene models and 1-11 unique  
 562 protein sequences (**Figure S1**). These sequences span nearly all the  
 563 Ascomycota, the largest phylum of Fungi, representing >600 million years of  
 564 evolution<sup>42</sup>. The 502 unique orthologs have variable lengths (**Figure 1A**), but  
 565 the DBD is at the C-terminus in 500 orthologs, and the distance between the  
 566 WxxLF motif and the DBD is very consistent (**Figure 1B**).

567 All species were from the Ascomycota except for five entries with three  
 568 unique sequences from Blastocladiomycota (**Figure S1**). The Blastocladiomycota  
 569 orthologs are the only proteins where the WxxLF motif does not align in the MSA.  
 570 The sequence context of their WxxLF motif is H-rich instead of acidic:

571 e.g. AAAQHVPAAADGQWLALFPHPSIDFDNFNSFHQSFSSPPPH

572 The Blastocladiomycota tiles with the WxxLF motif have high activity in  
 573 the assay. The regions of Blastocladiomycota orthologs that align to the WxxLF  
 574 motif in the MSA have low activity in the assay. We suspect the N-terminal  
 575 WxxLF in the Blastocladiomycota may have been gained by convergent evolution  
 576 (Supplemental note 1).

577 The Yeast Gene Order Browser has reconstructed the local synteny of the  
 578 Gcn4 locus for 37 genomes yielding a high-quality set of true homologs<sup>41</sup>. 36/37  
 579 species and the inferred ancestor contain one Gcn4 gene. *Kazachstania*  
 580 *saulgeensis* CLIB1764T is missing a Gcn4 homolog. All of the post whole genome  
 581 duplication species in this set contain only one Gcn4 homolog, suggesting there  
 582 is no advantage of retaining two copies. All but one of the 36 the orthologs,  
 583 *Zygosaccharomyces baili*ii ZYBA0L03268g, contain the WxxLF motif. Instead, *Z.*  
 584 *baili*ii has an insertion in the WxxLF motif yielding WPSLEPLF. This sequence was  
 585 not included in our experiment but was previously measured in a 44 AA tile,  
 586 LDQAVVDEFFVND DAPMFELDDGASGAWPSLEPLFGEDEERVAV, and had high activity in Replicate 2  
 587 of our previous paper<sup>16</sup>. This example further supports the observed  
 588 conservation of function without conservation of sequence.

589 Despite substantial sequence divergence, all homologs show negative  
 590 selection at the level of the full protein in the precomputed YGOB analysis. We  
 591 downloaded a list of 36 pairwise Ka, Ks, and omega coefficients calculated from  
 592 the yn00 output of Phylogenetic Analysis by Maximum Likelihood (PAML) (**Table**  
 593 **S14**, November 2024).

594 We confirmed that the WxxLF motif is well conserved in fungal TFs with  
 595 HMMER. We ran the web server for HMMER with default parameters, using  
 596 *S. cerevisiae* Gcn4 as the seed sequence and restricting our search to Fungi. In  
 597 the second, third, and fourth iterations of this search, the WxxLF motif was the  
 598 most prominent feature of the profile HMM in the central region of the TF and

599 always much more prominent than all other published motifs <sup>21,78</sup>. **Figure 2C**  
600 shows the pHMM from the fourth iteration.

601 For the full-length orthologs, MSAs were performed in Genious with the  
602 MAFFT algorithm (**Table S2**). We removed the two longest orthologs that had  
603 the DBD near the center. In the MSA, 54% of positions had less than 1% identity  
604 and 88% had less than 5% identity.

605 Short alignments were created with MUSCLE online (<https://www.ebi.ac.uk/Tools/msa/muscle/>) or with or with MAFFT v7.526 and visualized with  
606 weblogo.berkeley.edu or the LogoMaker Python package.  
607

## 608 **Design of the Gcn4 oligo library**

609 We took the 502 unique protein sequences and computationally chopped  
610 them into 40 AA tiles spaced every 5 AA (e.g. 1-40, 6-45, 11-50 etc.). As a result,  
611 if two closely related sequences contain identical regions, insertions or  
612 alternatives (start sites) that change the phasing, a single tile can map to  
613 multiple full-length orthologs. We removed duplicate tile sequences, yielding  
614 20679 unique tiles. We added 52 control sequences (controls were included  
615 twice in the oligo pool to increase the probability they were recovered in the  
616 plasmid pool during cloning). The controls included hand-designed mutants in  
617 control activation domains and a handful of sequences from our previous study <sup>16</sup>  
618 (**Table S3**, Control sequences). The final design file contained 20783 entries.

619 We reverse-translated tile sequences using *S. cerevisiae* preferred codons.  
620 We added primer sequences for PCR amplification and HiFi cloning ('ArrayDNA'  
621 column in **Table S5**). We also added four Stop codons in three reading frames to  
622 ensure translational termination, even if there were one or two bp deletions, the  
623 most common synthesis errors. We used synonymous mutations to remove  
624 instances where the same base occurred four or more times in a row to reduce  
625 DNA synthesis errors. The resulting oligo pool was ordered from Agilent  
626 Technologies. The final oligos were of the form (see primer sequences in **Table**  
627 **S4**):

628 `FullDNAseq = primer1 + ActivationDomainDNAseq + stopCodons + primer2`

## 629 **Plasmid Library construction**

630 The oligos were resuspended in 100 uL of water, yielding a 1 pM solution.  
631 The oligos were amplified with eight reactions of Q5 polymerase (NEB) using 1 ul  
632 of template, five cycles, T<sub>m</sub> = 72C and the LC3.P1 and LC3.P2 primers. The eight  
633 reactions were combined into a single PCR clean-up column (NEB Monarch).

634 The backbone was prepared by digesting 16 ug of pMVS219 with NheI-HF,  
635 PacI and AclI in eight reactions. We digested for seventeen hours at 37C and  
636 heat-inactivated for one hour at 80C. The desired 7025 bp fragment was run on a  
637 0.8% gel, visualized with SYBR Safe (Invitrogen), and gel purified (NEB Monarch  
638 Kit). Note pMVS219 and pMVS142 have the same sequence, but the pMVS142  
639 stock developed heteroplasmy, so we repurified it as pMVS219 and submitted  
640 the corrected stock to AddGene. Both pMVS219 and pMVS142 correspond to  
641 AddGene #99049.

642 We used NEB HiFi 2x mastermix to perform Gibson Isothermal Assembly to  
643 create the plasmid library. The 4x reaction volume had 328 ng of backbone and  
644 excess molar insert. We incubated at 50C for 15 min and assembled a backbone-  
645 only control in parallel. The assemblies were electroporated three times each  
646 into ElectroMax 10b E.coli (Invitrogen 18290-015) following the manufacturer's  
647 protocol. A dilution series was plated and the bulk of the cells grown overnight in  
648 140mL LB+Amp. These cultures overgrew, so they were spundown and frozen.  
649 The cultures were regrown with 105 mL LB+Amp and a MaxiPrep was performed

650 (Zymo). An estimated 4.2 million colonies were collected, covering the library  
651 200-fold.

652 To assess the quality of the plasmid library, we prepared an amplicon  
653 sequencing library (see below). Three independent amplicon libraries were  
654 prepared, and sequences present in all three were considered to be present in  
655 the plasmid pool with high confidence. GREP for the flanking NheI and Ascl sites  
656 was used to pull out the designed fragments. Only perfect matches were used in  
657 this analysis. 20717 of 20731 designed sequences were detected (99.9%). The  
658 vast majority sequence abundances were within 4-fold of each other, indicating  
659 minimal skew in library member abundance.

## 660 **Yeast transformation**

661 The plasmid library was integrated into the DHY213 BY superhost strain,  
662 MATa his1Δ1 leu2Δ0 ura3Δ0 met15Δ0 MKT1(30G) RMEI(INS-308A) TAO3(1493Q),  
663 CAT5(91M), MIP(661T) SAL1+ HAP1+, a generous gift from Angela Chu and Joe  
664 Horecka. Requests for the parent strain are best directed to them. We integrated  
665 our library into the URA3 locus with a three-piece PCR<sup>80</sup>. The upstream  
666 homology between URA3 and the ACT1 promoter was created by PCR amplifying  
667 the pMVS295 (Strader 6161) with the primers YP18 and CP19.P6. The  
668 downstream homology between the TEF terminator of KANMX and URA3 was  
669 amplified from pMVS196 (Strader 6768) with the primers YP7 and YP19. These  
670 template plasmids were a generous gift from Nick Morffy and Lucia Strader. To  
671 avoid PCR, the plasmid library was digested with Sal I-HF and EcoRI-HF (NEB)  
672 overnight, but not cleaned up. The homology arms were in 3:1 molar excess.  
673 1.25 ug of total DNA was used (225 ng of upstream homology 626 bp, 225 ng of  
674 downstream homology 665 bp, and 800 ng of digested plasmid 4583 bp). Cells  
675 were streaked out from the -80C on YP+Glycerol. Four transformation cells were  
676 grown overnight in YPD, diluted into YPD, and allowed to grow for at least two  
677 doublings. We performed a Lithium Acetate transformation with 30 minutes at 30  
678 C and 60 minutes at 42 C followed by a two hour recovery in synthetic dextrose  
679 minimal media without a nitrogen source, as recommended by Sasha Levy. We  
680 integrated plasmids in seven transformation batches, which were plated  
681 overnight on YPD and replica-plated onto YPD+G418 (200 ug/ml). Plates were  
682 stored at 4 C and then scraped with water, pooled, frozen into glycerol stocks,  
683 and mated. We collected an estimated 100,000 colonies, approximately five-fold  
684 coverage of the tiles. For 6/7 pools we sequenced tiles before and after mating,  
685 finding that 67-97% of tiles were detected both before and after mating,  
686 indicating that the mating sometimes reduced library complexity.

## 687 **Yeast Mating**

688 We mated each of the seven transformations independently to MY435  
689 (FY5, MATalpha, YBR032w::P3 GFP ClonNat-R (pMVS102)). Downstream  
690 sequencing revealed that transformations with modest numbers of colonies (e.g.  
691 4500) experienced no significant loss of complexity during mating, but  
692 transformations with more colonies (e.g. >20,000) experienced loss of  
693 complexity, up to 40% in one case. Subsequent matings were performed in  
694 larger volumes to avoid creating a bottleneck. Mated diploids were selected in  
695 liquid culture with YPD with 200 ug/ml G418 and 100 ug/ml ClonNat. After  
696 overnight selection, matings were concentrated and frozen as glycerol stocks.

## 697 **Cell Sorting**

698 The day before sorting, a glycerol stock of mated cells (~100 ul) was  
699 thawed into 5 mL SC+Glucose with 200 ug/ml G418 and 100 ug/ml ClonNat and  
700 grown overnight, shaking at 30 C. The morning, the culture was diluted 1:5 into

701 SC+Glucose with G418, ClonNat, and 10  $\mu$ M  $\beta$ -estradiol (Sigma). The culture was  
702 grown for 3.5-4 hours before sorting.

