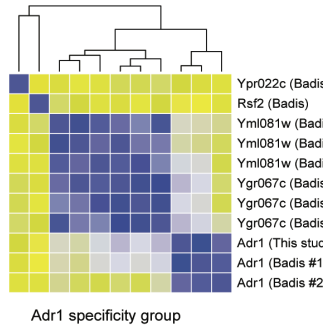


A

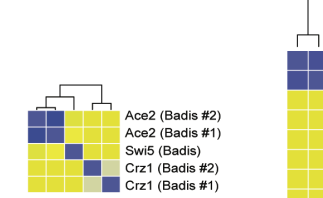
Canonical Residues	Specificity Group Members	Binding Scheme	
		ZF2	ZF1
ZF1 : ZF2		6 3 2 -1 6 3 2 -1	
<b>REHR:RDNQ</b>	MSN2, MSN4, COM2 USV1, RGM1	Q N D R R R H E R A G G G G G	
<b>TGHR:RDNQ</b>	NRG1, NRG2	Q N D R R R H G T A G G G T C C	
<b>SHHR:RDNQ</b>	GIS1, RPH1	Q N D R R R H H S A G G G G	
<b>RSWR:RDNQ</b>	MOT3	Q N D R R R W S T A G G C A C	
<b>REHR:RDLR</b>	RSF2, YML081W, ADR1 YGR067C, YPR022c	R L D R R R H E R G c G G G G	
<b>REHR:RDER</b>	MIG1, MIG2, MIG3	R E D R R R H E R G c G G G G	
<b>RSDR:RDAR</b>	MET31, MET32	R A D R R D S R G T G C G	
<b>IGYR:RDTT</b>	STP1, STP2	T T D R R Y G I g C G C g	
<b>RNDR:RDTN</b>	STP3, STP4	N T D R R D N R g c G G C T g	

