

SUPPLEMENTARY DATA

**TFBSshape: a motif database for DNA shape features of transcription factor binding sites**

Lin Yang<sup>1</sup>, Tianyin Zhou<sup>1</sup>, Iris Dror<sup>1,2</sup>, Anthony Mathelier<sup>3</sup>, Wyeth W. Wasserman<sup>3</sup>, Raluca Gordân<sup>4</sup>, and  
Remo Rohs<sup>1,\*</sup>

<sup>1</sup> Molecular and Computational Biology Program, University of Southern California, Los Angeles, CA 90089, USA

<sup>2</sup> Department of Biology, Technion – Israel Institute of Technology, Technion City, Haifa 32000, Israel

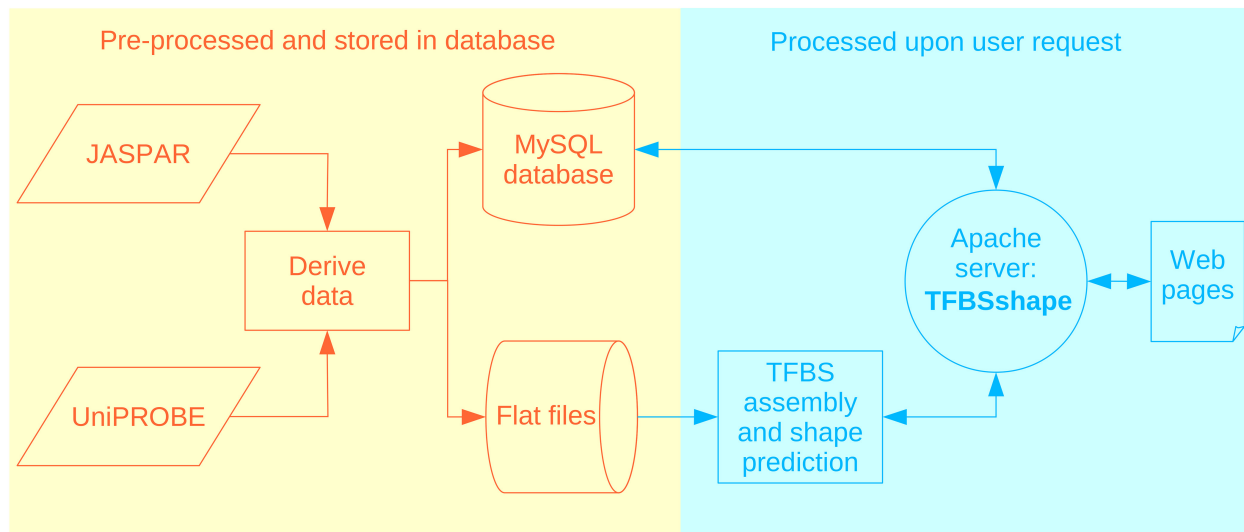
<sup>3</sup> Centre for Molecular Medicine and Therapeutics, University of British Columbia, Vancouver, BC, Canada

<sup>4</sup> Institute for Genome Sciences & Policy, Duke University, Durham, NC 27708, USA