703 Cells were sorted on a BD Aria Fusion equipped with four Lasers (488 blue,  
704 405 Violet, 561 Yellow-green and 640 Red) and eleven fluorescent detectors. We  
705 used two physical characteristics gates, first to enrich for live cells (FSC vs SSC)  
706 and second to enrich for single cells (FSC-Height vs FSC-Area). Cells were sorted  
707 by the GFP signal, the mCherry signal, or the ratio of GFP:mCherry signal. The  
708 ratio is a synthetic parameter that is very easy to saturate on the eighteen-bit  
709 scale available in the BD software. Great care was taken to change PMT voltage  
710 and the ratio scaling factor (5-10% depending on the day) to make the value of  
711 the top and bottom bins as different as possible. The dynamic range of our final  
712 estimate for activation domain activity is set by the value of the top and bottom  
713 bins. The maximum activation domain strength is 100% in the top bin, and  
714 assumes the value of the top bin. The minimum activation domain strength is  
715 100% in the bottom bin and assumes the value of the bottom bin.

716 We performed our sorting experiment twice. In the first run, we pooled all  
717 of the transformants into one sample and sorted it by GFP/mCherry ratio, GFP-  
718 only, mCherry-only. We sorted one million cells per bin. For the ratio sort, we  
719 split the ratio histogram in eight approximately equal bins <sup>16</sup>.

720 In the second round of sorting, we split the transformants into two pools,  
721 labeled A and B, so we could assess measurement reproducibility for  
722 independent transformants. Pool A and Pool B are true biological replicates. We  
723 sorted each pool by GFP/mCherry ratio, GFP-only, mCherry-only. We used the  
724 comparison of the A and B pool measurements to assess measurement  
725 reproducibility of true biological replicates. We have never previously measured  
726 this biological reproducibility. On this day, we sorted 250000 cells per bin.

727 Sorted cells were grown overnight in SC-glucose. The next morning, gDNA  
728 was extracted with the Zymo YeaSTAR D2002 kit, using Protocol I with  
729 chloroform according to the manufacturer instructions. We have previously  
730 shown that growing cells overnight makes the gDNA extraction easier but does  
731 not change the computed activation domain activity <sup>16</sup>.

### 732 **Amplicon Sequencing Library preparation**

733 Amplicon sequencing libraries were prepared from genomic DNA in three  
734 steps. First, the general vicinity of the tile sequence was amplified with CP21.P14  
735 and CP17.P12 using 100 ng of gDNA as template and yielding a 604 bp product  
736 that was cleaned up (Monarch PCR cleanup). In the second PCR, we added 1-4 bp  
737 of phasing on each end and the Illumina sequencing primer in 7-10 cycles with  
738 SL5.F[1-4] and SL5.R[1-3]. These seven phased primers were pooled and added  
739 to all samples. Four nanograms of the first PCR were used as template for the  
740 second PCR. Two microliters of the second PCR served as template for the third  
741 PCR. The third PCR added unique Index1 and Index2 sequences to each sample  
742 with an additional 7-10 cycles. These final products were cleaned up with PCR  
743 columns or magnetic beads (MacroLab at UC Berkeley) and submitted for  
744 sequencing. We performed 2x150 bp paired end sequencing in a shared Nova-  
745 Seq lane at the Washington University School of Medicine Genome Technology  
746 Access Center (GTAC). GTAC provided demultiplexed fastq files. We sequenced  
747 additional samples in shared Nova-seq lanes with MedGenome.

### 748 **Sequencing Analysis**

749 After demultiplexing samples and pairing reads with PEAR, we kept only  
750 the reads where the tile DNA sequence contained a perfect match to a designed  
751 tile. For each eight bin sort, we performed two normalizations. We first  
752 normalized the reads by the total number of reads in each bin. Then, we

753 normalized across the eight bins to calculate a relative abundance. We then  
 754 converted relative abundances to an activity score for each tile by taking the dot  
 755 product of the relative abundance with the median fluorescence value of each  
 756 bin (**Table S8**). This weighted average is the measured activation domain  
 757 activity. Tiles with fewer than forty-one reads were not included in the final  
 758 dataset. These analysis scripts are available at  
 759 [github.com/staller-lab/labtools/tree/main/src/labtools/adtools](https://github.com/staller-lab/labtools/tree/main/src/labtools/adtools). This preprocessing  
 760 computed an activity for each tile in each experiment. Activity is uncorrelated  
 761 with total reads (**Figure S5E**). The pooled ratio sort (BSY2) had 115.6 M reads.  
 762 The Replicate A ratio sort had 934.5 M reads, and the Replicate B ratio sort had  
 763 697 M reads. Replicate A GFP had 33.1 M reads, Replicate B GFP had 31.6 M  
 764 reads, Replicate A mCherry had 32.8 M reads, and Replicate B mCherry had 30.3  
 765 M reads.

## 766 **Measurement Reproducibility**

767 We used the two measurements of independent transformants to assess  
 768 the reproducibility of our measurements of true biological replicates ( $R = .870$ ;  
 769 **Figure S5A-D**). Reproducibility is higher ( $R = .919$ ) for highly abundant tiles  
 770 ( $>1000$  reads).

771 We combined data from the two biological replicates. For tiles present in  
 772 both populations ( $n = 11797$ ), we averaged the two measurements and used the  
 773 standard deviation as the error bar. For tiles present in only one population, we  
 774 used that measurement and did not report error bars. These combined data  
 775 agree very well with the pooled sort ( $R = .919$ ; **Figure S5C**). Activity was  
 776 saturated for forty-nine tiles, but most of these were measured with low fidelity  
 777 because they had low read depth, and forty-seven were present in only one  
 778 biological replicate. We identified forty-one tiles that were very highly active in  
 779 both replicates and had high read depth in both replicates (**Table S11**). These  
 780 we recommend for CRISPR Activation studies in yeast.

781 We assessed whether the mating introduced biological variability. We  
 782 remated seven pools of the integrated library to the same reporter line, selected  
 783 for diploids, pooled them, and resorted cells. This time we sorted 500,000 cells  
 784 per bin. This measurement agreed with the initial experiments ( $R = 0.920$ ;  
 785 **Figure S5D**).

786 Inferred activity was not correlated with read count, which, as previously  
 787 shown, is another indicator of high-quality data (**Figure S5E**).

788 We compared activity measurements to our previously published results  
 789 <sup>16</sup>. Previously, we used forty-four AA regions, and here we used forty AA tiles. We  
 790 considered any forty-four AA tile that contained one or our forty AA tiles to be  
 791 corresponding pairs. The extra four AA can modify activity, so the  
 792 correspondence of these measurements will not be perfect. The observed  
 793 Pearson correlation of 0.786 and Spearman correlation of 0.731 indicate the new  
 794 data are of high quality and consistent with previous measurements (**Figure**  
 795 **S5F**).

796 The technical reproducibility of our measurements at UC Berkeley are  
 797 lower than the published reproducibility from sorting at Washington University in  
 798 St. Louis <sup>16</sup>. In both cases, we sorted the same cell population twice and created  
 799 independent sequencing libraries. In 2018, the technical reproducibility was high,  
 800 Pearson  $R = 0.988$ . The 2018 work had a smaller library ( $<5000$  unique  
 801 sequences) and sorted more cells (1-2 million cells per bin). Sorting more cells  
 802 per library member increases the technical reproducibility of the measurement.  
 803 The sorter operator in the 2018 work was more experienced than the sorter  
 804 operator in this work (MVS), and the machine was maintained to a higher  
 805 standard of operation, so the sorted populations were purer.

806 The eight bin ratio activity measurements are primarily driven by the GFP  
 807 signal. Activity (ratio) is largely separable from abundance assessed by the  
 808 mCherry sort (**Figure S5G-I**) and well-correlated with the GFP sort (**Figure S5J-**  
 809 **L**).

## 810 **Determining a threshold for active tiles**

811 The full distribution of tile activities has a peak at low activity, which,  
 812 based on control sequences, is clearly inactive, with a heavy right shoulder and a  
 813 heavy right tail (**Figure 1F**). The tail contains the control sequences with known  
 814 high activity (**Figure S7**). We set out to fit the inactive sequences to a Gaussian  
 815 distribution and use this distribution to create a threshold for active sequences.  
 816 We first bin all tiles according to their activity score such that there are ~ 200  
 817 tiles per bin and plot a histogram. We hypothesized tile density is highest around  
 818 inactive tiles and thus refer to all tiles to the left of the resulting histogram's  
 819 peak as inactive tiles. We fit a one-sided Gaussian to these inactive tiles (**Figure**  
 820 **S8A**) and call the two-sided extension of this Gaussian the inactive tile  
 821 distribution (**Figure S8B**). Treating this Gaussian inactive tile distribution as our  
 822 null hypothesis, we calculate p-values for each tile (not including tiles earlier  
 823 used as inactive, **Figure S8D**). We then correct for multiple comparisons using  
 824 FDR<sup>81</sup> and Bonferroni<sup>82</sup> corrections. The 1% FDR threshold was 33821 (60.6% of  
 825 tiles active). The 1% FWER threshold was 45373 (46.6% of tiles active). As a  
 826 conservative threshold to call active sequences, we used the 1% FWER threshold  
 827 of 45,373. All of our designed inactive control sequences are below this  
 828 threshold.

829 After trying many thresholds (**Figure S8**), we ultimately chose the top  
 830 20% (94,031) as a threshold for high activity. The choice of threshold had very  
 831 little effect on our results. In particular, a wide range of threshold has almost no  
 832 effect on the number of orthologs with an active tile.

## 833 **Protein sequence parameters**

834 We computed protein sequence parameters (Net charge, local net charge,  
 835 Kyte Doolittle Hydrophobicity, Wimley White hydrophobicity, Kappa<sup>83</sup>) with  
 836 localCIDER<sup>84</sup>. The OmegaWFYL\_DE mixture parameter computes the mixture  
 837 statistic between W,F,Y,L residues and D,E residues using the  
 838 `seq.get_kappa_X(['D','E'],['W','F','Y','L'])` function in localCIDER<sup>85</sup>. We predicted  
 839 intrinsic disorder with MetaPredict2<sup>86</sup>. We counted motifs with regular  
 840 expressions in Python with the "re" package.

841 The MAFFT algorithm aligns the WxxLF motif for all but three orthologs.  
 842 For three orthologs, in the `Full_length_ortholog_dataframe`, we corrected the  
 843 "WxxLF motif location" parameter using the coordinates from the MSA. These  
 844 species are the only ones outside the Ascomycota that have the motif. We  
 845 suspect the WxxLF motif convergently evolved in these distance orthologs  
 846 because the context is very different and H rich. `Blastocladiomycota_jgi|Catan2|`  
 847 `1097078|CE97078_6759`, `Blastocladiomycota_jgi|Catan2|1466814|`  
 848 `fgenesh1_pg.199_#_9`, and `Blastocladiomycota_jgi|Catan2|1506241|`  
 849 `gm1.11555_g`.

