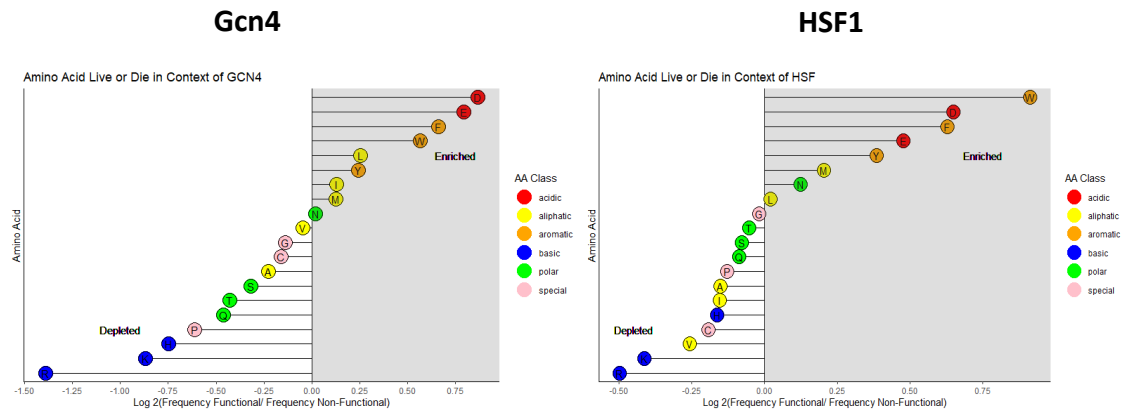# Supplemental information
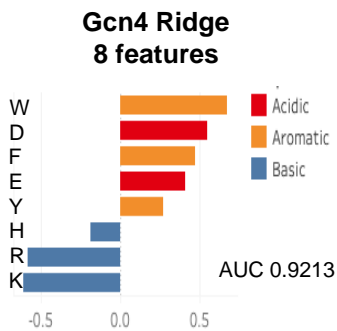
# Activation of gene expression by detergent-like

# protein domains

Bradley K. Broyles, Andrew T. Gutierrez, Theodore P. Maris, Daniel A. Coil, Thomas M. Wagner, Xiao Wang, Daisuke Kihara, Caleb A. Class, and Alexandre M. Erkine
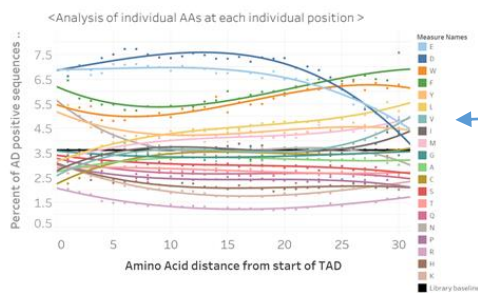
**Figure. S1. The profile of frequencies for individual amino acids is looking strikingly similar for both Gcn4 and HSF1 contexts. Related to Figure 2.** X-axis: Log2 of frequency in functional to frequency in non-functional sub-pools. Y-axis: individual amino acids.
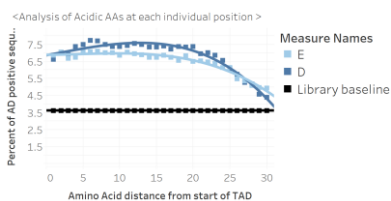
**Figure S2. ML coefficients of eight features determining the precision of the Ridge regression model. Related to figure 2.** Y-axis: amino acids as ML features. X-axis: ML feature coefficients values.
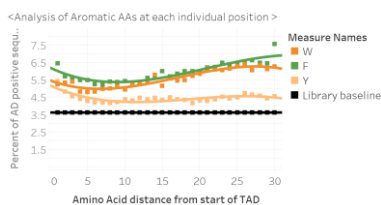
## All individual amino acids



Use group averages
(like panel C)

## Acidic (D, E)



## Aromatic (W, F, Y)



## Basic (K, R, H)



## Aliphatic (L, V, I M, A)



## Polar not acidic or basic (C, S, T, Q, N)



## Special cases (N, P, C, G)



**Figure S3. Probability of contribution to functionality for individual amino acids as a function of position in tAD. Related to Figure 6.** Data for indicated amino acids and amino acid groups are shown. Y-axis: Percent of functional sequences within the pool. X-axis: position of the amino acid with the 30 amino acid stretch tested (C-terminus is position 30).

PRRHKRCSFVTDYDPCFGNFWQSDCYMAIV
KVRVNQQGLNATVEDVVADELWAMGLFELYI
VKKCSMESERLDVVETVHEWLESINQVYIL

ADTDYPGFMLADDFDYIDDLFFEAWFFLEL
EDVDDLCFIADGGEGLEFDAIFDEFIMFEL
EYFVLESNYDDAEEFFVDEDWHFATGDIYE

Ridge Regression

|  |  | FN | TP |
|---|---|---|---|
| LSTM Network | TP | 395 | 6677 |
|  | FN | 368 | 145 |