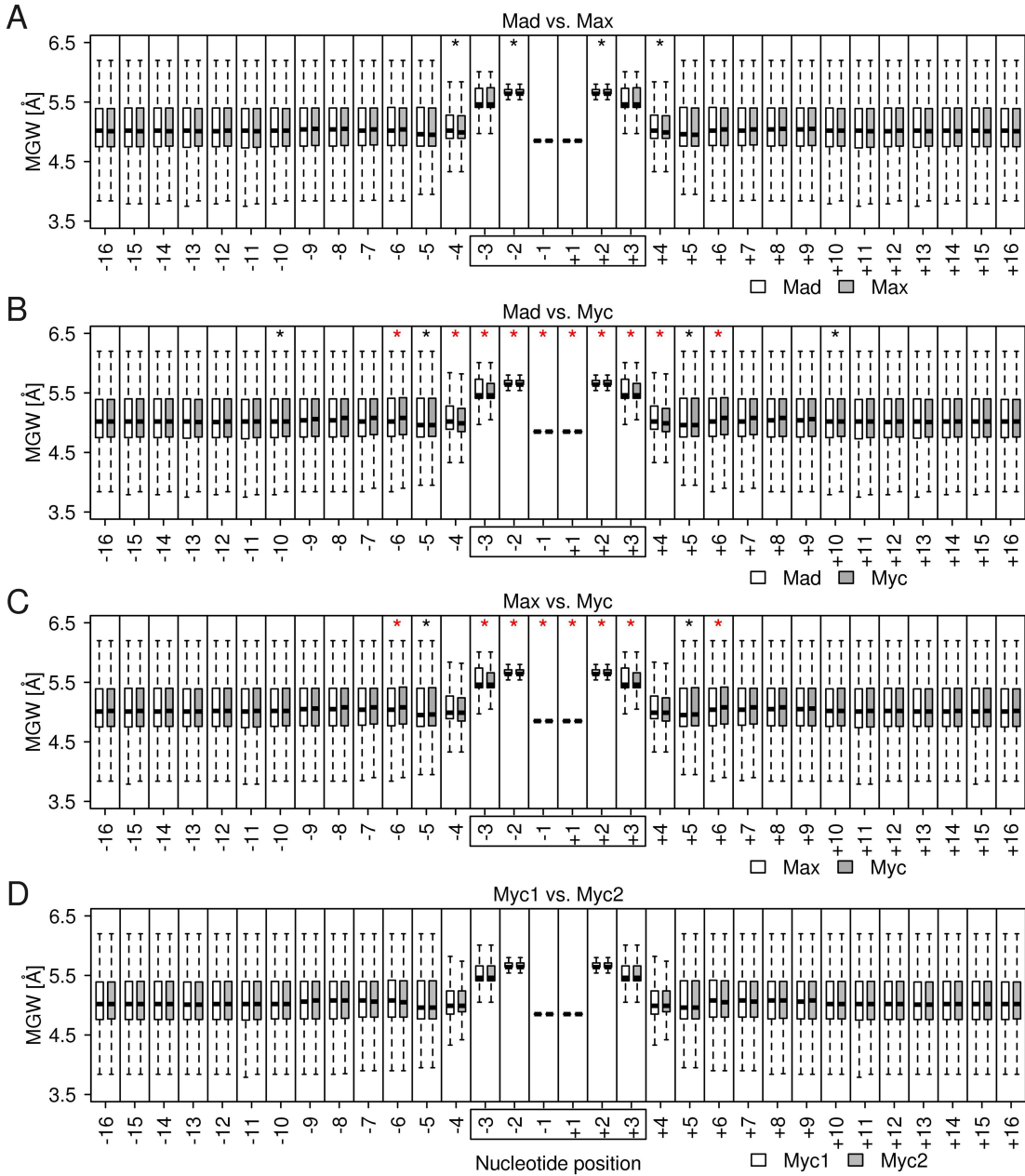
\*To whom correspondence may be addressed:

Remo Rohs, Ph.D.  
Molecular and Computational Biology Program  
1050 Childs Way RRI 404C  
University of Southern California  
Los Angeles, CA 90089  
United States  
Tel: +1-213-740-0552  
Fax: +1-213-821-4257  
Email: [rohs@usc.edu](mailto:rohs@usc.edu)