850 To predict helical propensity of ortholog sequences, we used the Sparrow  
 851 package in Python [<https://github.com/idptools/sparrow>]. A region was called  
 852 helical if it contained five adjacent residues with over 50% chance of being  
 853 helical. A large proportion of sequences have no residues with a >50%  
 854 probability of being helical in this region. We consider this predictor to capture  
 855 the propensity to form a helix in some context. To count proline residues in the  
 856 region homologous to the known helix, we used the five AA upstream and five AA  
 857 downstream of the WxxLF motif. From the 500 orthologs in the MSA, there are

858 115 unique 15 AA regions around the WxxLF motif; twenty-three contain three  
859 prolines (20%) and three contain four prolines (2.6%).

### 860 **Imputing activity in the full-length orthologs**

861 We used the tile data to impute the activity of each position in each of the  
862 full-length orthologs. The 19099 recovered tiles mapped to 68577 locations on  
863 the orthologs (each tile matched to 3.6 orthologs on average). We used a second  
864 order Loess smoothing (20 nearest points with the `loess.loess_1d.loess_1d()`  
865 function) across tiles to impute the activities of all positions in the 502 unique  
866 orthologs. This quadratic smoothing can cause artifacts on the extreme ends of  
867 the protein, such as predicting negative activity. To remove this artifact, we  
868 constrained the imputed activity to be no more than the maximum measured  
869 and no less than the minimum measured in that ortholog.

870 To validate the Loess smoothing, we averaged together all activities for all  
871 tiles that overlapped a position, equally weighing all tiles. These averages were  
872 more jagged because of the stepwise nature of the tiles. This simple average  
873 also created artifacts at the ends of the protein where only one tile is present.  
874 The Loess and average smoothing methods agreed well (97% had Pearson  $R >$   
875 0.80) (**Figure S33**).

876 We used the imputed activities to create the heatmaps to visualize activity  
877 across the orthologs. We tried many variations of these heatmaps but ultimately  
878 found that aligning the sequences on the start of the DBD or on the WxxLF motif  
879 was most informative. In the main text, we removed the twenty-seven longest  
880 sequences to make the visualization easier to display but added most of them  
881 back in Figure S9.

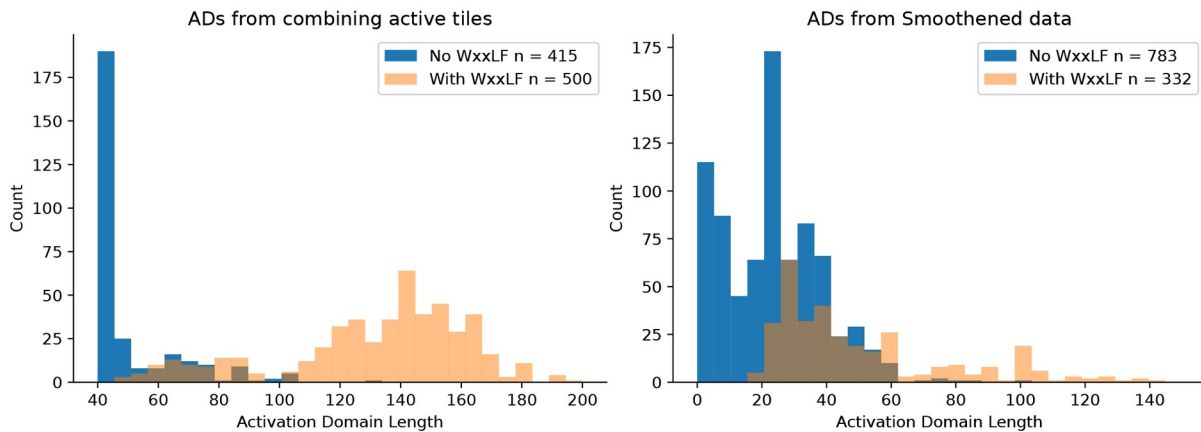
882 We tested the hypothesis that insertions are enriched for active tiles by  
883 projecting activity onto the MSA. We defined insertions as the positions in the  
884 MSA with residues (non-gaps) in less than 1% of sequences ( $n < 5$ ), which yielded  
885 880/2690 (32.7%) of positions. In a two-sided t-test of the imputed activities of the  
886 insertion positions compared to all other positions, insertions were less active ( $p <$   
887  $1e-52$ ). We concluded that insertions are depleted for sequences with activation  
888 domain activity in *S. cerevisiae*.

889 To estimate the activity at the WxxLF motif, we used the integral of the  
890 imputed activity from -10 to +10 around the W of the WxxLF motif. When this  
891 integral was below our activity threshold, we called sequences inactive in this  
892 region. Using this integral, ninety-two unique sequences had high activity  
893 ( $>150000$ ) and thirteen unique sequences had low activity, less than our activity  
894 threshold. Thirty-three had intermediate activity.

895 For motif enrichment, we performed a Welch's t-test assuming unequal  
896 variances `stats.ttest_ind(Sequences_WITH_Motif,Sequences_WITHOUT_Motif,`  
897 `equal_var=False)`.

898 To count activation domains on each TF, we combined active overlapping  
899 tiles, taking the union. With this method, we found 500 ADs with the WxxLF motif  
900 and 415 ADs without the WxxLF motif. We required more than forty residues  
901 between activation domains before they were called as two separate domains.  
902 Calling activation domains from the imputed activity map gives slightly different  
903 results because some very close double peaks are split. With the smoothed data,  
904 there are 332 ADs with the WxxLF motif and 783 ADs without the WxxLF motif.





905

## 906 ANOVA

907 We used ordinary least squares regression (OLS) to create a baseline  
 908 model for how composition controls activation domain function. We used ANOVA,  
 909 OLS, and adjusted R-squared to compare models. See the Composition\_ANOVA  
 910 jupyter notebook for the full analysis. Briefly, we used the `ols(formula,`  
 911 `ANOVA_DF).fit()` function from the `statsmodels` package to fit the model, find  
 912 coefficients, and compute adjusted R-squared values. We used the  
 913 `anova_lm(model, typ=2)` function to find the sum of squares explained by each  
 914 parameter. We used a Bonferroni multiple hypothesis correction to remove non-  
 915 significant parameters and refit the model. In most cases, one iteration was  
 916 sufficient to get a model where all parameters were significant. For the  
 917 dipeptides, we used two interaction terms. All ANOVA parameters are in **Table**  
 918 **S9**.

919 OLS regression on single amino acids explains 49.9% of variance in  
 920 activity (**Table 1**, AUC = 0.9346, PRC = 0.7620, **Table S9**). Iteratively removing  
 921 non-significant parameters led to sixteen residues which explain 49.9% of  
 922 variance. We repeated the regression with 400 dipeptides and found 69  
 923 significant parameters that explain 60.2% of the variance in activity (**Table 1**,  
 924 AUC = 0.9472, PRC = 0.8190). Half the variance in activity could be explained by  
 925 composition alone and dipeptides offered ~10% improvement.

926 We predicted *de novo* motifs using the DREAM suite and then repeated  
 927 the OLS ANOVA analysis using the motifs. We performed *de novo* motif searching  
 928 on multiple slices of the data, but highly active (n=3524) vs. inactive (n=15575)  
 929 were the most interpretable and gave the clearest signal in the ANOVA analysis.  
 930 First, we ran the package STREME from the MEME suite to discover motifs that  
 931 are enriched in a list of sequences relative to a user-provided control list.

932 For the OLS on *de novo* motifs, we used the motif counts provided by the  
 933 DREAM motif prediction software (**Table S10**). For simplicity, in the parameter  
 934 table, we refer to each motif as a string, but we used the PWM for actually  
 935 finding motifs in each sequence with FIMO.

## 936 Machine learning

937 We predicted activities on full length orthologs using publicly available  
 938 models, TADA, ADpred, and PADDLE<sup>18,21,23,24</sup>. All models were run on the SAVIO  
 939 high performance computing cluster at UC Berkeley. TADA uses 40 AA windows,  
 940 ADpred, 30 AA windows, and PADDLE 53 AA windows. For each TF, we tiled at 1  
 941 AA increments, spanning the full proteins (e.g. 1-40, 2-41 etc). For full length TF  
 942 analysis, we corrected the inferred activity at each position (Loess smoothing)  
 943 with the predictions at each position. The smoothed data averages out some

944 measurement noise so all the model performance is improved on smoothed data.  
 945 For individual tile analysis, we used the center aligned score. We also tried  
 946 maximum scores, average scores, and other variations, but chose center  
 947 aligned. ROC and PRC analyses were performed with the sklearn python  
 948 package.

949 Predicting the impact of mutating F residues in the central activation  
 950 domains. We tile the 138 unique 70AA central regions into 40AA tiles spaced  
 951 every 1 amino acid. For each tile, we computationally mutated each F  
 952 individually, all pairs, all triplets, and all sets of four or more. For each mutant,  
 953 we predicted activity. The mutants are predicted to have less activity. For each  
 954 mutant, we also computed the change in activity. Finally, we grouped the  
 955 changes in activity based on the conservation of each F residue.

## 956 **Pax6 alignments**

957 BLAST alignment of mouse Pax6 (P63015) and *D. melanogaster* Eyeless  
 958 (O18381) was performed with the Uniprot canonical sequences. We calculated  
 959 the DBD percent identity using the longest aligned region that encompassed the  
 960 annotated DBD (5-135 and 157-187, respectively). We realigned the regions C-  
 961 terminal to the end of this DBD alignment and found three regions with modest-  
 962 to-high scores:  $(79+16+7)/287 = 35.5\%$  residues identical and  $(88+28+11)/287$   
 963  $= 44.3\%$  residues similar in the three regions. We summed the number of  
 964 identical or similar residues to compute similarity. We used the shorter mouse  
 965 IDR length as the denominator, overstating conservation. Alignments are in  
 966 **Figure S34**. Using the more permissive BLOSSUM90 matrix yielded a fourth  
 967 small aligned region that increased the similarities:  $(79+16+14+7)/287 = 40.4\%$   
 968 residues identical and  $(88+26+18+11)/287 = 50\%$  residues similar.

## 969 **Datafiles**

970 All the raw sequencing data has been deposited at NIH SRA Accession

971 #PRJNA1186961: <http://www.ncbi.nlm.nih.gov/bioproject/1186961>

972 All the analysis scripts are deposited on github via Zenodo:

973 10.5281/zenodo.14201918

974 <https://github.com/staller-lab/Gcn4-evolution>

975 [github.com/staller-lab/labtools/tree/main/src/labtools/adtools](https://github.com/staller-lab/labtools/tree/main/src/labtools/adtools)

976 <https://github.com/staller-lab/Gcn4-evolution>

977 All the processed data is attached in supplemental tables (**Tables S5 - S7**).

978 Processed sequencing read counts are in **Table S13**.

979

980 The 'masterDF' dataframe contains each designed tile (**Table S5**). Tiles  
 981 that were not measured have activity recorded as nan or 0. The 'orthologDF'  
 982 dataframe contains all tiles associated with each original full-length ortholog  
 983 (**Table S6**). As a result, tiles occur multiple times because they map to multiple  
 984 orthologs. The 'NativeLocation' is the position of the tile relative to the first  
 985 amino acid. The 'NormLocation' is the position of the tile relative to the WxxLF  
 986 motif. Finally, the 'FullOrthoDF' dataframe contains one entry for each full-length  
 987 ortholog, and each column contains an array with values for each position  
 988 (**Table S7**), such as imputed activity at each position and local charge from

989 localCIDER. The location of the bZIP DNA-binding domain was identified with the  
990 InterPro signature (IPR004827).