B



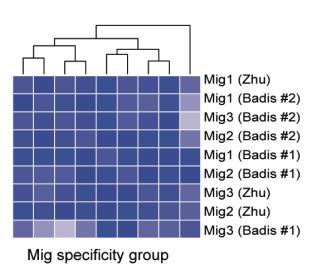
Adr1 specificity group

D

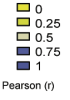


Swi5 specificity group

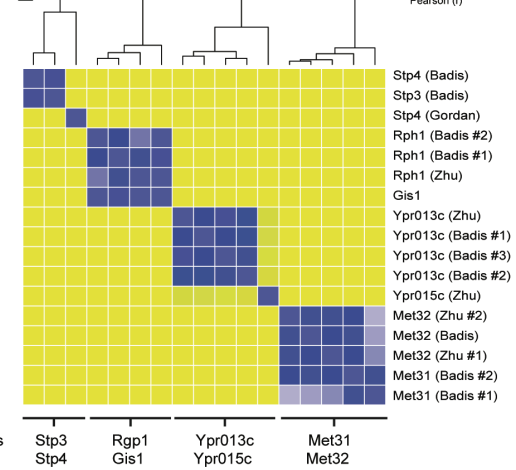
C



Mig specificity group

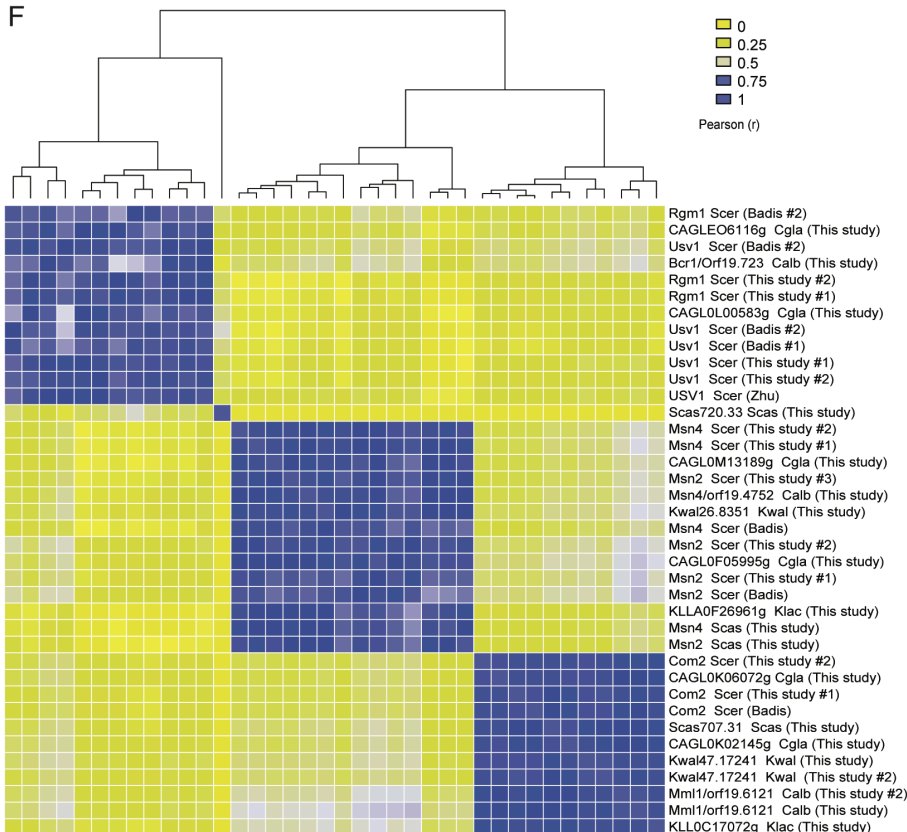


E

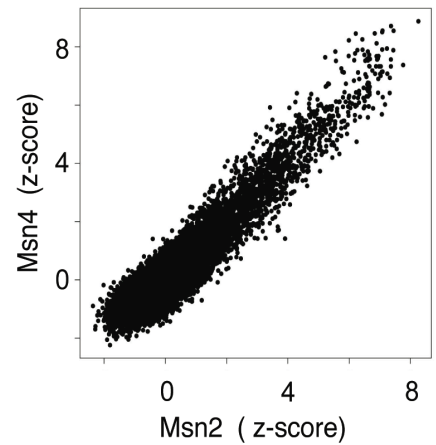


Specificity groups

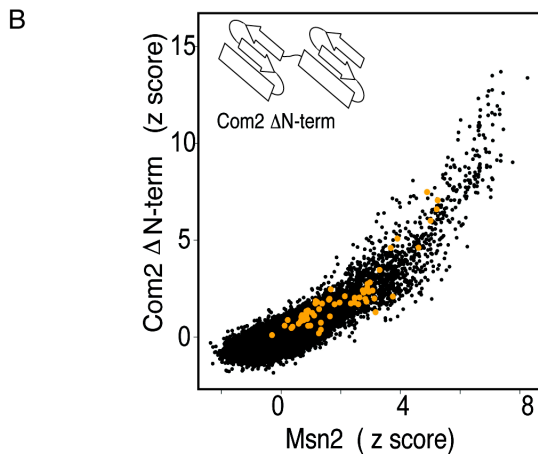
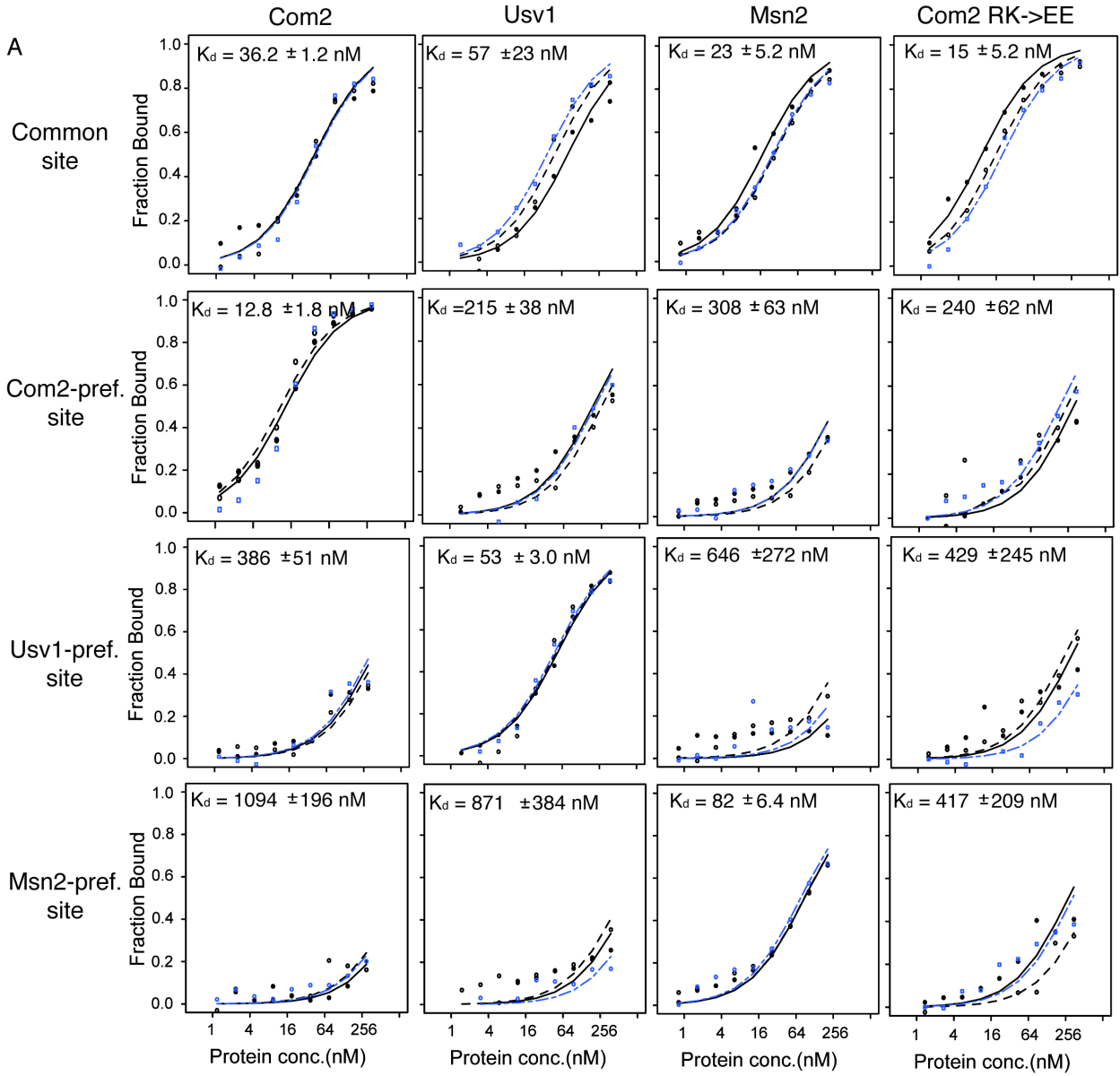
F



