

## Supplementary Text S1

### The classifier's performance on the datasets used by Bhardwaj et al.<sup>1</sup> and Stawiski et al.<sup>2</sup>

Stawiski et al., developed a neural net classifier on a dataset of 54 DBPs<sup>3</sup> and 250 nDBPs. Using leave one out cross validation, the classifier reached sensitivity of 81% and specificity of 94%. We examined our descriptors on the same dataset and reached sensitivity of 87% at specificity of 94%, using 10-fold cross validation.

Bhardwaj et al.<sup>1</sup> examined their classifier on a dataset of DBPs collected from previous studies<sup>2,4,5</sup> and the set nDBPs used by Stawiski et al.<sup>2</sup> The dataset was further filtered by allowing maximum of 20% sequence identity between each pair of proteins. Using 5-fold cross validation on the filtered dataset, their support vector machine reached sensitivity of 67.4% and specificity of 94.9%.

Using 5-fold cross validation on this dataset we reached sensitivity of 73.6% at specificity of 94.9%.

### Culling the Bhardwaj dataset

Using the PISCES server we culled the PDB chains of the entries given at the web site which accompanied the paper (<http://proteomics.bioengr.uic.edu/pro-dna/>).<sup>1</sup> From the resulting list, we removed redundancy using PSI-CD-HIT<sup>6</sup> with a sequence identity cutoff of 20%. Finally we had a list of 87 DBPs and 216 nDBPs.

## References:

1. Bhardwaj, N., Langlois, R. E., Zhao, G. & Lu, H. (2005). Kernel-based machine learning protocol for predicting DNA-binding proteins. *Nucleic Acids Res.* 33, 6486-93.
2. Stawiski, E. W., Gregoret, L. M. & Mandel-Gutfreund, Y. (2003). Annotating nucleic acid-binding function based on protein structure. *J. Mol. Biol.* 326, 1065-79.
3. Luscombe, N. M., Austin, S. E., Berman, H. M. & Thornton, J. M. (2000). An overview of the structures of protein-DNA complexes. *Genome Biol.* 1, REVIEWS001.1-001.37.
4. Ahmad, S., Gromiha, M. M. & Sarai, A. (2004). Analysis and prediction of DNA-binding proteins and their binding residues based on composition, sequence and structural information. *Bioinformatics* 20, 477-86.
5. Jones, S., Shanahan, H. P., Berman, H. M. & Thornton, J. M. (2003). Using electrostatic potentials to predict DNA-binding sites on DNA-binding proteins. *Nucleic Acids Res.* 31, 7189-98.
6. Li, W. & Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 22, 1658-9.

**Supplementary Table S1:**

Score cutoff	Sensitivity	Specificity
0.49	0.906	0.713
0.50	0.906	0.725
0.51	0.891	0.734
0.52	0.884	0.745
0.53	0.884	0.757
0.54	0.877	0.779
0.55	0.870	0.790
0.56	0.862	0.800
0.57	0.855	0.810
0.58	0.848	0.821
0.59	0.833	0.830
0.61	0.819	0.840
0.62	0.812	0.851
0.63	0.804	0.861
0.64	0.775	0.868
0.65	0.732	0.873
0.66	0.717	0.883
0.68	0.717	0.893
0.69	0.710	0.904
0.70	0.667	0.921
0.72	0.645	0.929
0.75	0.623	0.937
0.77	0.587	0.943
0.78	0.580	0.954
0.80	0.536	0.958
0.83	0.464	0.969
0.85	0.428	0.975
0.87	0.377	0.979
0.90	0.333	0.985
0.91	0.283	0.987
0.93	0.217	0.989

**The sensitivity and specificity of the classifier at different score cutoffs.** The values were measured on a dataset of 138 DBPs and 843 nDBPs.

## Supplementary Table S2:

The list of PDB entries in N-Func predicted as DBPs by the classifier, sorted in decreasing order of the confidence measure.