991

## 992 **Description of python analysis scripts**

- 993 ● Step2\_AddSeqFeaturestoDataFrame\_Oct\_2024.ipynb
  - 994 ○ Combines the data from the two replicates.
  - 995 ○ Computes many sequence features, like net charge.
- 996 ● AD\_AlignmentDists.ipynb
  - 997 ○ This script looks at the Edit distances between pairs of sequences. It
  - 998 shows that many changes in sequence do not change activity.
- 999 ● AD\_properties Fall 2024.ipynb
  - 1000 ○ This script explores how sequence properties, like AA abundance or
  - 1001 motif locations, contribute to activation domain activity.
  - 1002 ○ Contains main figure panels
- 1003 ● Composition\_ANOVA Fall 2024.ipynb
  - 1004 ○ ANOVA analysis of OLS regression on composition and dipeptides
- 1005 ● Controls\_oct024.ipynb
  - 1006 ○ Barplots for control sequences
  - 1007 ○ Reproducibility analysis
- 1008 ● Full\_Length\_TFs\_Heatmaps\_Fall 2024.ipynb
  - 1009 ○ Script to make heatmaps of full-length orthologs
- 1010 ● Sensu strictu v2.ipynb
  - 1011 ○ Plot activity traces of *S. cerevisiae* and closest species
- 1012 ● Gaussian\_Threshold.ipynb
  - 1013 ○ Analysis of inactive sequences to find activity threshold
- 1014 ● YeastAnalysisfunctions.py
  - 1015 ○ Support functions for visualizing data

## 1017 **Acknowledgments**

1018 We would like to thank Nick Ingolia, Zeba Wunderlich, Rachel Brem, Alex  
1019 Holehouse, Shahar Sukenik, Micheal Botchen, and Ashley Wolf for helpful  
1020 comments on the manuscript. Sumanth Mutte for finding the initial orthologs. We  
1021 thank Lucia Strader, Nicholas Morffy, Ross Sozzani, Lisa Van den Broeck, Hunter  
1022 Nisonoff, and Jennifer Listgarten for helpful discussions, and Nick Morffy and  
1023 Lucia Strader for the yeast genome targeting plasmids. Igor Grigoriev identified  
1024 the deprecated *Tortispora caseinolytica* gene models. Weijing Tang performed  
1025 exploratory analyses not included in the final manuscript. The Regents of the  
1026 University of California have filed a patent based on the findings of this study.  
1027 The DHY213 BY super host strain used for library construction was a generous  
1028 gift from Angela Chu and Joe Horecka, and requests for this strain should be  
1029 directed to them.

## 1030 Funding

1031 CJL training grant T32HG4725. AL UC Berkeley URAP. MAZ T32GM148378.  
 1032 MS and SRK UC Berkeley SEED Scholars Program. SRK UC Berkeley SURF. AF  
 1033 biophysics training grant T32GM146614. GPS UC Berkeley BSP scholar, McNair  
 1034 Scholar, and UC Berkeley SURF. This work was supported by the Burroughs  
 1035 Wellcome Fund PEDP, Simons Foundation grant 1018719 to MVS, NSF grant  
 1036 2112057 to MVS, and NIH grant R35GM150813 to MVS. MVS is a Chan  
 1037 Zuckerberg Biohub – San Francisco Investigator.  
 1038  
 1039

## 1040 Supplementary Note: Additional analysis of the 1041 orthologs

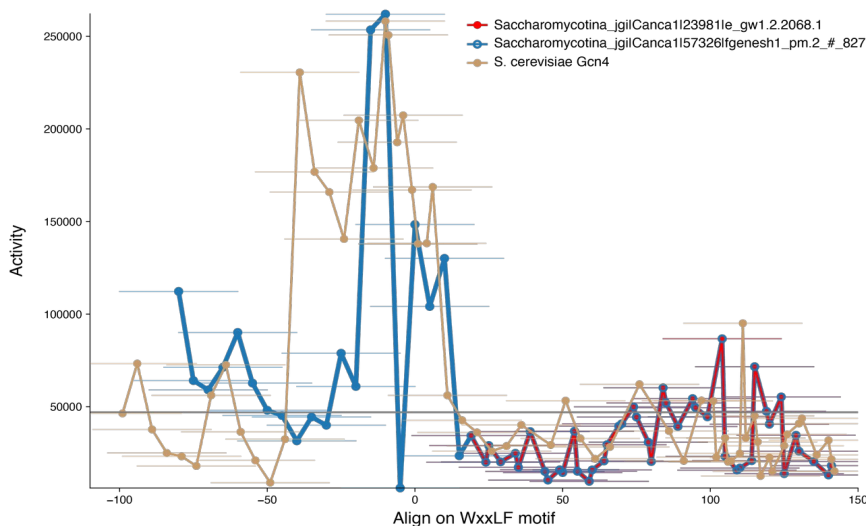
### 1042 Selection of the Gcn4 orthologs

1043 We chose a diverse set of orthologous Gcn4 protein sequences for  
 1044 functional characterization in *S. cerevisiae*. We started with a set of forty-nine  
 1045 previously identified orthologs<sup>16,41,55</sup>. In these, 48/49 contain an WxxLF motif.  
 1046 Next, we scanned 207 representative proteomes from the MycoCosm database,  
 1047 sampling the diversity of fungal genomes (**Figure S1, S2**). To distinguish Gcn4  
 1048 orthologs from other basic-leucine zipper (bZIP) domain TFs, we required the  
 1049 presence of both a bZIP DNA-binding domain (IPR004827) and the WxxLF motif.  
 1050 This computational screen yielded 1188 hits in 129 genomes. There are 502  
 1051 unique Gcn4 ortholog sequences that we used for all our experiments and  
 1052 analyses (**Figure S1**). These sequences span nearly all the Ascomycota, the  
 1053 largest phylum of Fungi, representing >600 million years of evolution<sup>42</sup>. The 502  
 1054 unique orthologs have variable lengths (**Figure 1A**), but the DBD is at the C-  
 1055 terminus in 500, and the distance between the WxxLF motif and the DBD is very  
 1056 consistent (**Figure 1B**).

1057 The Gcn4 MSA typifies eukaryotic TF evolution, with a highly conserved  
 1058 DBD and lower conservation in the rest of the protein (**Figure 1C**). Sequence  
 1059 divergence is driven by insertions: 88% of columns in the MSA contain fewer  
 1060 than 5% of sequences ( $n < 25$ ) and 54% of columns contain <1% of sequences  
 1061 ( $n < 5$ ) (**Figure S4**). Without user input, the MAFT algorithm aligned the WxxLF  
 1062 motif in nearly all sequences (Methods). We suspect that MAFT aligned nearly all  
 1063 WxxLF motifs because the distance between this motif and the DBD is highly  
 1064 consistent. Distant pairs of sequences do not align outside of the DBD, but we  
 1065 have enough sequences to bridge the full diversity of the collection. The central  
 1066 activation domain shows intermediate levels of conservation largely driven by  
 1067 the WxxLF motif. Since we required all the orthologs to contain a WxxLF motif,  
 1068 the conservation of this motif is overstated in **Figure 1C**, but we independently  
 1069 verified that this motif is the most conserved sequence outside the DNA-binding  
 1070 domain using a HMMER search of fungal TFs (**Figure 2C**).

## 1071 All orthologs are activators

1072 To show that all the orthologs contain at least one active tile, we used  
 1073 multiple thresholds. As an unbiased threshold for modest activity, we fit a  
 1074 Gaussian distribution to the inactive sequences. Using this highly permissive  
 1075 threshold, all orthologs have at least one tile that is active. As a stringent  
 1076 threshold for activity we doubled this threshold, or used the top 20% of  
 1077 sequences, which yielded very similar values. At the stringent threshold, there is  
 1078 only one ortholog with no active tiles, Canca1\_23981 from *Tortispora*  
 1079 *caseinolytica*. This ortholog is an alternative gene model for the Canca1\_57326  
 1080 protein, which contains an additional 99 N-terminal residues with twenty-three  
 1081 overlapping active tiles that comprise two activation domains, the second of  
 1082 which overlaps the WxxLF motif. The short form of the protein starts at the  
 1083 WxxLF motif. Based on improved, transcript-based gene models, the short  
 1084 version, Canca1\_23981, is likely a computational annotation error. There is more  
 1085 support for the long version, Canca1\_57326. Given the relatively weak evidence  
 1086 supporting the one potential exception, we conclude all of the Gcn4 orthologs  
 1087 are activators.



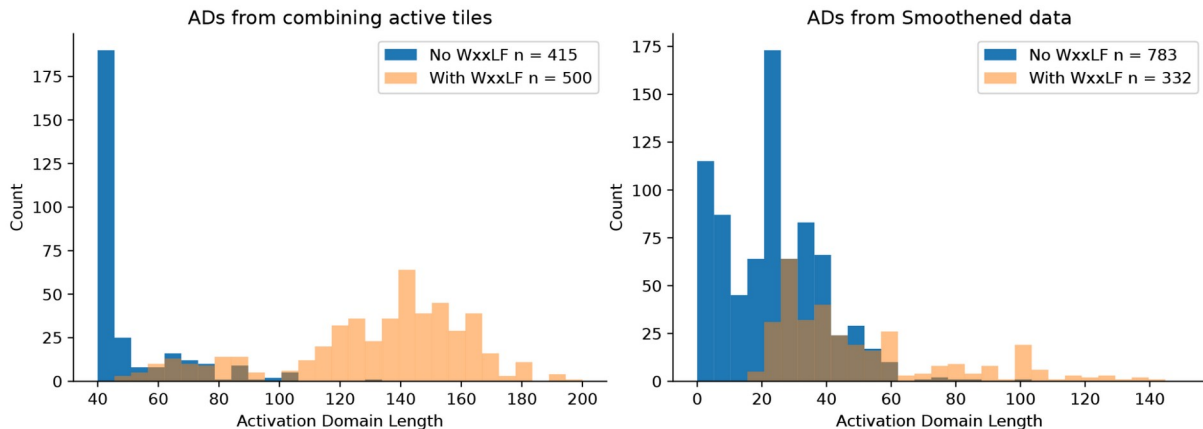
1088 Alternative gene models from *Tortispora caseinolytica*. Canca1\_23981 (red) and  
 1089 Canca1\_57326 (dark blue) are alternative gene models for the same locus.  
 1090 Importantly, they are identical, so the red overlaps the dark blue.  
 1091

## 1092 Activation domains per ortholog

1093 Longer TFs often have more active tiles (**Figure 2E**). When we merged  
 1094 overlapping active tiles, most orthologs had more than one activation domain  
 1095 (**Figure 2F**). The lengths of the merged activation domains are bimodal, but  
 1096 they are generally <200 AA (**Figure 2G**).

