

Prospects for a Protein-DNA “Recognition Code”

Gary Stormo
ITP, UCSB
February 26, 2003

What is the value of a Protein-DNA recognition code?

Given a protein sequence, predict its binding sites and regulated genes in the genome.

Given a DNA binding site, predict which of the genome’s proteins bind to it.

Design a protein to bind a specific DNA site, or DNA site to bind a specific protein.

Proc. Nat. Acad. Sci. USA
Vol. 73, No. 3, pp. 804-808, March 1976
Biophysics

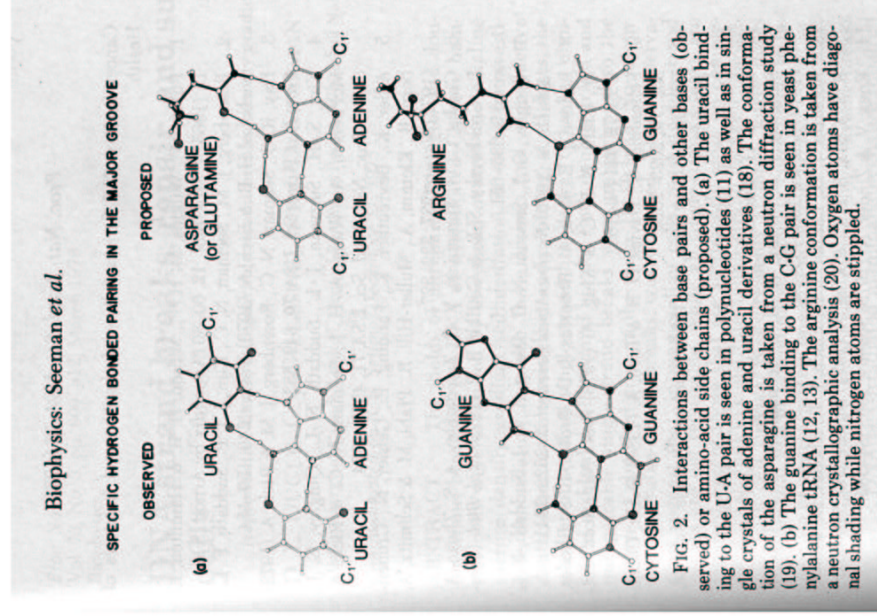
Sequence-specific recognition of double helical nucleic acids by proteins

(base pairs/hydrogen bonding/recognition fidelity/ion binding)

NADRIAN C. SEEMAN, JOHN M. ROSENBERG*, AND ALEXANDER RICH

Department of Biology, Massachusetts Institute of Technology, Cambridge, Mass. 02139

Contributed by Alexander Rich, December 10, 1975



Protein-DNA interaction**No code for recognition**

Brian W. Matthews

THIS has been something of a banner year for repressor-operator complexes. On page 321 of this issue¹, Sigler and colleagues describe the structure of *trp* repressor complexed with its synthetic operator. Other complexes involving the DNA-binding domains of λ -repressor² and of phage 434 repressor³ as well as 434 Cro protein⁴ have also now been determined. The structure at near-atomic resolution of another complex of 434 repressor was described last year⁵. It is an opportune time to take stock of the field and to consider implications for the future.

that the binding geometry of the recognition helix on the DNA is similar for many repressor proteins¹⁷. The complexes now available show that the geometry of binding is variable. Even in the case of 434 Cro⁴ and 434 repressor headpiece^{3,5}, two proteins with 50 per cent amino-acid sequence identity and very similar three-dimensional structures, and which are bound to the same operator, the respective helix-turn-helix units interact with the DNA in similar, but distinctly different, ways⁴.

Specificity. The early studies of Cro, CAP

the protein and the phosphate backbone contribute to the overall affinity of binding; sequence-specific contacts with the exposed parts of the operator base pairs allow the repressor to discriminate between its operator site and other DNA sequences.