PDB id and chain	The proportion of votes as DBPs				
<a href="#">2DBBA</a>	0.9822	<a href="#">1WOZA</a>	0.7498	<a href="#">1XA0A</a>	0.6551
<a href="#">1PC6A</a>	0.9565	<a href="#">1ZN6A</a>	0.7446	<a href="#">1KU9A</a>	0.6544
<a href="#">1G2RA</a>	0.9453	<a href="#">1YKUA</a>	0.7443	<a href="#">2GMQA</a>	0.6534
<a href="#">2HJ3A</a>	0.9308	<a href="#">1RW1A</a>	0.7423	<a href="#">1RFZA</a>	0.651
<a href="#">2GTVX</a>	0.9305	<a href="#">1L9GA</a>	0.7417	<a href="#">2GZ4A</a>	0.6506
<a href="#">2NRKA</a>	0.9118	<a href="#">2B6CA</a>	0.7377	<a href="#">2HZTA</a>	0.6501
<a href="#">2ESHA</a>	0.9107	<a href="#">1V8DA</a>	0.737	<a href="#">2IBDA</a>	0.648
<a href="#">1ZG2A</a>	0.9067	<a href="#">2FRNA</a>	0.7339	<a href="#">1J26A</a>	0.6479
<a href="#">2O38A</a>	0.8953	<a href="#">1MWWA</a>	0.7264	<a href="#">2FEFA</a>	0.647
<a href="#">1VKWA</a>	0.8832	<a href="#">2B78A</a>	0.7228	<a href="#">1UANA</a>	0.6453
<a href="#">1VBKA</a>	0.8775	<a href="#">1VDYA</a>	0.7225	<a href="#">2DJ6A</a>	0.6449
<a href="#">1NOGA</a>	0.8697	<a href="#">2AEUA</a>	0.72	<a href="#">2H3RA</a>	0.6429
<a href="#">1WVTA</a>	0.8563	<a href="#">1WW1A</a>	0.7198	<a href="#">2OHWA</a>	0.6392
<a href="#">1WJ9A</a>	0.8526	<a href="#">1TO0A</a>	0.7182	<a href="#">1XPJA</a>	0.6338
<a href="#">1XD7A</a>	0.844	<a href="#">2AR1A</a>	0.7174	<a href="#">2OD0A</a>	0.6331
<a href="#">1SDIA</a>	0.8417	<a href="#">2IL5A</a>	0.7168	<a href="#">1V30A</a>	0.6322
<a href="#">2CPXA</a>	0.8278	<a href="#">2HXJA</a>	0.7088	<a href="#">2B25A</a>	0.6299
<a href="#">1VQSA</a>	0.8247	<a href="#">1S4KA</a>	0.7036	<a href="#">2HH8A</a>	0.6268
<a href="#">2I3FA</a>	0.8245	<a href="#">1Y0KA</a>	0.7018	<a href="#">1RTYA</a>	0.6206
<a href="#">2CYYA</a>	0.8237	<a href="#">1XBWA</a>	0.6983	<a href="#">1UGJA</a>	0.6196
<a href="#">1TLJA</a>	0.8194	<a href="#">2ATZA</a>	0.697	<a href="#">2CU5A</a>	0.6195
<a href="#">1K3RA</a>	0.8181	<a href="#">1YDMA</a>	0.6943	<a href="#">1VGJA</a>	0.6162
<a href="#">1LN4A</a>	0.8152	<a href="#">2GUKA</a>	0.6929	<a href="#">2G7ZA</a>	0.6138