1097 We used two methods to count activation domains on each ortholog. First,  
 1098 we aggregated overlapping active tiles. This method biases towards fewer longer  
 1099 activation domains because there must be more than forty AA between active  
 1100 regions for them to be called as separate activation domains. With this method,  
 1101 245 orthologs (48.8% ) have only one activation domain, and for all these  
 1102 orthologs, the AD overlaps the WxxLF motif. In total 500 activation domains  
 1103 contained the WxxLF motif, and these were longer than N-terminal activation

1104 domains. There were also many single-tile activation domains. Second, we used  
 1105 the smoothed data to find activation domains. This method averages out some  
 1106 experimental noise and shortens active regions. In this approach, there are only  
 1107 332 orthologs with an activation domain that contains the WxxLF motif,  
 1108 consistent with the peak of activity being upstream of this motif. There are more  
 1109 N-terminal activation domains, and they are shorter than activation domains with  
 1110 the WxxLF motif. In both methods, the sequences of the N-terminal activation  
 1111 domains are diverse.



1112

### 1113 Clusters of aromatic and leucine residues make large 1114 contributions to function

1115 In the control activation domains, all published motifs of aromatic and  
 1116 leucine residues made large contributions to activity, but no individual motif was  
 1117 sufficient for full activity. Historically, *S. cerevisiae* Gcn4 is annotated with two  
 1118 activation domains: the CAD is residues 101-140, while the N terminal activation  
 1119 domain (NAD) is residues 1-100 (**Figure 2A**)<sup>78,87,88</sup>. There are six published  
 1120 motifs, F9 F16 (FxxxxxxF), F45 F48 (FxxF), F67 F69 (FxF), F97 F98 (FF), M107  
 1121 Y110 L113 (MxxYxxL or MFxYxxL), and W120 L123 F124 (WxxLF)<sup>78,88</sup>. The CAD  
 1122 has two motifs that make large contributions to activity<sup>78,88</sup> (**Figure 2B**). The  
 1123 strongest tile from Gcn4 was the junction of the NAD and CAD (residues 90-129),  
 1124 which we call the altCAD, a region with three motifs<sup>78</sup> that make large  
 1125 contributions to function (**Figure 2B, S7**). All published motifs are enriched in  
 1126 active tiles (**Figure S20A**), and tiles with multiple motifs are more likely to have  
 1127 high activity (**Figure S20B**). However, in our sequences and an independent set  
 1128 of fungal orthologs, only the WxxLF motif is well conserved (**Figure 1C, 2C, S3**).  
 1129 We do not see reemergence of any published motifs. The hydrophobic motifs  
 1130 essential for function in *S. cerevisiae* are not conserved and do not experience  
 1131 evolutionary turnover.

1132

### 1133 Sequence features of strongly active tiles

1134 The Gcn4 ortholog tiles efficiently detected known sequence features of  
 1135 strong yeast activation domains. Acidic, aromatic, leucine, and methionine  
 1136 residues make the largest contributions to activity<sup>16,18-23,28,31</sup> (**Figure 4A, C**).  
 1137 Aromatics generally increase activity, but too many aromatic residues reduces  
 1138 activity (**Figure S17F,G**), a non-monotonic trend previously seen only in

1139 synthetic peptides<sup>18</sup> and mutant activation domains<sup>28</sup>. This non-monotonicity is a  
 1140 key piece of evidence supporting the acidic exposure model because it shows  
 1141 how too many hydrophobic residues can overwhelm the exposure capacity of the  
 1142 acidic residues<sup>26-28</sup>. Moreover, aspartic acid (D) makes much stronger  
 1143 contributions to activity than glutamic acid (E) (**Figure 4C**), which has only been  
 1144 seen in mutants<sup>18</sup> and weakly in plant activation domains<sup>23</sup>. We suspect this  
 1145 effect occurs because the negative charge is closer to the peptide backbone,  
 1146 leading to a stronger solvation effect and more exposure of nearby hydrophobic  
 1147 residues<sup>43</sup>. This modestly sized dataset gave a much clearer picture of key  
 1148 sequence properties than much larger datasets<sup>18,21,23</sup>, indicating that orthologs  
 1149 provide a very efficient set of sequences for learning the sequence features that  
 1150 control function (**Figure S20**).

### 1151 **Evidence for negative (purifying) selection**

1152 The Yeast Gene Order Browser (YGOB) contains a high quality set of thirty-  
 1153 six true homologs inferred from chromosomal synteny. All of the species  
 1154 following the whole genome duplication contain only one Gcn4 homolog,  
 1155 suggesting there is no advantage of retaining two copies. This result suggests  
 1156 that most species will have just one true homolog. The YGOB analysis of full-  
 1157 length TFs shows negative selection (**Figure S1E**), implying there is pressure to  
 1158 maintain a functional protein. This weak negative selection and large protein  
 1159 diversity supports the idea that the neutral space is very large and that the Gcn4  
 1160 sequence can drift.

1161 Enforcing the presence of a strict WxxLF motif left out one true homolog  
 1162 from YGOB, *Zygosaccharomyces baili* ZYBA0L03268g, which has an insertion in  
 1163 the WxxLF motif yielding **WPSLEPLF**. This sequence was not included in our  
 1164 current experiment but was measured as highly active in one replicate in Staller  
 1165 et al. 2018, suggesting activation domain function is also conserved in this  
 1166 ortholog. This example reinforces the idea that motifs can be flexible.

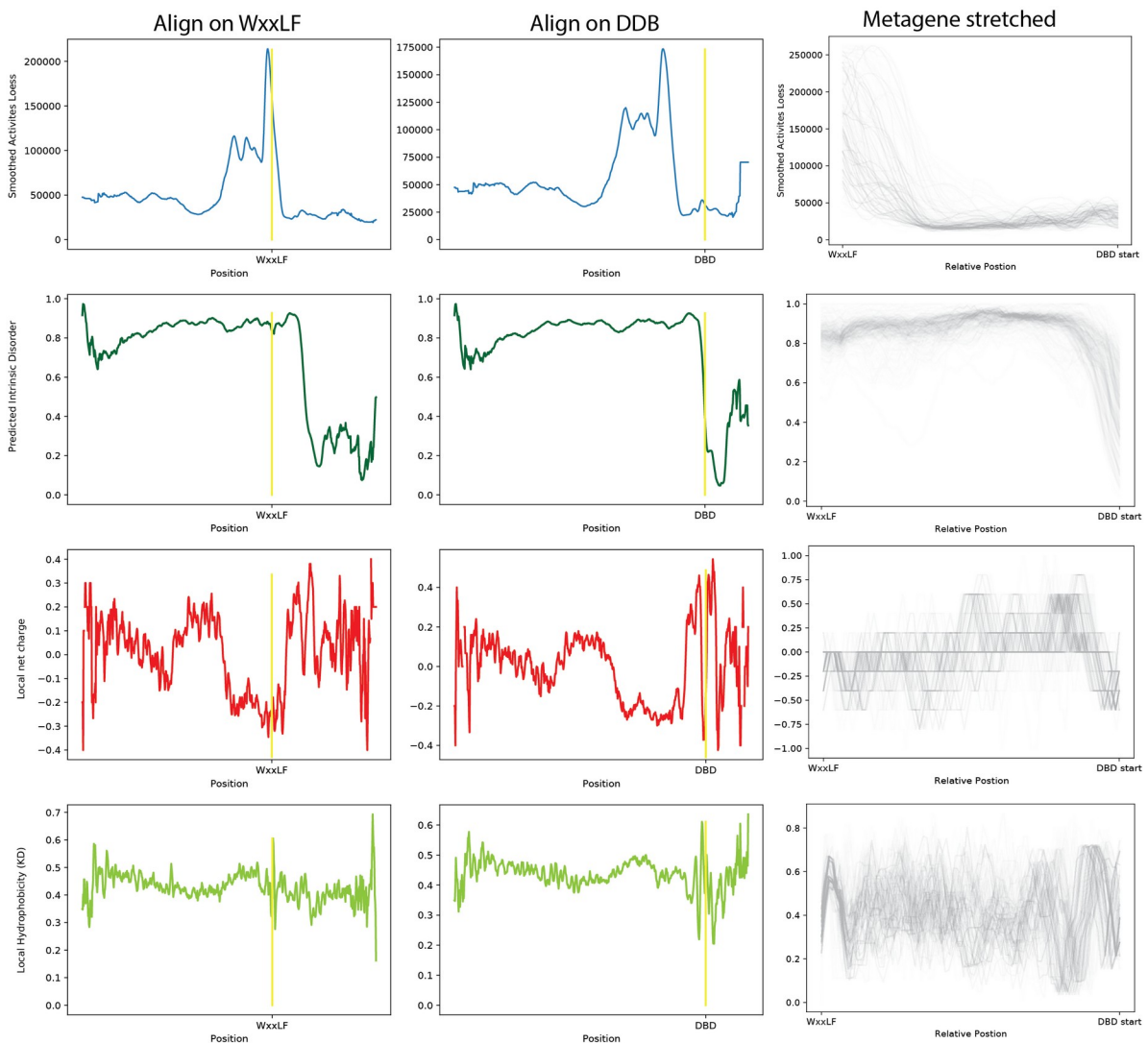
### 1167 **Analysis of the Gal11/Med15 coactivator**

1168 The best characterized coactivator of *S. cerevisiae* Gcn4 is Gal11/Med15.  
 1169 Med15 contains four regions that bind to Gcn4, the KIX domain and three  
 1170 activation domain binding domains (ABD1, ABD2, ABD3)<sup>44</sup>. Activity of our P3  
 1171 promoter is well correlated with *in vitro* binding to Med15<sup>18</sup>, indicating this  
 1172 promoter is a reliable reporter of binding to Med15. We collected a set of 653  
 1173 Gal11 orthologs from the Y1000+ genomes and created an MSA. The KIX, ABD1,  
 1174 and ABD3 domains are more conserved than the rest of the protein. ABD2  
 1175 approaches the rest of the protein. The residues of the ABD1 domain that  
 1176 contact Gcn4<sup>45</sup> are reasonably conserved, but not more conserved than the rest  
 1177 of ABD1 (**Figure S31**). Overall the conservation of Med15 is much higher than  
 1178 Gcn4.