G

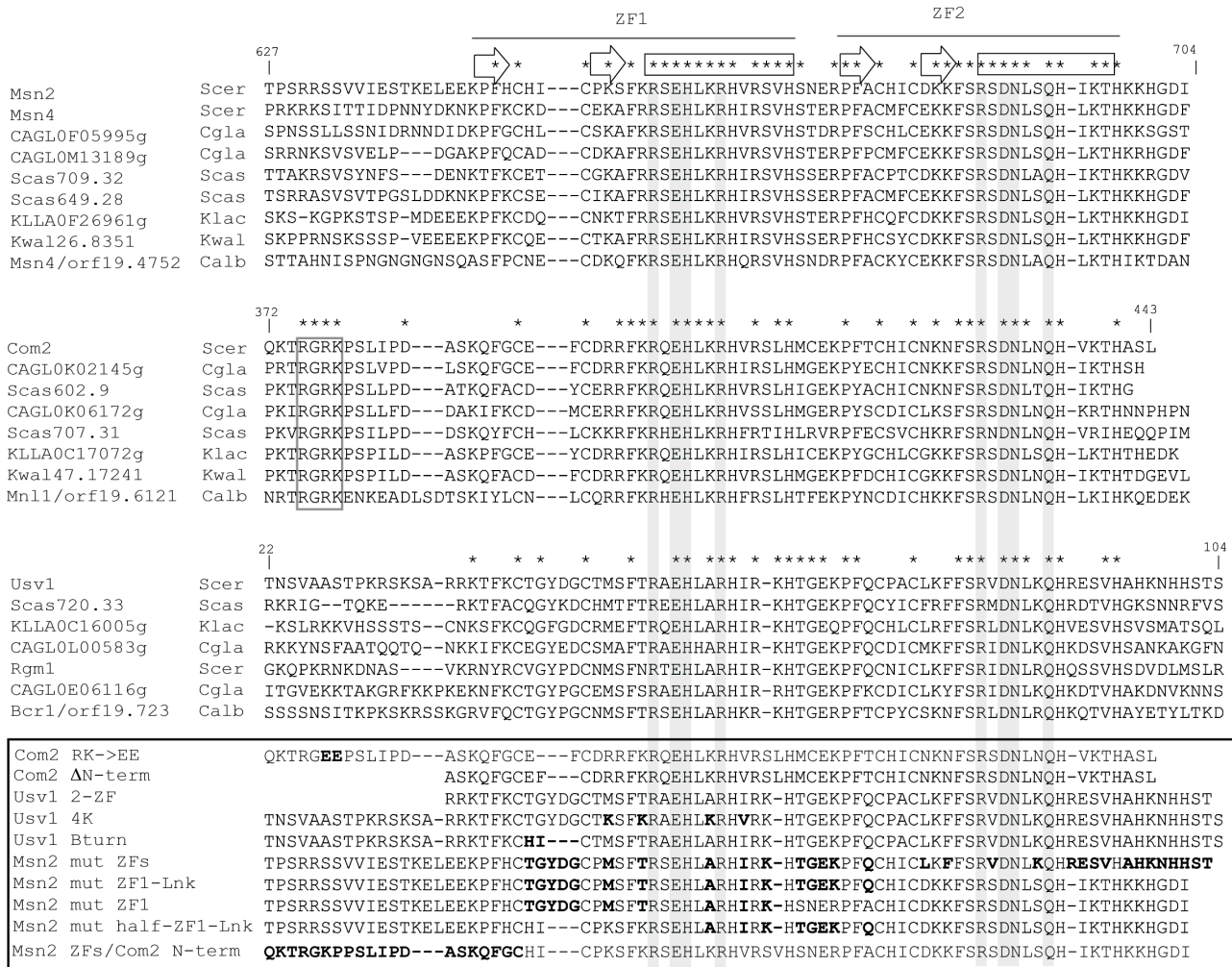


**Figure S1. DNA-binding Scheme and Profiles for Yeast C2H2 ZF proteins** (Related to Figure 1). **(A)** Shown is an amino acid-base interaction map proposed for the nine C2H2 ZF specificity groups identified in Table S1. Specificity groups are defined by identical amino acids at the eight canonical recognition residues in the two ZF domains (four residues in each ZF). ZF orientation and proposed amino acid-base interactions are based on (i) previously described binding schemes proposed for Adr1 (Schaufler and Klevit, 2003), (ii) published binding logos for the different families (Badis et al., 2008; Zhu et al., 2009), (iii) known recognition propensities for C2H2 ZF domains (Wolfe et al., 2000) (Elrod-Erickson et al., 1996) (Miller and Pabo, 2001), and (iv) recognized patterns between amino acids in the recognition positions and preferred bases at stereotyped positions. **(B-E)** Pairwise comparison of the DNA-binding profiles proteins from the ‘Specificity Groups’ defined in Table S1. Comparisons for groups with more than two member proteins are shown in **(B)**, **(C)** and **(D)**. Comparisons for groups with two member proteins are shown together in **(E)**. Pairwise binding similarity was assessed as Pearson correlation of binding z scores to the highest affinity 500 sites in each PBM experiment. Comparisons were performed for published datasets (Badis(Badis et al., 2008), Zhu(Zhu et al., 2009), or Gordân (Gordan et al., 2011)). **(F)** DNA-binding similarity of Msn2-family orthologs across six Ascomycota species. Pairwise comparison of the DNA-binding profiles for *S. cerevisiae* Msn2-family proteins (as in Figure 1A) and eighteen orthologs from *S. glabrata*, *S. castellii*, *K. lactis*, *K. waltii* and *C. albicans*. Binding comparison and hierarchical clustering are as described for Figure 1A. Clone IDs from species other than *S. cerevisiae* are from Fungal Orthogroups v1.1 database ([www.broadinstitute/regev/orthogroups](http://www.broadinstitute/regev/orthogroups)) (Wapinski et al., 2007). **(G)** Pairwise binding profile comparison for Msn2 and Msn4 (details as in Figure 1).



**Figure S2. EMSA-derived Saturation Binding Curves.** (Related to Figures 1 and 2)

**(A)** Saturation binding curves for Com2, Usv1, Msn2 and Com2 RK→EE binding to the common, Com2-preferred, Usv1-preferred and Msn2-preferred sites (see Extended Experimental Procedure). EMSA binding experiments were performed in triplicate or duplicate (as shown). Dissociation binding constants ( $K_d$ ) were determined for each replicate experiment; mean and standard deviation over the replicates are shown. **(B)** Pairwise binding profile comparisons (as in Figure 1) for Com2  $\Delta$ N-term mutant relative to Msn2. Com2-preferred sites (as in Figure 1B) are highlighted (orange).



Mutant Constructs

**Figure S3. Sequence Alignment of ZF Regions from Msn2-family Orthologs and Mutant Constructs** (Related to Figures 1 and 2). Multiple protein sequence alignment spanning the ZF domains for Msn2-family orthologs from six Ascomycota species, and mutant *S. cerevisiae* constructs (boxed). Clone IDs are from Fungal Orthogroups v1.1 database ([www.broadinstitute/regev/orthogroups](http://www.broadinstitute/regev/orthogroups)) (Wapinski et al., 2007) (see also Supplementary Table 1). Beta-strand and alpha-helix secondary structure elements for each ZF domain are indicated at the top (arrow and box, respectively); residues conserved across all Msn2-family orthologs are also indicated. Canonical ZF DNA-contacting residues, identical in all Msn2-family members (by definition), are highlighted with grey bars. A conserved RGRK motif N-terminal to ZF1 in all Com2 orthologs is indicated with a grey box. For mutant constructs amino acid changes are highlighted in bold.