<a href="#">2I76A</a>	0.811	<a href="#">2IPQX</a>	0.6925	<a href="#">2ETSA</a>	0.6134
<a href="#">2GSCA</a>	0.809	<a href="#">2I51A</a>	0.691	<a href="#">2I9IA</a>	0.6121
<a href="#">2HIYA</a>	0.7965	<a href="#">1Z85A</a>	0.6901	<a href="#">1IUKA</a>	0.6049
<a href="#">2I2OA</a>	0.7903	<a href="#">2AEGA</a>	0.6881	<a href="#">2ICUA</a>	0.6049
<a href="#">1U3EM</a>	0.7864	<a href="#">1X0MA</a>	0.688	<a href="#">2IKBA</a>	0.6039
<a href="#">1VQYA</a>	0.7832	<a href="#">1K7JA</a>	0.6871	<a href="#">1ZD0A</a>	0.6029
<a href="#">2GMYA</a>	0.783	<a href="#">1VAVA</a>	0.6828	<a href="#">2GKPA</a>	0.6027
<a href="#">1UILA</a>	0.7822	<a href="#">2GM3A</a>	0.6828	<a href="#">2HD9A</a>	0.6024
<a href="#">1P9QC</a>	0.7758	<a href="#">1YX1A</a>	0.6814	<a href="#">1NG6A</a>	0.6018
<a href="#">1J27A</a>	0.7732	<a href="#">2IM9A</a>	0.68	<a href="#">1VPVA</a>	0.6004
<a href="#">2HH7A</a>	0.7664	<a href="#">2A1VA</a>	0.6797	<a href="#">2BDVA</a>	0.6001
<a href="#">1ZEEA</a>	0.7663	<a href="#">1YW1A</a>	0.6784	<a href="#">2G7JA</a>	0.5994
<a href="#">2IRUA</a>	0.7604	<a href="#">2O3AA</a>	0.6783	<a href="#">1T06A</a>	0.5973
<a href="#">2D59A</a>	0.7596	<a href="#">2CPHA</a>	0.6735	<a href="#">1VMJA</a>	0.5973
<a href="#">2F6SA</a>	0.757	<a href="#">2BDTA</a>	0.6715	<a href="#">2CWQA</a>	0.5971
<a href="#">2NWIA</a>	0.7561	<a href="#">1XJCA</a>	0.6698	<a href="#">2I6HA</a>	0.5944
<a href="#">1WHXA</a>	0.7551	<a href="#">1WHRA</a>	0.6679	<a href="#">2AMHA</a>	0.5928
<a href="#">2AP3A</a>	0.7535	<a href="#">2IF6A</a>	0.6668	<a href="#">1SMBA</a>	0.5923
<a href="#">2A8EA</a>	0.7529	<a href="#">2GNXA</a>	0.6665	<a href="#">1WXXA</a>	0.5921
<a href="#">2FNAA</a>	0.7525	<a href="#">1JRMA</a>	0.6613	<a href="#">2FUPA</a>	0.5916
<a href="#">2FSWA</a>	0.7507	<a href="#">2DP9A</a>	0.66	<a href="#">1WOLA</a>	0.59
		<a href="#">1XNEA</a>	0.6589	<a href="#">1RTTA</a>	0.5896
		<a href="#">1VPYA</a>	0.658	<a href="#">1IUJA</a>	0.5888
		<a href="#">2HWJA</a>	0.6557	<a href="#">1UFOA</a>	0.5878
		<a href="#">1YLNA</a>	0.6556	<a href="#">2AV9A</a>	0.5872