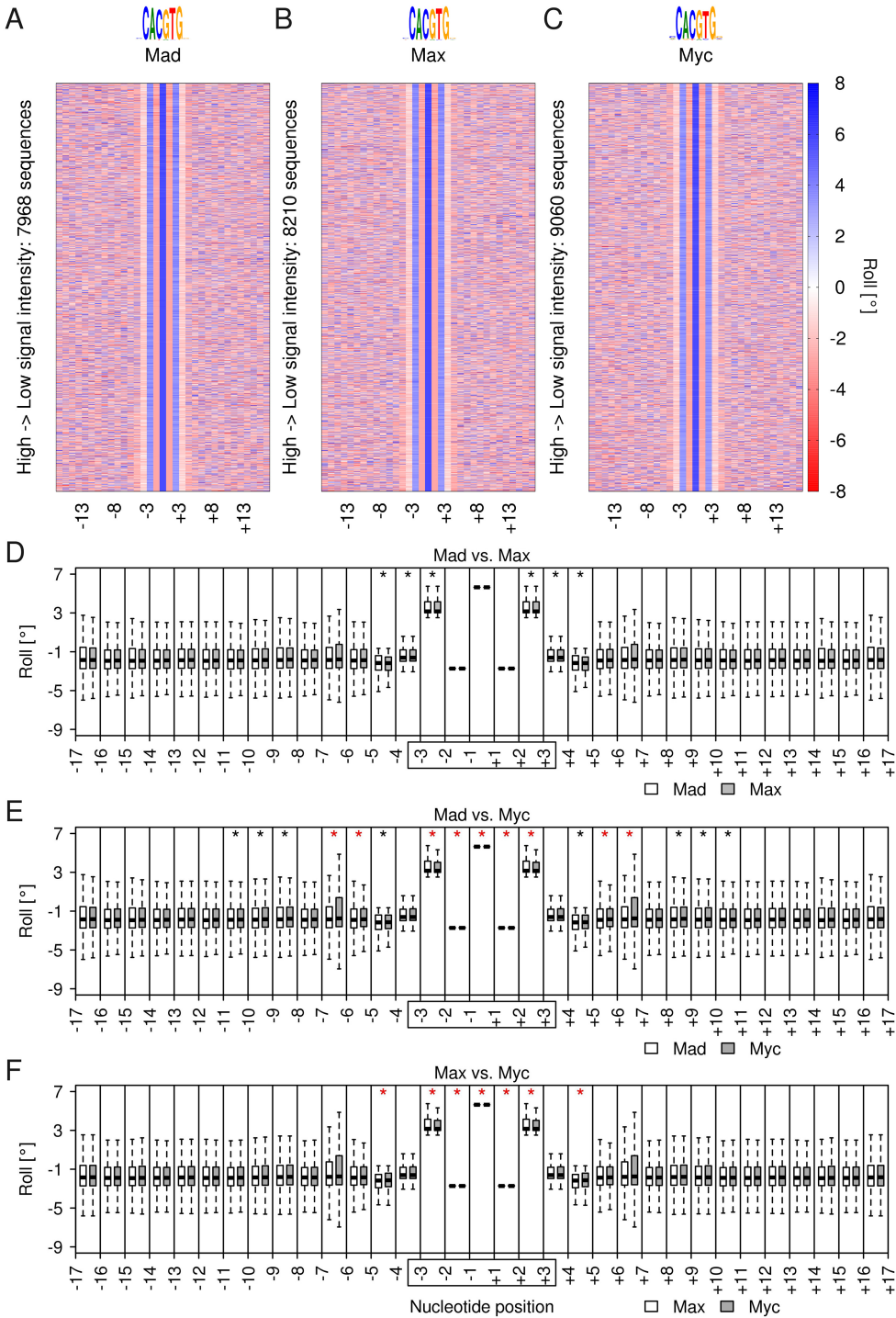
## SUPPLEMENTARY FIGURES



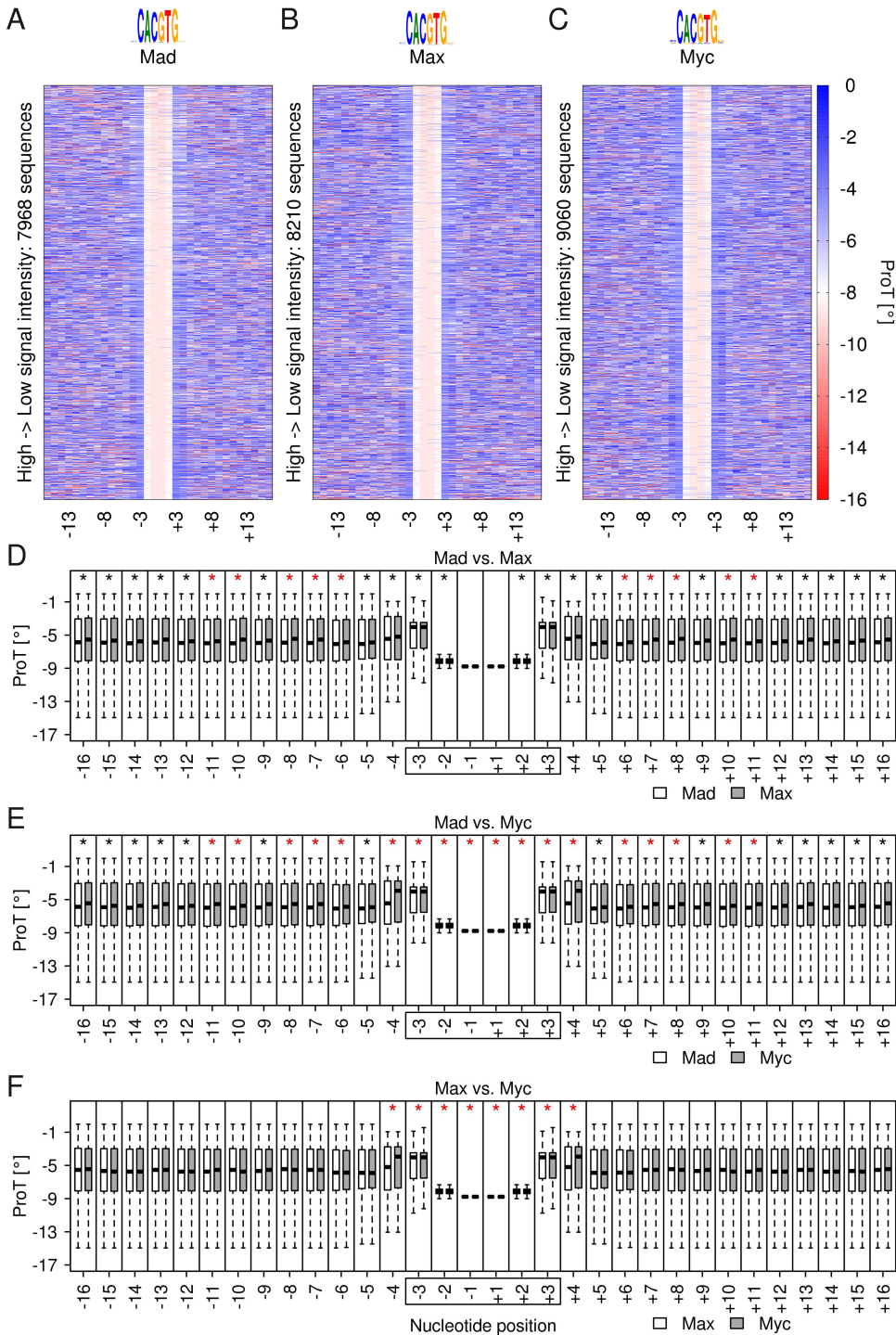
**Supplementary Figure S1.** TFBSshape database flowchart. Input data are nucleotide sequences derived from the motif databases JASPAR and UniPROBE, which are stored and managed using a MySQL database (yellow). Following this pre-processing, TFBSs are assembled, DNA shape features are predicted “on the fly”, and an Apache server provides the user interface (blue).