**Table S1. Specificity Groups of 2-ZF proteins** (Related to Figure 1).

Canonical DNA Recognition Residues (ZF1: ZF2)	Number of Members	Specificity Group Members
REHR : RDLR	5	Rsf2   Yml081w    Adr1    Ygr067c    Ypr022c
REHR : RDNQ	5	Msn2l Msn4    UsvlRgm1    Com2
REHR : RDER	3	Mig2   Mig3    Mig1
RYNS : RHDR	3	Swi5   Ace2    Crz1
RDTT : VSNR	2	Ypr013c    Ypr015c
RNDR : RDAR	2	Met31   Met32
RNDR : RDTN	2	Stp3   Stp4
SHHR : RDHQ	2	Gis1   Rph1
IGYR : RDTT	2	Stp1   Stp2
TGHR : RDNQ	2	Nrg1   Nrg2

Listed are the twenty-eight C2H2 ZF proteins from *S. cerevisiae* that have 2 adjacent ZF domains. Proteins are organized into ‘Specificity Groups’ based on the identity of the 8 canonical ZF DNA-contacting residues. The canonical residues 4 from each ZF domain separated by a colon for each group are shown in column 1. Columns 2 and 3 list the number of proteins in each group, and their common names or SGD identifiers, respectively. Paralogs generated by the WGD event in the *S. cerevisiae* lineage (Wapinski et al., 2007) are shown separated by a ‘|’.

**Table S2. ZF construct sequences** (Related to Experimental Procedures). Construct name (column 1), source gene (column 2), yeast species (column 3) and construct sequence (column 4) are shown for all *S. cerevisiae* constructs/mutants (A) and all related species (B). Clone IDs for non-*S. cerevisiae* species are from Fungal Orthogroups v1.1 database ([www.broadinstitute.org/regev/orthogroups](http://www.broadinstitute.org/regev/orthogroups)) (Wapinski et al., 2007). Homology relation to *S. cerevisiae* ortholog is indicated in column 3.

**Table S3. PBM 8-mer data** (Related to Experimental Procedures). The universal PBM 8-mer median fluorescence intensity values are provided for all PBM experiments discussed in the paper.

**Table S4. PBM raw probe data** (Related to Experimental Procedures). The raw PBM probe fluorescence intensity values are provided for all PBM experiments discussed in the paper.

## Supplemental Experimental Procedures

**Cloning and Preparation of Protein Samples.** Open reading frames (see Supplementary Table S2) were cloned into Gateway pDEST15 (N-terminal GST-tag) expression vectors (Invitrogen). Clones for the species *C. glabrata*, *S. castellii*, *K. lactis*, *K. waltii* and *C. albicans* were made by PCR amplification from cDNA libraries (generously provided by I. Wapinski, A. Regev). Mutant *S. cerevisiae* clones were made by gene synthesis (Integrated DNA technologies (IDT)). Protein samples for PBM experiments were produced by *in vitro* transcription and translation (IVT) using the PURExpress kit (New England BioLabs) from purified plasmids. Protein concentrations were quantified by Western blot using GST standards (Zhu et al., 2009). Protein samples for EMSA binding assays were expressed in bacteria and purified using affinity column. Briefly, proteins were expressed in *E. coli* BL21 (DE3) pLYS cells. Bacterial cultures were grown in 200mL of Terrific Broth (TB) medium (in 2-L Erlenmeyer Flasks) supplemented with 75 µg/mL carbenicillin (Sigma-Aldrich), 34 µg/mL chloramphenicol (Sigma), and 2 mM MgSO<sub>4</sub> and were incubated at 30°C with shaking at 250 rpm. Protein expression was induced with 0.5 mM IPTG (FisherScientific) after cultures achieved an OD<sub>600</sub> of ~0.55, and bacterial cultures were incubated at the above conditions for an additional 3 hours before being harvested via centrifugation (Sorvall GSA rotor; 5,000 rpm; 4°C, for 15min). Cell pellets were stored in an -80°C freezer until further use. Bacterial cells were lysed GST via a French press. Bacterial cell pellets were resuspended in 25 mL of Wash Buffer A (1x PBS, pH 7.4, 0.02% Triton-X-100 (Acros Organics), 2 mM DTT (Fisher), 1 mM PMSF (Sigma), and 50 µM zinc acetate (Sigma); pH 7.42) at room temperature. Resuspended cells were lysed in a French press homogenizer at ~17,000 psi. The crude bacterial extract was clarified of cellular debris by centrifugation (Sorvall SS-34 rotor; 12,500 rpm; 4°C, for 45 min). The sample was