### 1179 **Analysis of the spacer sequence between the CAD and DBD**

1180 The distance between the WxxLF motif (CAD) and the DBD is highly  
 1181 conserved and may be an entropic spacer. The amino acid sequence of this  
 1182 spacer is very poorly conserved, but both the undulating charge pattern and the

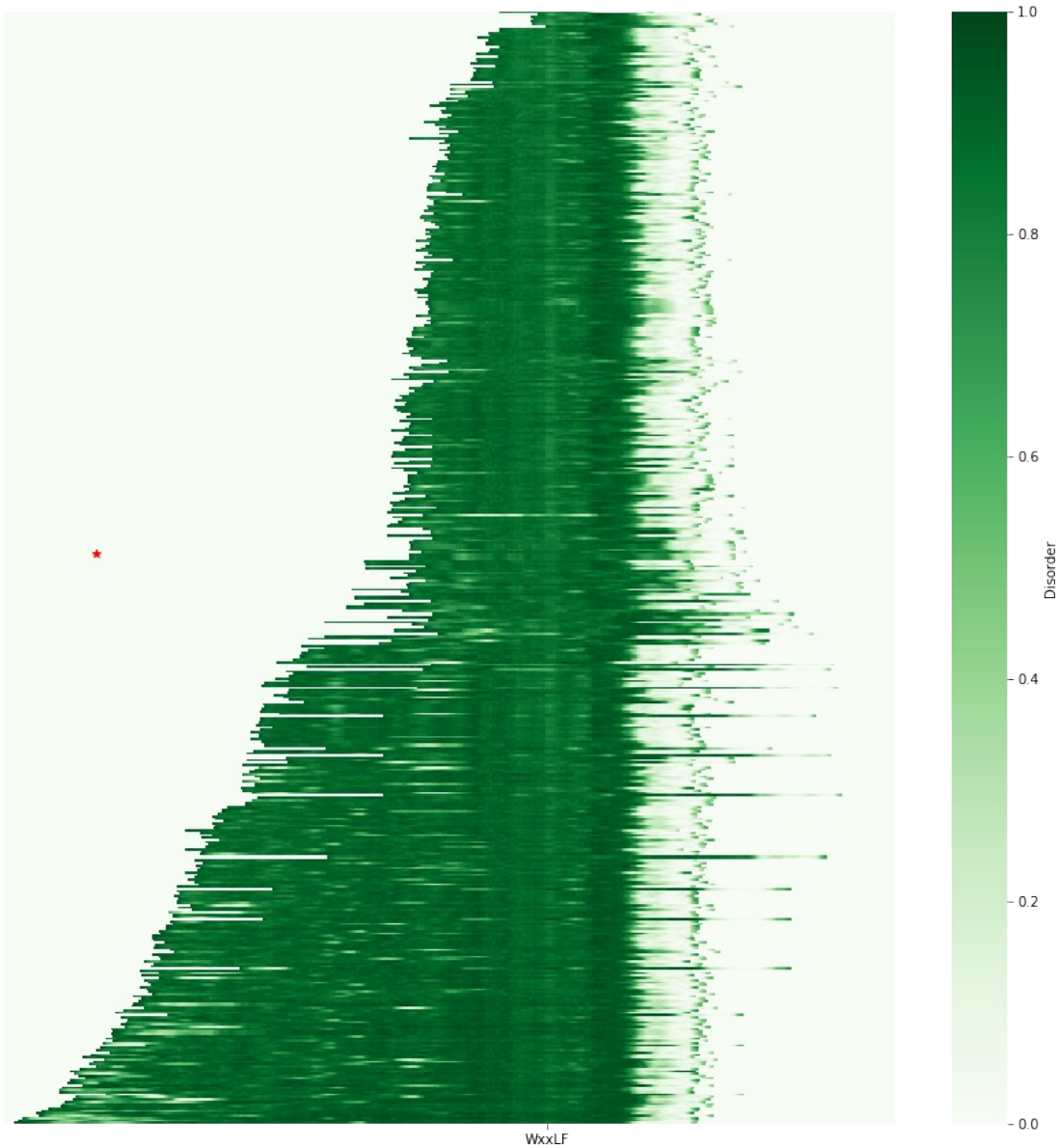
1183 high degree of predicted intrinsic disorder are conserved. NMR data clearly  
 1184 indicates that *S. cerevisiae* Gcn4 is fully disordered in solution and that the DBD  
 1185 folds upon binding DNA and the CAD folds upon binding Med15. Predicting this  
 1186 pattern is difficult, and Gcn4 has become a stringent test for intrinsic disorder  
 1187 prediction algorithms. AlphaFold predicts the DBD correctly. AlphaFold predicts  
 1188 many short, low-confidence helices outside the DBD, but none overlap the CAD  
 1189 NMR helix. To predict intrinsic disorder of the orthologs, we used Metapredict,  
 1190 which carefully examined performance on Gcn4 during algorithm development  
 1191 <sup>86</sup>. Based on this analysis, the most disordered region in all orthologs is the  
 1192 sequence between the CAD and DBD. This region has a positive to negative  
 1193 charge undulation just before the DBD.  
 1194



1195  
 1196 **Analysis of the spacer sequence between the WxxLF motif and the DBD**  
 1197 Left panels align position on the WxxLF motif. Middle panels align position on the  
 1198 DBD. The spacer is the sequence between these landmarks. Imputed activity of  
 1199 the spacer is low. Predicted intrinsic disorder of the spacer is high  
 1200 (Metapredict2). Negative charge undulates between the landmarks. The region  
 1201 right after the WxxLF is negatively charged, followed by a positively charged

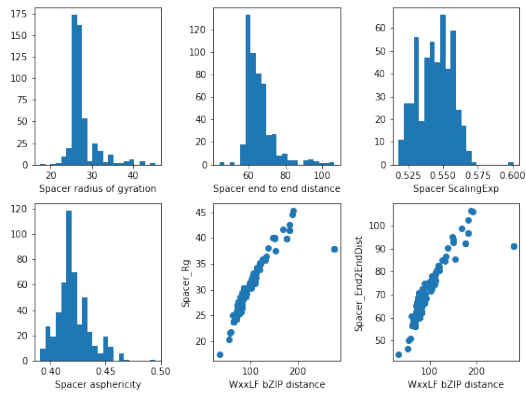


1202 region and another net negative region just before the positively charged DBD.  
 1203 Hydrophobicity is high throughout.  
 1204



1205  
 1206 **Predicted disorder in the spacer sequence peaks between the WxxLF**  
 1207 **motif and the DBD**

1208  
 1209  
 1210 We speculate this region is a conserved entropic spacer that keeps the  
 1211 activation domain away from the DBD and exposed to partners. *S. cerevisiae* has  
 1212 uncommonly long spacing between the WxxLF and DBD (**Figure 1B**, red arrow).  
 1213 We tested this idea by predicting biophysical parameters with Albatross<sup>47</sup>. We  
 1214 see that the predicted radius of gyration (estimate of ensemble size) and end-to-  
 1215 end distance distributions are very tight, implying that there might be some  
 1216 selection to maintain a specific 3D spacing distance.



1217

1218

## 1218 **Computationally predicted scaling exponents and biophysical properties**

### 1219 **of the spacer from Sparrow.**

1219

1220

1220 The highly consistent predicted dimensions support they hypothesis that this  
1221 spacer is keeping the central activation domain away from the DBD.

1222

1223

1224

1225

### 1225 **Additional analysis of tile sequence properties**

1226

1227

1228

1229

1230

1231

1232

1233

1234

1226 Yeast activation domains are more reliant on aromatic residues than  
1227 leucine residues. This difference is illustrated by the human CITED2 activation  
1228 domain. In human cells, the aromatic residues make small contributions to  
1229 CITED2 function, but in yeast, these residues make large contributions to  
1230 function. Leucine residues contribute to CITED2 function in both yeast and  
1231 human cells. The mutant of CITED2 without aromatic residues was the strongest  
1232 sequence with no aromatic residues (**Figure S7B**). It is mildly surprising that  
1233 CITED2 works in yeast because its primary coactivator partner, TAZ1, is not  
1234 present in yeast.

1235

1236

1237

1238

1239

1240

1241

1242

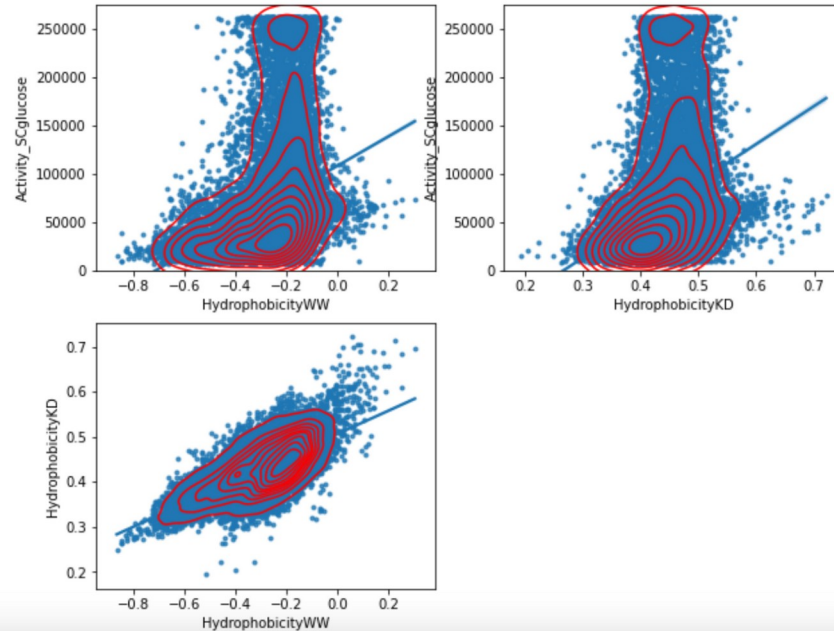
1243

1244

1235 Sanborn et al. argued that the Wimley White hydrophobicity (WW) score  
1236 was well correlated with AD activity <sup>18</sup>. We had previously used the Kyte Doolittle  
1237 hydrophathy (KD) score and found no correlation in designed mutants <sup>16</sup>. The  
1238 largest difference between these tables is tryptophan, W, which has a high value  
1239 on WW and moderate value on KD. Since W makes large contributions to activity,  
1240 we believe that the number of W's drives the conclusion by Sanborn et al. 2021.  
1241 In our Gcn4 ortholog tiles, the two hydrophobicity scores are well correlated with  
1242 each other. Both have similar, low correlations with activity. Some  
1243 hydrophobicity is required for activity. The combination of acidity and  
1244 hydrophobicity is more predictive than hydrophobicity alone.

Out[18]:

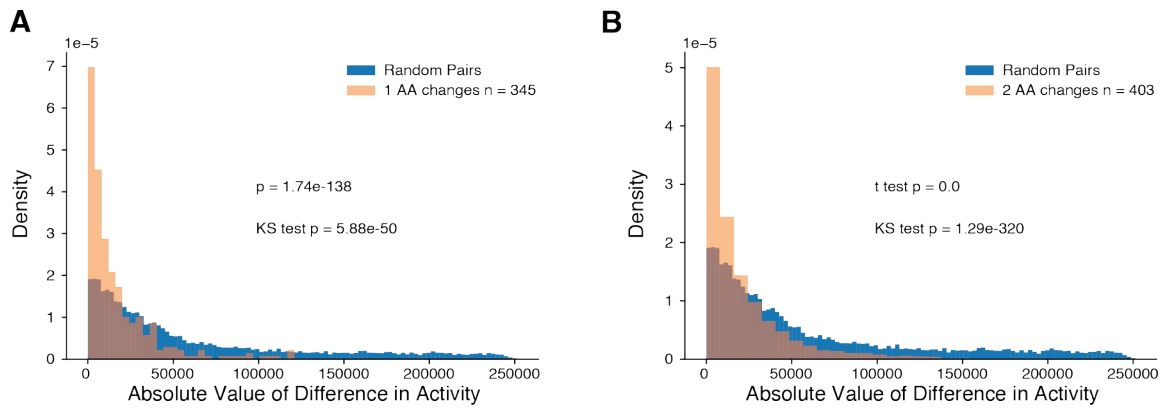
	HydrophobicityKD	HydrophobicityWW	Activity_SCglucose
HydrophobicityKD	1.000000	0.689229	0.341205
HydrophobicityWW	0.689229	1.000000	0.354076
Activity_SCglucose	0.341205	0.354076	1.000000

1245  
1246

## 1247 **Naturally occurring changes in sequence generally do not change** 1248 **activity**

1249 Most naturally occurring sequence changes do not change activity.  
1250 Starting with the altCAD as an anchor, we identified related sequences with  
1251 increasing edit distance. As sequence divergence increased, all the natural  
1252 sequences maintained high activity. In contrast, designed mutants show that  
1253 small changes in sequence can cause loss of activity. Large effect changes are  
1254 absent from the evolutionary record. This result supports a model where neutral  
1255 drift and weak negative selection maintain activation domain activity.