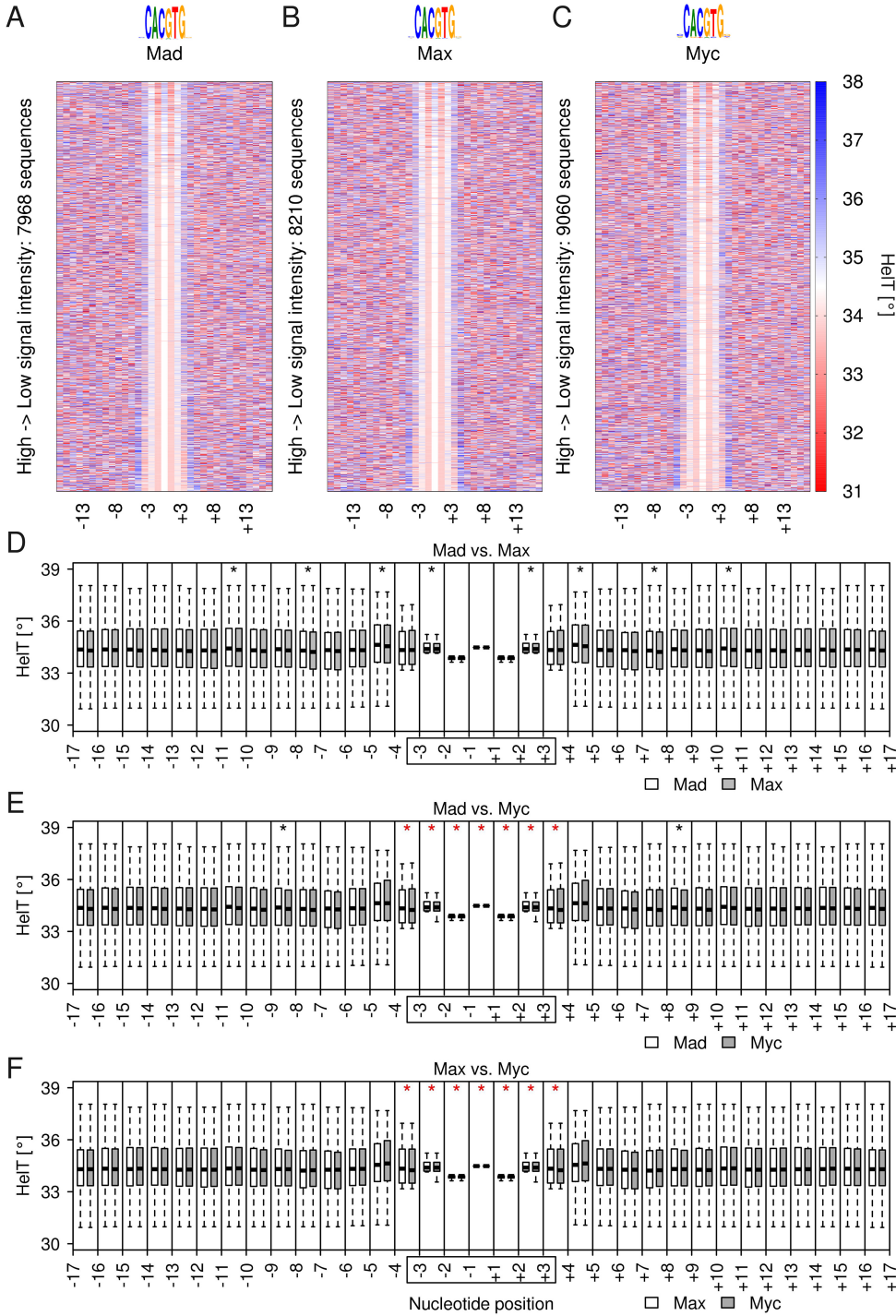
**Supplementary Figure S2.** DNA shape analysis of human bHLH TFBSs. (A-C) MGW preferences of the three TFs were compared using box plots. Boxes represent the median (line inside the box), 1<sup>st</sup> and 3<sup>rd</sup> quartiles (edges of the box), and the whiskers define the furthest data points within 1.5 x inter-quartile range from the edges of the box. Asterisks indicate nucleotide positions where differences in MGW distributions selected by (A) Mad vs. Max, (B) Mad vs. Myc, and (C) Max vs. Myc were significant based on a K-S test (black asterisk,  $P < 0.05$ ; red asterisk,  $P < 0.001$ ). (D) As a negative control, two independent experiments for Myc ('Myc1' and 'Myc2') were also compared and DNA shape selections were found to be essentially identical. MGW features were symmetrized based on the palindromic E-box.