further clarified by passing it through a 0.45  $\mu\text{m}$  nylon filter (Fisher). The final sample volume was increased to 50 mL using Wash Buffer A. Proteins were purified from the crude bacterial extract using an FPLC and GSTrap FF (GE Healthcare) column. Briefly, the 1 mL GST column was equilibrated in 5 column volumes (CVs) of Wash Buffer A (see above) @ 0.2 mL/min (used throughout). The entire sample (50 mL) was loaded onto the column, and unbound protein was washed away using 15 CVs of Wash Buffer A. The desired protein was eluted using 5 CVs of Elution Buffer B (50 mM Tris [pH 8.0] (Fisher), 0.02% Triton-X-100, 2 mM DTT, 1 mM PMSF, 15 mM reduced L-Glutathione (Sigma), and 50  $\mu\text{M}$  zinc acetate; pH 7.22 @ Room temperature), with 500  $\mu\text{L}$  fractions taken. Fractions were analyzed for protein yield and purity using a Bradford Assay kit (ThermoScientific) and SDS-PAGE followed by Coomassie staining. Desired fractions were pooled, flash-frozen, and stored at  $-80^{\circ}\text{C}$  until further use.

**Electromobility Shift Assays (EMSA) for  $K_d$  Determination**

**Electromobility Shift Assays (EMSA) for  $K_d$  Determination:** Double-stranded (ds) DNA oligonucleotides (60 bp) were generated by primer extension (45  $\mu\text{L}$  total) were prepared with 8  $\mu\text{M}$  60-bp ssDNA template strand (sequences provided below), 8  $\mu\text{M}$  24-bp extension primer (5'-CGCGTCGCAC CCTACCTTTC GTTA), 1.6 mM dNTPs (New England BioLabs), and 1x ThermoPol Reaction Buffer (NEB). Separate enzyme mixtures (5  $\mu\text{L}$  total) were prepared with 4 units *Bst* DNA polymerase, Large fragment (NEB) and 1x ThermoPol Reaction Buffer. Reaction mixtures were heated in a thermocycler to  $95^{\circ}\text{C}$ , and gradually cooled to  $63^{\circ}\text{C}$ , at a rate of  $-0.1^{\circ}\text{C}/\text{s}$ . The enzyme mix was placed in the thermocycler for 1 min to equilibrate, and then 5  $\mu\text{L}$  of this mix was pipetted into each 45  $\mu\text{L}$  reaction mixture. The resulting extension reactions were incubated at  $63^{\circ}\text{C}$  for an additional 90 min. The double-stranded DNA was purified from the extension reaction using a MinElute PCR Purification Kit (Qiagen), as per the manual. DNA probes were radioactively labeled using T4 polynucleotide kinase (PNK) (NEB) and ATP [ $\gamma$   $^{32}\text{P}$ ] (Perkin Elmer). Briefly, 2 pmol of oligonucleotide was incubated with 10 units of T4 PNK and 20  $\mu\text{Ci}$  ATP [ $\gamma$   $^{32}\text{P}$ ] in 1x PNK Buffer (15  $\mu\text{L}$  total reaction) @  $37^{\circ}\text{C}$  for 2 hr. The labeled probe was purified from the labeling reaction using a QIAquick Nucleotide



Removal Kit (Qiagen), as per the manual. DNA probes used in EMSA reactions: Com2-preferred (5'- GATAAGCGCC AAATAGGAGA CCACAGTTCA CGTAGTTAAC GAAAGGTAGC GTGCGACGCG), Usv1-preferred (5'-GATAAGCGCC ATTCAGGTAC CCACAGTTCA CGTAGTTAAC GAAAGGTAGC GTGCGACGCG), Msn2-preferred (5'- GATAAGCGCC AAACGGGGT CCACAGTTCA CGTAGTTAAC GAAAGGTAGC GTGCGACGCG). common (5'- GATAAGCGCC ATTCAGGGGT CCACAGTTCA CGTAGTTAAC GAAAGGTAGC GTGCGACGCG).