1256 Next, we compared pairs of sequences that differed by one or two amino  
1257 acids. As a null model for differences in tile activities, we chose 10000 random  
1258 pairs of tiles and computed the difference between their activities. The  
1259 distribution of activity differences between tiles that differ at 1-2 amino acids is  
1260 much smaller.



1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

In most cases, there was little-to-no change in activity. We imposed a strong threshold for change in activity: either one member of the pair was active and the other inactive, or both were active but differed in activity by more than 50%. In the majority of cases that change activity, the sequence change was interpretable by our acidic exposure model: the stronger tile had additional acidic or hydrophobic residues. Of the 345 pairs of tiles that differ at a single position, 15 pairs (2.5%) had different activities and 9 supported the acidic exposure model. In four cases, an L or M was added that increased activity. In one case, an E>D change increased activity. In three cases, adding an S or G, which promotes disorder and expansion, increased activity. Of the 403 pairs of tiles that differ at two positions, 27 changed activity (7%). Two of these were designed mutants in the altCAD, FF>AA and LL>AA, both of which caused large decreases in activity (**Figure S7**). 17/27 cases (or 15/25 natural cases) supported the acidic exposure model. These data further support the mounting evidence that activation domains are robust enough to maintain because most single and double AA changes do change activity.

1279

## References

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1. Onuma, Y., Takahashi, S., Asashima, M., Kurata, S. & Gehring, W. J. Conservation of Pax 6 function and upstream activation by Notch signaling in eye development of frogs and flies. *Proceedings of the National Academy of Sciences* **99**, 2020–2025 (2002).
2. Lynch, V. J. & Wagner, G. P. Revisiting a classic example of transcription factor functional equivalence: are Eyeless and Pax6 functionally equivalent or divergent? *J. Exp. Zool. B Mol. Dev. Evol.* **316B**, 93–98 (2011).
3. Halder, G., Callaerts, P. & Gehring, W. J. Induction of Ectopic Eyes by Targeted Expression of the eyeless Gene in *Drosophila*. *Science* **267**, 1788–1792 (1995).
4. Andersson, L. S. *et al.* Mutations in DMRT3 affect locomotion in horses and spinal circuit function in mice. *Nature* **488**, 642–646 (2012).
5. Lynch, V. J., May, G. & Wagner, G. P. Regulatory evolution through divergence of a phosphoswitch in the transcription factor CEBPB. *Nature* **480**, 383–386 (2011).
6. Gao, Y. *et al.* The emergence of Sox and POU transcription factors predates the origins of animal stem cells. *Nat. Commun.* **15**, 1–16 (2024).
7. Chothia, C. & Finkelstein, A. V. The classification and origins of protein folding patterns. *Annu. Rev. Biochem.* **59**, 1007–1039 (1990).
8. Lim, W. A. & Sauer, R. T. Alternative packing arrangements in the hydrophobic core of lambda repressor. *Nature* **339**, 31–36 (05 1989).

- 1299 9. Metcalf, P., Blum, M., Freymann, D., Turner, M. & Wiley, D. C. Two variant surface  
 1300 glycoproteins of *Trypanosoma Brucei* of different sequence classes have similar 6 Å  
 1301 resolution X-ray structures. *Nature* **325**, 84–86 (1987).
- 1302 10. Chin, A. F., Zheng, Y. & Hilser, V. J. Phylogenetic convergence of phase separation  
 1303 and mitotic function in the disordered protein BuGZ. *Protein Sci.* **31**, 822–834 (2022).
- 1304 11. Beh, L. Y., Colwell, L. J. & Francis, N. J. A core subunit of Polycomb repressive complex  
 1305 1 is broadly conserved in function but not primary sequence. *Proceedings of the  
 1306 National Academy of Sciences* **109**, E1063–71 (05 2012).
- 1307 12. Schmidt, H. B., Barreau, A. & Rohatgi, R. Phase separation-deficient TDP43 remains  
 1308 functional in splicing. *Nat. Commun.* **10**, 4890 (2019).
- 1309 13. Langstein-Skora, I. *et al.* Sequence- and chemical specificity define the functional  
 1310 landscape of intrinsically disordered regions. *bioRxiv* 2022.02.10.480018 (2022)  
 1311 doi:10.1101/2022.02.10.480018.
- 1312 14. Mindel, V. *et al.* Intrinsically disordered regions of the Msn2 transcription factor  
 1313 encode multiple functions using interwoven sequence grammars. *Nucleic Acids Res.*  
 1314 **52**, 2260–2272 (2024).
- 1315 15. Sigler, P. B. Transcriptional activation. Acid blobs and negative noodles. *Nature* **333**,  
 1316 210–212 (05 1988).
- 1317 16. Staller, M. V. *et al.* A high-throughput mutational scan of an intrinsically disordered  
 1318 acidic transcriptional activation domain. *Cell Syst.* **6**, 444–455.e6 (2018).
- 1319 17. Kumar, M. *et al.* ELM-the Eukaryotic Linear Motif resource-2024 update. *Nucleic Acids  
 1320 Res.* **52**, D442–D455 (2024).
- 1321 18. Sanborn, A. L. *et al.* Simple biochemical features underlie transcriptional activation  
 1322 domain diversity and dynamic, fuzzy binding to Mediator. *Elife* **10**, e68068 (2021).
- 1323 19. Ravarani, C. N. *et al.* High-throughput discovery of functional disordered regions:  
 1324 investigation of transactivation domains. *Mol. Syst. Biol.* **14**, e8190 (2018).
- 1325 20. Broyles, B. K. *et al.* Activation of gene expression by detergent-like protein domains.  
 1326 *iScience* **24**, 103017 (2021).
- 1327 21. Erijman, A. *et al.* A High-Throughput Screen for Transcription Activation Domains  
 1328 Reveals Their Sequence Features and Permits Prediction by Deep Learning. *Mol. Cell*  
 1329 **78**, 890–902.e6 (2020).
- 1330 22. Arnold, C. D. *et al.* A high-throughput method to identify trans-activation domains  
 1331 within transcription factor sequences. *EMBO J.* **37**, e98896 (2018).
- 1332 23. Morffy, N. *et al.* Identification of plant transcriptional activation domains. *Nature* **632**,  
 1333 166–173 (2024).
- 1334 24. Mahatma, S. *et al.* Prediction and functional characterization of transcriptional  
 1335 activation domains. in *2023 57th Annual Conference on Information Sciences and  
 1336 Systems (CISS)* 1–6 (2023).
- 1337 25. Erkina, T. Y. & Erkin, A. M. Nucleosome distortion as a possible mechanism of  
 1338 transcription activation domain function. *Epigenetics Chromatin* **9**, 40 (2016).
- 1339 26. Kotha, S. R. & Staller, M. V. Clusters of acidic and hydrophobic residues can predict  
 1340 acidic transcriptional activation domains from protein sequence. *Genetics* **225**,  
 1341 (2023).
- 1342 27. Udupa, A., Kotha, S. R. & Staller, M. V. Commonly asked questions about  
 1343 transcriptional activation domains. *Curr. Opin. Struct. Biol.* **84**, 102732 (2024).
- 1344 28. Staller, M. V. *et al.* Directed mutational scanning reveals a balance between acidic  
 1345 and hydrophobic residues in strong human activation domains. *Cell Systems* **13**,  
 1346 334–345.e5 (2022).
- 1347 29. Cress, W. D. & Triezenberg, S. J. Critical structural elements of the VP16  
 1348 transcriptional activation domain. *Science* **251**, 87–90 (01 1991).
- 1349 30. Shen, F., Triezenberg, S. J., Hensley, P., Porter, D. & Knutson, J. R. Critical amino acids  
 1350 in the transcriptional activation domain of the herpesvirus protein VP16 are solvent-  
 1351 exposed in highly mobile protein segments. An intrinsic fluorescence study. *J. Biol.  
 1352 Chem.* **271**, 4819–4826 (03 1996).
- 1353 31. DelRosso, N. *et al.* Large-scale mapping and mutagenesis of human transcriptional  
 1354 effector domains. *Nature* (2023) doi:10.1038/s41586-023-05906-y.
- 1355 32. Alerasool, N., Leng, H., Lin, Z.-Y., Gingras, A.-C. & Taipale, M. Identification and  
 1356 functional characterization of transcriptional activators in human cells. *Mol. Cell* **82**,  
 1357 677–695.e7 (2022).