<a href="#">1T95A</a>	0.5855
<a href="#">1VPHA</a>	0.5839
<a href="#">2CRRA</a>	0.5836
<a href="#">1UJRA</a>	0.5808
<a href="#">1IHNA</a>	0.5788
<a href="#">2IM8A</a>	0.5788
<a href="#">1SH8A</a>	0.5782
<a href="#">1U5WA</a>	0.5749
<a href="#">1WY7A</a>	0.5748
<a href="#">1O6DA</a>	0.5726
<a href="#">2ALIA</a>	0.5712
<a href="#">1ZMBA</a>	0.5707
<a href="#">2CV9A</a>	0.5707
<a href="#">1TUHA</a>	0.57
<a href="#">2DLXA</a>	0.5685
<a href="#">1T6SA</a>	0.5679
<a href="#">2GBSA</a>	0.5668
<a href="#">1O3UA</a>	0.5659
<a href="#">1VLMA</a>	0.5646
<a href="#">1VKMA</a>	0.5627
<a href="#">1MK4A</a>	0.5626
<a href="#">2FBLA</a>	0.5626
<a href="#">1P90A</a>	0.5625
<a href="#">2FUJA</a>	0.5609
<a href="#">2IJCA</a>	0.5534
<a href="#">2D13A</a>	0.5533
<a href="#">2IAYA</a>	0.5527
<a href="#">1Y81A</a>	0.5525
<a href="#">1T6AA</a>	0.5519
<a href="#">2NRQA</a>	0.5512
<a href="#">1Y6IA</a>	0.5476
<a href="#">2GX8A</a>	0.546
<a href="#">2F46A</a>	0.5453
<a href="#">2IMJA</a>	0.545
<a href="#">1WICA</a>	0.5425
<a href="#">1XM7A</a>	0.5422
<a href="#">1R3DA</a>	0.5401
<a href="#">2NWUA</a>	0.5394
<a href="#">1PZXA</a>	0.5375
<a href="#">2IVYA</a>	0.5374
<a href="#">2GJGA</a>	0.5365
<a href="#">1VPQA</a>	0.5339
<a href="#">1ZXUA</a>	0.5327
<a href="#">2AH6A</a>	0.5312
<a href="#">1LFPA</a>	0.5307
<a href="#">1WEKA</a>	0.5292
<a href="#">1O50A</a>	0.5283
<a href="#">2HQYA</a>	0.528
<a href="#">2NLVA</a>	0.528
<a href="#">1SJ5A</a>	0.5276
<a href="#">1W8IA</a>	0.5275
<a href="#">2F06A</a>	0.5269
<a href="#">1VJLA</a>	0.5266
<a href="#">2ETDA</a>	0.5264