The dissociation constants ( $K_d$ ) were determined by EMSA. DNA-binding reactions were performed with 0.5 nM radioactively labeled DNA and 1x EMSA Binding Buffer (1x PBS [pH 7.4], 0.02% Triton-X-100, 1 mM dithiothreitol (DTT), 0.2 mg/mL bovine serum albumin (BSA) (NEB), 5% v/v glycerol (Sigma), and 50  $\mu$ M zinc acetate). Protein samples were thawed rapidly and used for serial dilutions (dilution factor =  $\frac{1}{2}$ ), on ice. 2  $\mu$ L of these protein dilutions were added to their respective DNA-binding reactions (samples were used within 30 min of being diluted). The DNA-binding reactions were incubated at room temperature for 30 min. Samples were electrophoresed for 1.7 hr at 70V on prerun, non-denaturing polyacrylamide gels (6% [29:1] acryl-bisacrylamide [Fisher], 0.5x TB [pH 8.3]; 1.5 mm W x 8.25 cm H x 10 cm L). Gels were dried for 1.5 hr at 80°C with suction, using a BioRad Gel Dryer. Dried gel films were placed onto a PhosphorImager (GE Healthcare) and left to expose overnight (~12 hours).

DNA bands were visualized using a Typhoon Trio scanner (GE Healthcare) with a 100  $\mu$ m pixel size. Resulting .gel files were analyzed using Image J, where peak integration of lane intensity histograms was performed to determine Free DNA band intensity. The fraction of bound DNA ( $[DNA_{bound}]/[DNA_{tot}]$ ) was estimated as 1 minus the fraction of unbound probe in a lane, relative to lane in which no protein was added ( $1 - [DNA_{unbound}]/[DNA_{free}]$ ). The R Statistical Package was used to generate binding plots and calculate  $K_d$ , using a non-linear least-squares fit.

**Protein Binding Microarray (PBM) Experiments and Analysis.** PBM experiments were performed using custom-designed, universal 'all-10mer' microarrays (Agilent Technologies Inc., AMADID #016060, 4x44K array format (Zhu et al., 2009)) described previously (Berger and Bulyk, 2006). Microarrays were converted to double-stranded

DNA arrays by primer extension and used in PBM experiments as described previously (Berger and Bulyk, 2006, 2009). Protein samples were incubated on the microarrays at ~200 nM for 1 h in binding buffer (PBS, pH 7.4; 0.2 µg/µl bovine serum albumin (BSA) (New England BioLabs #B9001S); 0.3 ng/µl salmon testes DNA (Sigma, #D7656); 2% non-fat dry milk (Stop & Shop brand); 0.02% Triton X-100; 3 mM dithiothreitol (DTT); 0.02% Triton-X-100; 50 µM zinc acetate dihydrate). Protein-bound arrays were then washed and incubated with antibody for 20 min (0.05 mg/ml Alexa 488-conjugated anti-GST antibody (Invitrogen, #A11131, 2mg/ml); PBS, pH 7.4; 2% non-fat dried milk (Stop & Shop brand), 50 µM Zinc acetate dihydrate). PBM wash protocol is as previously described (Berger and Bulyk, 2009). Microarray scanning, quantification, and data normalization were performed using GenePix Pro ver. 6 (Axon) and masliner (MicroArray LINEar Regression) (Dudley et al., 2002) software as previously described (Berger and Bulyk, 2006, 2009). Full PBM probe intensities (Table S4) and derived 8-mer median intensities (Table S3) are provided.

### **Comparing PBM binding profiles.**

Data normalization and binding profile correlation. For each of the ~32,000 8-mer sequences the median intensity fluorescence values were determined by averaging over all universal PBM probes containing that 8-mer, as previously described (Berger et al., 2008; Berger et al., 2006). Natural log of the median signal intensity values for all 8-mers were then transformed into z-scores ( $z = (x - m)/s$ , m: distribution mean, s: distribution standard deviation). The pair-wise DNA-binding similarity between two proteins was quantified as the Pearson correlation coefficient calculated using the z-scores for the 500 top-scoring 8-mers from each PBM experiment (500 top-scoring 8-mers were selected independently from each experiment then examined in aggregate). Pearson correlation was calculated using the R software package (function: cor; method:pearson).

Hierarchical clustering of PBM profile comparisons (i.e., matrix of pair-wise Pearson correlation scores) was used to determine specificity clusters for the ZF proteins (as in Figure 1). Hierarchical clustering and visualization performed with the R statistical software packing, using the heatmap function in R, with a 'euclidean' distance function and a 'complete' clustering function. *TF-preferred* k-mers were identified as 8-mers