VSCGTSSCTTRCRRRRAVALALGRRRRDNM
VRSKLGASGFEPPNSDRFPAKQSTLGTLNK
NSNRDMVHERYTVVDNRERVTVERTVNGET

VIAFGEDVVEDFAEEFYLWMLLVDIAYAET
SGALWFAMKVSVFMVEYWFSHETLWWEASA
DEVCEGTVTSETLFTDDGEYEYDFRADITV
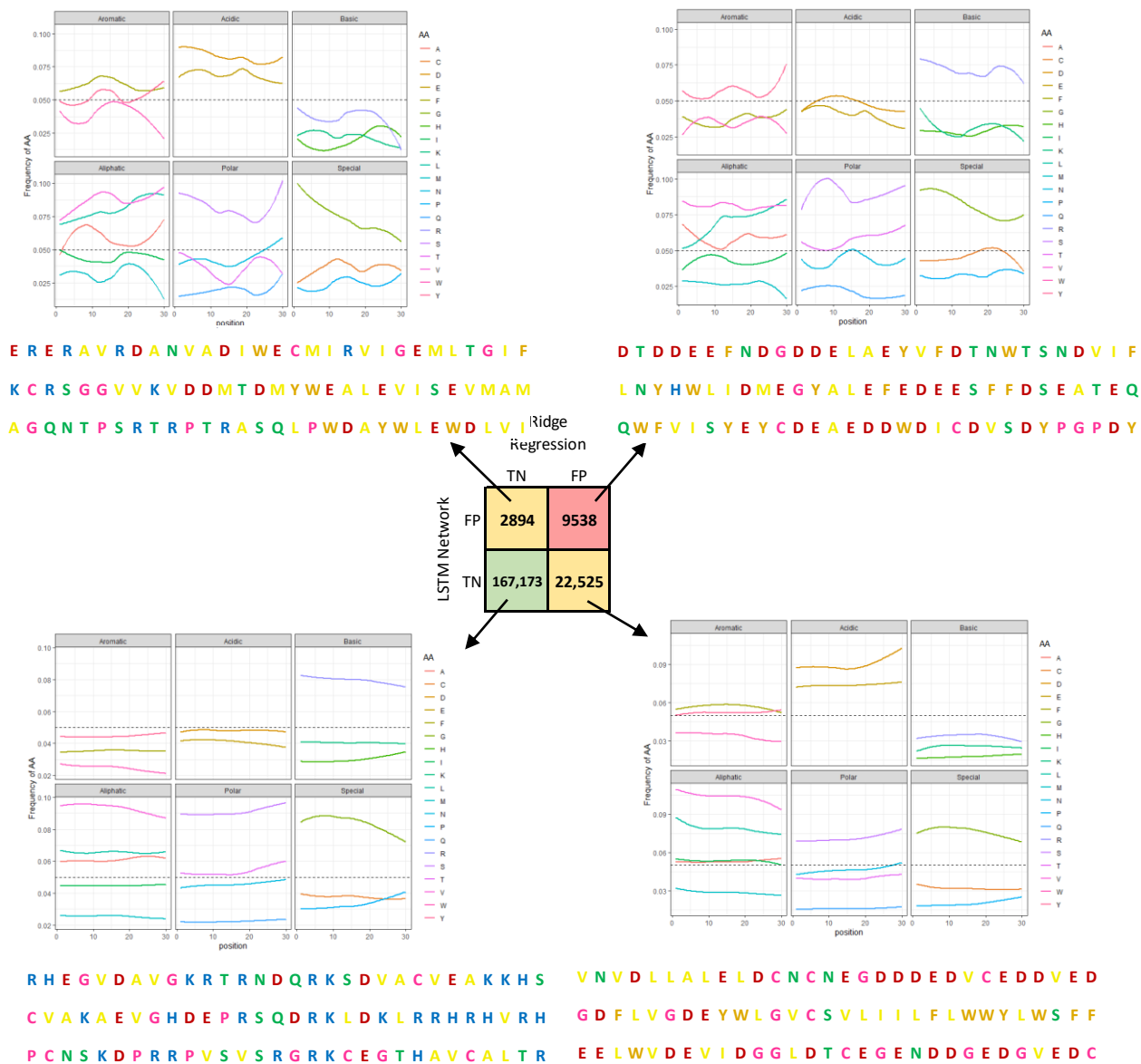
**Figure S4. Comparisons of Ridge Regression and LSTM Network predictions for True Functional sequences. Related to Figure 7.** Matrix in center: Number of true functional sequences identified correctly (TP = True Positive) or incorrectly (FN = False Negative) by LSTM network and ridge regression. A consensus plot is provided for the sequences falling in each of the four boxes, along with three example sequences. LSTM network analysis performs better in identifying functional sequences, by identifying more complicated factors (such as locations of relevant amino acids) that contribute to functionality. Sequences mis-identified by both methods could demonstrate unique patterns of functional sequences, or simply experimental error.

ERERAVRDANVADIWECMIRVIGEMLTGIF
KCRSGGVVKVDDMTDMYWEALEVISEVMAM
AGQNTPSRTRPTRASQLPWDAYWLEWDLVI

DTDDEEFNDGDDELAEYVFDTNWTSNDVIF
LNYHWLIDMEGYALEFEDEESFFDSEATEQ
QWFVISYEYCDEAEDDWDICDVSDYPGPDY

**Ridge Regression**

|  | TN | FP |
|---|---|---|
| FP | 2894 | 9538 |
| TN | 167,173 | 22,525 |

LSTM Network

RHEGVDAVGKRTRNDQRKSDVACVEAKKHS
CVAKAEVGHDEPRSQDRKLDKLRRHRHVRH
PCNSKDPRRPVSVSRGRKCEGTHAVCALTR

VNVDLLALELDCNCNEGDDDEDVCEDDVED
GDFLVGDEYWLGVCSVLIILFLWWYLWSFF
EELWVDEVIDGGLDTCEGENDDGEDGVEDC

**Figure S5. Comparisons of Ridge Regression and LSTM Network predictions for True Non-Functional sequences. Related to Figure 7.** Matrix in center: Number of true non-functional sequences identified correctly (TN = True Negative) or incorrectly (FP = False Positive) by LSTM network and ridge regression. A consensus plot is provided for the sequences falling in each of the four boxes, along with three example sequences. LSTM network analysis performs better in identifying non-functional sequences, by identifying more complicated factors (such as combinations of aromatic & acidic amino acids) that contribute to functionality.