- 1358 33. Kato, S. *et al.* Understanding the function-structure and function-mutation  
 1359 relationships of p53 tumor suppressor protein by high-resolution missense mutation  
 1360 analysis. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 8424– 8429 (07 2003).  
 1361 34. Sadowski, I., Ma, J., Triezenberg, S. & Ptashne, M. GAL4-VP16 is an unusually potent  
 1362 transcriptional activator. *Nature* **335**, 563–564 (10 1988).  
 1363 35. Burz, D. S. & Hanes, S. D. Isolation of Mutations that Disrupt Cooperative DNA  
 1364 Binding by the Drosophila Bicoid Protein☆. *J. Mol. Biol.* **305**, 219–230 (2001).  
 1365 36. Lebrecht, D. *et al.* Bicoid cooperative DNA binding is critical for embryonic patterning  
 1366 in Drosophila. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 13176– 13181 (09 2005).  
 1367 37. Hummel, N. F. C., Markel, K., Stefani, J., Staller, M. V. & Shih, P. M. Systematic  
 1368 identification of transcriptional activation domains from non-transcription factor  
 1369 proteins in plants and yeast. *Cell Syst* (2024) doi:10.1016/j.cels.2024.05.007.  
 1370 38. Hummel, N. F. C. *et al.* The trans-regulatory landscape of gene networks in plants.  
 1371 *Cell Syst* **14**, 501–511.e4 (2023).  
 1372 39. Tsong, A. E., Tuch, B. B., Li, H. & Johnson, A. D. Evolution of alternative transcriptional  
 1373 circuits with identical logic. *Nature* **443**, 415–420 (2006).  
 1374 40. Lynch, M. The evolution of genetic networks by non-adaptive processes. *Nat. Rev.*  
 1375 *Genet.* **8**, 803– 813 (2007).  
 1376 41. Byrne, K. P. & Wolfe, K. H. The Yeast Gene Order Browser: combining curated  
 1377 homology and syntenic context reveals gene fate in polyploid species. *Genome Res.*  
 1378 **15**, 1456–1461 (2005).  
 1379 42. Bennett, R. J. & Turgeon, B. G. Fungal Sex: The Ascomycota. *Microbiol Spectr* **4**,  
 1380 (2016).  
 1381 43. Roesgaard, M. A. *et al.* Deciphering the Alphabet of Disorder-Glu and Asp Act  
 1382 Differently on Local but Not Global Properties. *Biomolecules* **12**, (2022).  
 1383 44. Tuttle, L. M. *et al.* Gcn4-Mediator Specificity Is Mediated by a Large and Dynamic  
 1384 Fuzzy Protein-Protein Complex. *CellReports* **22**, 3251– 3264 (03 2018).  
 1385 45. Brzovic, P. S. *et al.* The acidic transcription activator Gcn4 binds the mediator subunit  
 1386 Gal11/Med15 using a simple protein interface forming a fuzzy complex. *Mol. Cell* **44**,  
 1387 942– 953 (12 2011).  
 1388 46. Scholes, N. S. & Weinzierl, R. O. J. Molecular Dynamics of ‘Fuzzy’ Transcriptional  
 1389 Activator-Coactivator Interactions. *PLoS Comput. Biol.* **12**, e1004935 (2016).  
 1390 47. Lotthammer, J. M., Ginell, G. M., Griffith, D., Emenecker, R. J. & Holehouse, A. S.  
 1391 Direct prediction of intrinsically disordered protein conformational properties from  
 1392 sequence. *Nat. Methods* **21**, 465–476 (2024).  
 1393 48. Shemer, R., Meimoun, A., Holtzman, T. & Kornitzer, D. Regulation of the transcription  
 1394 factor Gcn4 by Pho85 cyclin PCL5. *Mol. Cell. Biol.* **22**, 5395– 5404 (2002).  
 1395 49. Chi, Y. *et al.* Negative regulation of Gcn4 and Msn2 transcription factors by Srb10  
 1396 cyclin-dependent kinase. *Genes Dev.* **15**, 1078– 1092 (05 2001).  
 1397 50. Conti, M. M. *et al.* Phosphosite Scanning reveals a complex phosphorylation code  
 1398 underlying CDK-dependent activation of Hcm1. *Nat. Commun.* **14**, 310 (2023).  
 1399 51. Raj, N. & Attardi, L. D. The Transactivation Domains of the p53 Protein. *Cold Spring*  
 1400 *Harb. Perspect. Med.* **7**, a026047–19 (2017).  
 1401 52. Dyson, H. J. & Wright, P. E. Role of Intrinsic Protein Disorder in the Function and  
 1402 Interactions of the Transcriptional Coactivators CREB-binding Protein (CBP) and p300.  
 1403 *J. Biol. Chem.* **291**, 6714–6722 (2016).  
 1404 53. Ludwig, C. H. *et al.* High-throughput discovery and characterization of viral  
 1405 transcriptional effectors in human cells. *Cell Syst* **14**, 482–500.e8 (2023).  
 1406 54. Piskacek, M., Vasku, A., Hajek, R. & Knight, A. Shared structural features of the  
 1407 9aaTAD family in complex with CBP. *Mol. Biosyst.* **11**, 844– 851 (2015).  
 1408 55. Warfield, L., Tuttle, L. M., Pacheco, D., Klevit, R. E. & Hahn, S. A sequence-specific  
 1409 transcription activator motif and powerful synthetic variants that bind Mediator using  
 1410 a fuzzy protein interface. *Proceedings of the National Academy of Sciences* **111**,  
 1411 E3506– E3513 (08 2014).  
 1412 56. Pacheco, D. *et al.* Transcription activation domains of the yeast factors Met4 and  
 1413 Ino2: tandem activation domains with properties similar to the yeast Gcn4 activator.  
 1414 *Mol. Cell. Biol.* MCB.00038–18 – 39 (03 2018).  
 1415 57. Tuttle, L. M. *et al.* Mediator subunit Med15 dictates the conserved ‘fuzzy’ binding  
 1416 mechanism of yeast transcription activators Gal4 and Gcn4. *Nat. Commun.* **12**, 1–11

- 1417 (2021).
- 1418 58. Schuler, B. *et al.* Binding without folding - the biomolecular function of disordered  
1419 polyelectrolyte complexes. *Curr. Opin. Struct. Biol.* **60**, 66-76 (2020).
- 1420 59. Dunker, A. K., Bondos, S. E., Huang, F. & Oldfield, C. J. Intrinsically disordered  
1421 proteins and multicellular organisms. *Semin. Cell Dev. Biol.* **37**, 44-55 (2015).
- 1422 60. Tenthorey, J. L., Young, C., Sodeinde, A., Emerman, M. & Malik, H. S. Mutational  
1423 resilience of antiviral restriction favors primate TRIM5 $\alpha$  in host-virus evolutionary  
1424 arms races. *Elife* **9**, (2020).
- 1425 61. Koonin, E. V. & Dolja, V. V. A virocentric perspective on the evolution of life. *Curr.*  
1426 *Opin. Virol.* **3**, 546-557 (2013).
- 1427 62. Dalal, C. K. & Johnson, A. D. How transcription circuits explore alternative  
1428 architectures while maintaining overall circuit output. *Genes Dev.* **31**, 1397-1405  
1429 (2017).
- 1430 63. Fowler, K. R., Leon, F. & Johnson, A. D. Ancient transcriptional regulators can easily  
1431 evolve new pair-wise cooperativity. *Proc. Natl. Acad. Sci. U. S. A.* **120**, e2302445120  
1432 (2023).
- 1433 64. Liu, Y. *et al.* Evolution of the activation domain in a Hox transcription factor. *Int. J.*  
1434 *Dev. Biol.* **62**, 745- 753 (2018).
- 1435 65. Zarin, T. *et al.* Proteome-wide signatures of function in highly diverged intrinsically  
1436 disordered regions. *eLife* **xx**, xxx-45 (03 2019).
- 1437 66. Zarin, T. *et al.* Identifying molecular features that are associated with biological  
1438 function of intrinsically disordered protein regions. *Elife* **10**, e60220 (2021).
- 1439 67. Zarin, T., Tsai, C. N., Ba, A. N. N. & Moses, A. M. Selection maintains signaling  
1440 function of a highly diverged intrinsically disordered region. *Proc. Natl. Acad. Sci. U.*  
1441 *S. A.* **114**, E1450-E1459 (2017).
- 1442 68. Parker, M. W. *et al.* A new class of disordered elements controls DNA replication  
1443 through initiator self-assembly. *Elife* **8**, e48562 (2019).
- 1444 69. Parker, M. W., Kao, J. A., Huang, A., Berger, J. M. & Botchan, M. R. Molecular  
1445 determinants of phase separation for Drosophila DNA replication licensing factors.  
1446 *Elife* **10**, (2021).
- 1447 70. Davey, N. E., Cyert, M. S. & Moses, A. M. Short linear motifs - ex nihilo evolution of  
1448 protein regulation. *Cell Commun. Signal.* **13**, 43 (2015).
- 1449 71. Wong, E. S. *et al.* Deep conservation of the enhancer regulatory code in animals.  
1450 *Science* **370**, (2020).
- 1451 72. Ludwig, M. Z., Bergman, C., Patel, N. H. & Kreitman, M. Evidence for stabilizing  
1452 selection in a eukaryotic enhancer element. *Nature* **403**, 564-567 (2000).
- 1453 73. Ludwig, M. Z. *et al.* Functional evolution of a cis-regulatory module. *PLoS Biol.* **3**, e93  
1454 (2005).
- 1455 74. Hare, E. E., Peterson, B. K., Iyer, V. N., Meier, R. & Eisen, M. B. Sepsid even-skipped  
1456 Enhancers Are Functionally Conserved in Drosophila Despite Lack of Sequence  
1457 Conservation. *PLoS Genet.* **4**, e1000106 (2008).
- 1458 75. Peterson, B. K. *et al.* Big genomes facilitate the comparative identification of  
1459 regulatory elements. *PLoS One* **4**, e4688 (2009).
- 1460 76. Kaplow, I. M. *et al.* Relating enhancer genetic variation across mammals to complex  
1461 phenotypes using machine learning. *Science* **380**, eabm7993 (2023).
- 1462 77. Arnosti, D. N. & Kulkarni, M. M. Transcriptional enhancers: Intelligent enhanceosomes  
1463 or flexible billboards? *J. Cell. Biochem.* **94**, 890- 898 (2005).
- 1464 78. Jackson, B. M., Drysdale, C. M., Natarajan, K. & Hinnebusch, A. G. Identification of  
1465 seven hydrophobic clusters in GCN4 making redundant contributions to  
1466 transcriptional activation. *Mol. Cell. Biol.* **16**, 5557-5571 (1996).
- 1467 79. Ginell, G. M., Emenecker, R. J., Lotthammer, J. M., Usher, E. T. & Holehouse, A. S.  
1468 Direct prediction of intermolecular interactions driven by disordered regions.  
1469 *bioRxiv* 2024.06.03.597104 (2024).
- 1470 80. Amberg, D. C., Burke, D. & Strathern, J. N. *Methods in Yeast Genetics: A Cold Spring*  
1471 *Harbor Laboratory Course Manual.* (CSHL Press, 2005).
- 1472 81. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and  
1473 powerful approach to multiple testing. *J. R. Stat. Soc. Series B Stat. Methodol.* **57**,  
1474 289-300 (1995).
- 1475 82. Dunn, O. J. Multiple Comparisons among Means. *J. Am. Stat. Assoc.* **56**, 52-64 (1961).

- 1476 83. Das, R. K. & Pappu, R. V. Conformations of intrinsically disordered proteins are  
1477 influenced by linear sequence distributions of oppositely charged residues.  
1478 *Proceedings of the National Academy of Sciences* **110**, 13392–13397 (08 2013).  
1479 84. Ginell, G. M. & Holehouse, A. S. Intrinsically Disordered Proteins, Methods and  
1480 Protocols. *Methods Mol. Biol.* **2141**, 103–126 (2020).  
1481 85. Martin, E. W. *et al.* Sequence determinants of the conformational properties of an  
1482 intrinsically disordered protein prior to and upon multisite phosphorylation. *J. Am.*  
1483 *Chem. Soc.* **138**, 15323–15335 (2016).  
1484 86. Emenecker, R. J., Griffith, D. & Holehouse, A. S. Metapredict V2: An update to  
1485 metapredict, a fast, accurate, and easy-to-use predictor of consensus disorder and  
1486 structure. *bioRxiv* 2022.06.06.494887 (2022) doi:10.1101/2022.06.06.494887.  
1487 87. Hope, I. A., Mahadevan, S. & Struhl, K. Structural and functional characterization of  
1488 the short acidic transcriptional activation region of yeast GCN4 protein. **333**, 635–  
1489 640 (06 1988).  
1490 88. Drysdale, C. M. *et al.* The transcriptional activator GCN4 contains multiple activation  
1491 domains that are critically dependent on hydrophobic amino acids. *Mol. Cell. Biol.* **15**,  
1492 1220–1233 (1995).