<a href="#">2IA0A</a>	0.5241
<a href="#">2F4NA</a>	0.5239
<a href="#">2O8QA</a>	0.523
<a href="#">2D4GA</a>	0.5214
<a href="#">1Y7PA</a>	0.5212
<a href="#">1UFBA</a>	0.5209
<a href="#">1PULA</a>	0.5191
<a href="#">2GUUA</a>	0.5182
<a href="#">1XXLA</a>	0.517
<a href="#">2O8IA</a>	0.5167
<a href="#">1S4CA</a>	0.5166
<a href="#">2DCEA</a>	0.5164
<a href="#">1XMTA</a>	0.5155
<a href="#">1T3UA</a>	0.5138
<a href="#">1VE0A</a>	0.5125
<a href="#">1RCUA</a>	0.5115
<a href="#">2F20A</a>	0.5095
<a href="#">1SK7A</a>	0.5082
<a href="#">1PVMA</a>	0.5068
<a href="#">1ZTCA</a>	0.5062
<a href="#">1V96A</a>	0.5053
<a href="#">1YYVA</a>	0.5014
<a href="#">2G40A</a>	0.5007
<a href="#">2B0AA</a>	0.5003

**Supplementary Table S3:**

PDB entry	Prediction score	ProFunc <sup>1</sup>	Dali <sup>2</sup>	Additional support	Predicted as RNA binding	Other
2DBB	0.9822	probable (HTH)	+	InterPro: <sup>3</sup> Winged helix repressor DNA-binding; Transcription regulator, AsnC/Lrp		
1PC6	0.9565			InterPro: Recombinase NinB		
1G2R	0.9453	long shot (DBP Template)			J.Osipiuk et al. (2001) <sup>4</sup>	
2HJ3	0.9308					ERV/ALR Sulfhydryl Oxidase- E.Vitu et al. (2006)
2GTV	0.9305	Possible (RBP Template) long shot (DBP Template)				Chorismate mutase
2NRK	0.9118	long shot (DBP Template)	+			
2ESH	0.9107	long shot (DBP Template) long shot (RBP Template)	+	InterPro: Winged helix repressor DNA-binding; Transcriptional regulator PadR-like		
1ZG2	0.9067		+	InterPro: Excinuclease ABC, C subunit, N-terminal; Pfam: GIY-YIG catalytic domain		
2O38	0.8953	probable (HTH)	+			
1VKW	0.8832	long shot (RBP Template)				InterPro: FMN-dependent nitroreductase-like superfamily
1VBK	0.8775		+		Pfam: <sup>5</sup> THUMP domain - a predicted RNA-binding domain	
1NOG	0.8697		+			InterPro:ATP:cob(I)alamin adenosyltransferase
1WVT	0.8563		+			InterPro:ATP:cob(I)alamin adenosyltransferase
1WJ9	0.8526				A.Ebihara et al. (2006) <sup>6</sup>	
1XD7	0.844	probable (HTH)	+	Prosite: rrf2-type HTH domain. A putative DNA-binding domain.		

1SDI	0.8417	long shot (DBP Template) long shot (RBP Template)				
2CPX	0.8278	probable (RBP Template)			InterPro: RNA recognition motif, RNP-1	
1VQS	0.8247					
2I3F	0.8245					j.g.mccoy et al. - Crystal Structure of a Glycolipid transfer-like protein from Galdieria sulphuraria; InterPro: Glycolipid transfer protein, GLTP superfamily
2CYY	0.8237	probable (HTH)	+	InterPro: Winged helix repressor DNA-binding; Transcription regulator, AsnC/Lrp; Bacterial regulatory protein, ArsR		
1TLJ	0.8194				InterPro: methyltransferase TYW3 (tRNA-yW- synthesising protein 3)	
1K3R	0.8181				T.I.Zarembinski et al. (2003)	
1LN4	0.8152		+		predicted novel class of RNA binding proteins. G.J.Ostheimer et al. (2002) <sup>7</sup>	
2I76	0.811					(Crystallized with NDP)
2GSC	0.809	long shot (DBP Template)			Pfam: S23 ribosomal protein	
2HIY	0.7965					
2I2O	0.7903				e.bitto et al. eif4g-like protein	
1U3E	0.7864	probable (HTH)	+	(Crystallized with dsDNA)		
1VQY	0.7832		+			
2GMV	0.783					
1UIL	0.7822	probable (RBP Template)	+		t.nagata et al. - double-stranded RNA-binding motif	

### Table S3:

**Analysis of the top ranked PDB entries from N-Func predicted as DBPs.** Prediction score, the score assigned by the classifier; ProFunc, indicates helix-turn-helix<sup>8</sup> motifs or similarity with DNA/RNA binding templates<sup>9</sup> (DBP/RBP respectively) found in the query protein; Dali, a '+' sign indicates a significant fold similarity<sup>2</sup> between the query protein and a known DBP; Additional support, specify additional support for the prediction. The last two columns specify whether there is evidence that the protein binds RNA or other molecules.