bound significantly better in one PBM experiment than in another. The approach to determine preferential binding is as follows. The PBM probe sequences for each experiment were independently ranked based on their fluorescence intensity values. On the universal PBM design, each 8-mer is present on ~32 different PBM probe sequences (16 probes for palindromic 8-mers), as previously described in detail (Berger and Bulyk, 2006). For each 8-mer the corresponding probe rank values were determined for the two experiments:  $P_i^{\text{Exp1}}$  and  $P_i^{\text{Exp2}}$  (e.g., for the 8-mer AATAGGGG two lists containing the 32 probe rank values from each experiment were determined). Rank values were then compared between these two lists to determine if the ranks were statistically higher on one experiment or another. In practice, the lists  $P_i^{\text{Exp1}}$  and  $P_i^{\text{Exp2}}$ , were combined (maintaining their order) and re-ranked 1 through 64. Using these *normalized* rank values, the enrichment of the rank values for one experiment over the other was determined using the two-sample, unpaired Wilcoxon test (a.k.a. Mann-Whitney test), and p-values determined. Calculations were performed using the R statistical software package (`wilcox.test`). Comparisons were restricted to all strongly bound 8-mers that scored significantly in at least one experiment (8-mer PBM E-score > 0.4); using PBM E-scores as described previously (Berger et al., 2008; Berger et al., 2006) and represents a metric for 8-mer binding significance in a universal PBM experiment. After the statistical enrichment (i.e., TF preference) was determined for all strongly bound 8-mers, p-values were Bonferroni corrected, and those 8-mers with a corrected p-value < 0.001 were deemed TF-preferred. *Common* 8-mers were determined separately for each pairwise comparison and defined as any 8-mer with an E-value > 0.48 in both experiments.

**Sequence and Motif Analysis.** Multiple protein alignment was performed with the MUSCLE algorithm ([www.ebi.ac.uk/Tools/msa/muscle](http://www.ebi.ac.uk/Tools/msa/muscle)) (Edgar, 2004). DNA binding site motifs for TF-preferred and common sites (as in Figure 1) were determined by running the Priority 2.1.0 motif finding algorithm (Gordan et al., 2010) on the 8-bp sequences. Graphical sequence logos were generated using enoLOGOS (Workman et al., 2005).

**Genome Analysis.** Enrichment of TF-preferred sites common to Usv1 and Rgm1 was examined in *S. cerevisiae* gene promoters identified as bound by Rgm1 ( $P < 0.01$ , 22 promoters matching genes in UCSC genome build SacCer1 ((Wang et al., 2011)). Briefly, Usv1- and Rgm1-preferred 8-mers were identified relative to Msn2 using the procedure described above. TF-preferred 8-mers common to both Usv1 and Rgm1 were used for genomic analysis. The presence of Usv1/Rgm1-preferred 8-mers was identified in gene promoter regions (600-bp upstream of the gene transcription start sites, UCSC genome build sacCer1) using custom Perl scripts (T.S, available upon request) for all 5739 SGD genes assayed in the Wang *et al.* dataset. Promoters were labeled as ‘bound’ (containing 1 or more TF-preferred 8-mers, 15/22 Rgm1-bound promoters, 1291/5739 total promoters) or ‘not bound’. An enrichment p-value was calculated using the hypergeometric distribution (i.e., Fisher’s exact test) using the R statistical package `dhyper()` function –  $dhyper(15,1291,5739-1291,22) = 5.2 \times 10^{-6}$ .

### Supplemental References

- Berger, M.F., Badis, G., Gehrke, A.R., Talukder, S., Philippakis, A.A., Pena-Castillo, L., Alleyne, T.M., Mnaimneh, S., Botvinnik, O.B., Chan, E.T., *et al.* (2008). Variation in homeodomain DNA binding revealed by high-resolution analysis of sequence preferences. *Cell* 133, 1266-1276.
- Dudley, A.M., Aach, J., Steffen, M.A., and Church, G.M. (2002). Measuring absolute expression with microarrays with a calibrated reference sample and an extended signal intensity range. *Proc Natl Acad Sci U S A* 99, 7554-7559.
- Edgar, R.C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic acids research* 32, 1792-1797.
- Elrod-Erickson, M., Rould, M.A., Nekludova, L., and Pabo, C.O. (1996). Zif268 protein-DNA complex refined at 1.6 Å: a model system for understanding zinc finger-DNA interactions. *Structure* 4, 1171-1180.
- Gordan, R., Narlikar, L., and Hartemink, A.J. (2010). Finding regulatory DNA motifs using alignment-free evolutionary conservation information. *Nucleic Acids Res* 38, e90.
- Miller, J.C., and Pabo, C.O. (2001). Rearrangement of side-chains in a Zif268 mutant highlights the complexities of zinc finger-DNA recognition. *Journal of molecular biology* 313, 309-315.
- Schaufler, L.E., and Klevit, R.E. (2003). Mechanism of DNA binding by the ADR1 zinc finger transcription factor as determined by SPR. *Journal of molecular biology* 329, 931-939.
- Workman, C.T., Yin, Y., Corcoran, D.L., Ideker, T., Stormo, G.D., and Benos, P.V. (2005). enoLOGOS: a versatile web tool for energy normalized sequence logos. *Nucleic Acids Res* 33, W389-392.