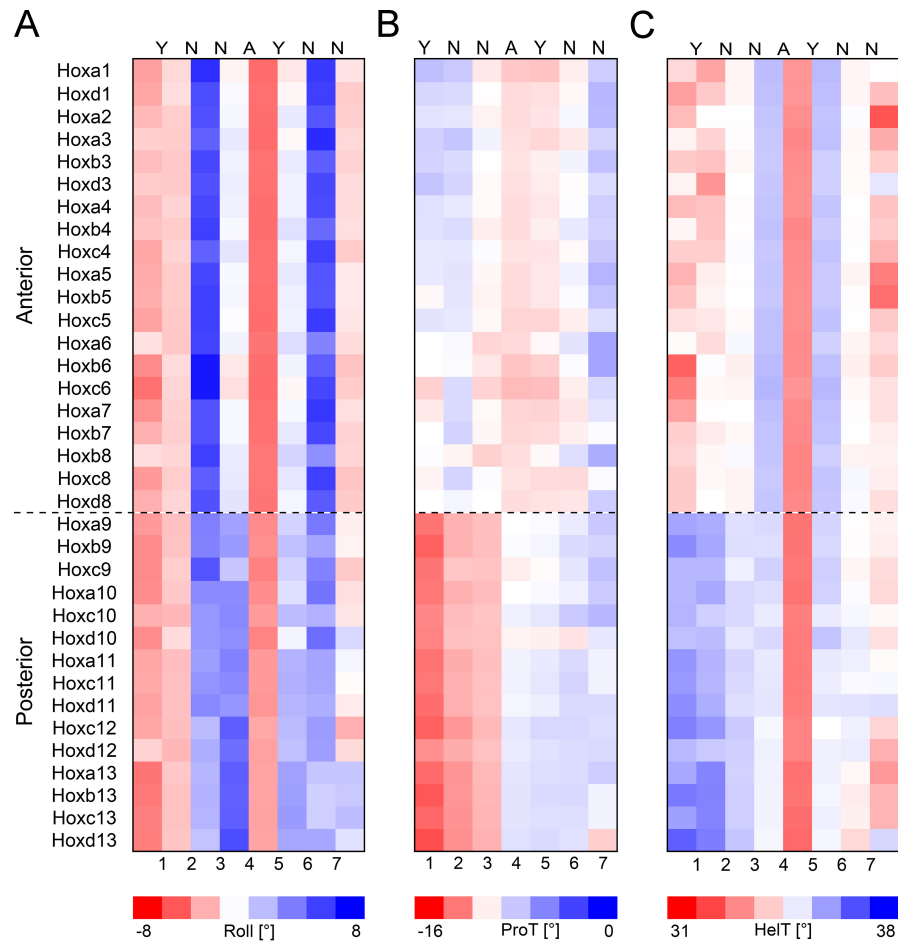
**Supplementary Figure S3.** DNA shape analysis of human bHLH TFBSs. Heat maps illustrate Roll selections of (A) Mad, (B) Max, and (C) Myc. (D-F) Roll preferences of the three TFs were compared using box plots. Boxes represent the median (line inside the box), 1<sup>st</sup> and 3<sup>rd</sup> quartiles (edges of the box), and the whiskers define the furthest data points within 1.5 x inter-quartile range from the edges of the box. Asterisks indicate nucleotide positions where differences in Roll distributions selected by (D) Mad vs. Max, (E) Mad vs. Myc, and (F) Max vs. Myc were significant based on a K-S test (black asterisk,  $P < 0.05$ ; red asterisk,  $P < 0.001$ ). Roll features were symmetrized based on the palindromic E-box.



**Supplementary Figure S4.** DNA shape analysis of human bHLH TFBSs. Heat maps illustrate ProT selections of (A) Mad, (B) Max, and (C) Myc. (D-F) ProT preferences of the three TFs were compared using box plots. Boxes represent the median (line inside the box), 1<sup>st</sup> and 3<sup>rd</sup> quartiles (edges of the box), and the whiskers define the furthest data points within 1.5 x inter-quartile range from the edges of the box. Asterisks indicate nucleotide positions where differences in ProT distributions selected by (D) Mad vs. Max, (E) Mad vs. Myc, and (F) Max vs. Myc were significant based on a K-S test (black asterisk,  $P < 0.05$ ; red asterisk,  $P < 0.001$ ). ProT features were symmetrized based on the palindromic E-box.



**Supplementary Figure S5.** DNA shape analysis of human bHLH TFBSs. Heat maps illustrate HelT selections of (A) Mad, (B) Max, and (C) Myc. (D-F) HelT preferences of the three TFs were compared using box plots. Boxes represent the median (line inside the box), 1<sup>st</sup> and 3<sup>rd</sup> quartiles (edges of the box), and the whiskers define the furthest data points within 1.5 x inter-quartile range from the edges of the box. Asterisks indicate nucleotide positions where differences in HelT distributions selected by (D) Mad vs. Max, (E) Mad vs. Myc, and (F) Max vs. Myc were significant based on a K-S test (black asterisk,  $P < 0.05$ ; red asterisk,  $P < 0.001$ ). HelT features were symmetrized based on the palindromic E-box.



**Supplementary Figure S6.** DNA shape analysis of mouse Hox TFBSs. Heat maps illustrate (A) Roll, (B) ProT, and (C) HelT preferences of monomeric mouse Hox TFs determined by universal PBM experiments (26).





## **AUTHOR CONTRIBUTIONS**

L.Y. developed the methodology for generating DNA shape data from TFBSs provided by JASPAR and from PBM probes provided by UniPROBE, designed and generated the TFBSshape database, and analyzed DNA binding specificity data for human bHLH TFs. T.Z. developed the methodology for high-throughput DNA shape prediction and implemented the multiple linear regression models. I.D. analyzed mouse Hox TFBSs. A.M. and W.W. provided unpublished JASPAR data, and enabled the integration of TFBSshape with JASPAR2014. R.G. gave advice on the analysis of human bHLH TFBSs and the detection of TFBSs on probes derived from UniPROBE. L.Y. and R.R. wrote the manuscript. R.R. conceived, designed, and supervised the project.