### References:

1. Laskowski, R. A., Watson, J. D. & Thornton, J. M. (2005). ProFunc: a server for predicting protein function from 3D structure. *Nucleic Acids Res.* 33, W89-93.
2. Holm, L. & Park, J. (2000). DaliLite workbench for protein structure comparison. *Bioinformatics* 16, 566-7.
3. Mulder, N. J., Apweiler, R., Attwood, T. K., Bairoch, A., Bateman, A., Binns, D., Bork, P., Buillard, V., Cerutti, L., Copley, R., Courcelle, E., Das, U., Daugherty, L., Dibley, M., Finn, R., Fleischmann, W., Gough, J., Haft, D., Hulo, N., Hunter, S., Kahn, D., Kanapin, A., Kejariwal, A., Labarga, A., Langendijk-Genevaux, P. S., Lonsdale, D., Lopez, R., Letunic, I., Madera, M., Maslen, J., McAnulla, C., McDowall, J., Mistry, J., Mitchell, A., Nikolskaya, A. N., Orchard, S., Orengo, C., Petryszak, R., Selengut, J. D., Sigrist, C. J., Thomas, P. D., Valentin, F., Wilson, D., Wu, C. H. & Yeats, C. (2007). New developments in the InterPro database. *Nucleic Acids Res.* 35, D224-8.
4. Osipiuk, J., Gornicki, P., Maj, L., Dementieva, I., Laskowski, R. & Joachimiak, A. (2001). Streptococcus pneumonia YlxR at 1.35 Å shows a putative new fold. *Acta Crystallogr. D Biol. Crystallogr.* 57, 1747-51.
5. Finn, R. D., Tate, J., Mistry, J., Coghill, P. C., Sammut, S. J., Hotz, H. R., Ceric, G., Forslund, K., Eddy, S. R., Sonnhammer, E. L. & Bateman, A. (2008). The Pfam protein families database. *Nucleic Acids Res.* 36, D281-8.
6. Ebihara, A., Yao, M., Masui, R., Tanaka, I., Yokoyama, S. & Kuramitsu, S. (2006). Crystal structure of hypothetical protein TTHB192 from Thermus thermophilus HB8 reveals a new protein family with an RNA recognition motif-like domain. *Protein Sci.* 15, 1494-9.
7. Ostheimer, G. J., Barkan, A. & Matthews, B. W. (2002). Crystal structure of E. coli YhbY: a representative of a novel class of RNA binding proteins. *Structure* 10, 1593-601.
8. Ferrer-Costa, C., Shanahan, H. P., Jones, S. & Thornton, J. M. (2005). HTHquery: a method for detecting DNA-binding proteins with a helix-turn-helix structural motif. *Bioinformatics* 21, 3679-80.
9. Laskowski, R. A., Watson, J. D. & Thornton, J. M. (2005). Protein function prediction using local 3D templates. *J. Mol. Biol.* 351, 614-26.

**Supplementary Table S4:**

Method PDB id	DBS- PRED <sup>1</sup> (sequence based)	DB-MOM <sup>2</sup>	PreDs <sup>3</sup>	Szilágyi and Skolnick <sup>4</sup>	Nimrod et al., 2008
2IVHA	+	+	+	+	+
2JG3A	-	-	+	+	+
2BGWA	-	-	-	+	+
2NQJA	-	-	-	-	+
2UZKA	+	+	+	+	+
2VLAA	+	-	+	+	+
3CLZA	+	-	+	-	+
2V6EA	+	-	+	+	+
2BSQE	-	+	-	+	+
2V1UA	-	-	-	+	+
2IVKA	+	-	+	+	+
Correct predictions	6	3	7	9	11

**The success in the identification of 11 new structures of DBPs by 5 different methods.** Each column represents a method and each row is an examined PDB entry. The '+' sign (grey shading), represents successful prediction, and a '-' stands for a false prediction. On this small set, our methods identified all the new structure correctly.

## References

1. Ahmad, S., Gromiha, M. M. & Sarai, A. (2004). Analysis and prediction of DNA-binding proteins and their binding residues based on composition, sequence and structural information. *Bioinformatics* 20, 477-86.
2. Ahmad, S. & Sarai, A. (2004). Moment-based prediction of DNA-binding proteins. *J. Mol. Biol.* 341, 65-71.
3. Tsuchiya, Y., Kinoshita, K. & Nakamura, H. (2004). Structure-based prediction of DNA-binding sites on proteins using the empirical preference of electrostatic potential and the shape of molecular surfaces. *Proteins* 55, 885-94.
4. Szilagyi, A. & Skolnick, J. (2006). Efficient prediction of nucleic acid binding function from low-resolution protein structures. *J. Mol. Biol.* 358, 922-33.