Direct or indirect readout? The complexes involving λ -repressor headpiece, 434 Cro and 434 repressor headpiece all display multiple contacts both to the DNA backbone and to the parts of the base pairs that are exposed within the grooves of the DNA^{3,5}. Although deviations from uniform B-form are observed in the conformations of the DNA, the recognition of the specific operator sequences appears to occur primarily by direct readout. In other words, the DNA retains essentially B-type conformation and the protein directly reads the sequence infor-

doi:10.1006/jmbi.2000.3918 available online at <http://www.idealibrary.com on> **IDEAL** J. Mol. Biol. (2000) 301, 597-624**JMB****Geometric Analysis and Comparison of Protein-DNA Interfaces: Why is there no Simple Code for Recognition?**

Carl O. Pabo* and Lena Nekludova

Howard Hughes Medical Institute, Department of Biology 68-580, Massachusetts Institute of Technology Cambridge, MA 02139, USA

Structural studies of protein-DNA complexes have shown that there are many distinct families of DNA-binding proteins, and have shown that there is no simple "code" describing side-chain/base interactions. However, systematic analysis and comparison of protein-DNA complexes has been complicated by the diversity of observed contacts, the sheer number of complexes currently available and the absence of any consistent method of comparison that retains detailed structural information about the protein-DNA interface. To address these problems, we have developed geometric methods for characterizing the local structural environment in which particular side-chain/base interactions are observed. In particular, we develop methods for analyzing and comparing spatial relationships at the protein-DNA interface. Our method involves attaching local coordinate systems to the DNA bases and to the C α atoms of the peptide backbone (these are relatively rigid structural units). We use these tools to consider how the position and orientation of the polypeptide backbone (with respect to the DNA) helps to determine what contacts are possible at any given position in a protein-DNA complex. Here, we focus on base contacts that are made in the major groove and we use spatial relationships in analyzing: (i) the observed patterns of side-chain/base interactions; (ii) observed helix docking orientations; (iii) family/subfamily relationships among DNA-binding proteins; and (iv) broader questions about evolution, altered specificity mutants and the limits for the design of new DNA-binding proteins. Our analysis, which highlights differences in spatial relationships in different complexes and at different positions in a complex, helps explain why there is no simple, general code for protein-DNA recognition.

Keywords: protein-DNA complexes; recognition code; recognition helix; major groove; α helix

*Corresponding author

© 2000 Academic Press

Kinds of Codes

Deterministic (Simple, one-to-one)

Genetic Code: CCA _ Pro

Probabilistic (Degenerate, many-to-one)

Genetic Code (reversed)

Pro _ {CCA, CCC, CCG, CCT}

Probabilistic (many-to-many)

{cl, cro} _ {OR1, OR2, OR3, ...}

Existing models...

A. Qualitative models

• Desjarlais and Berg (1992)

• Pabo *et al.* (1997)

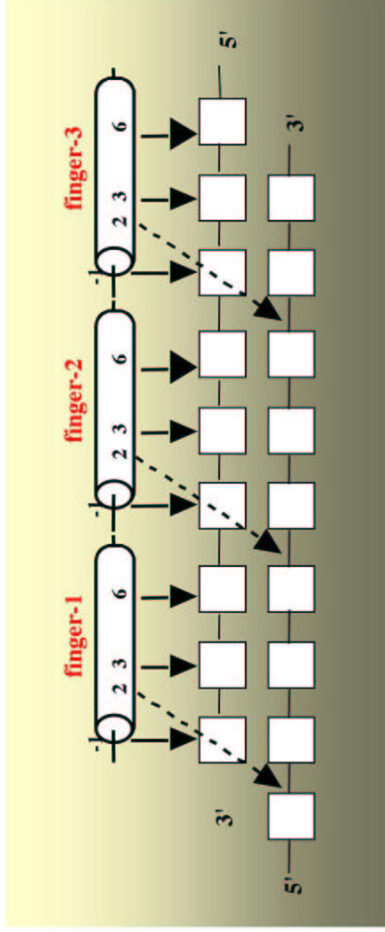
• Choo and Klug (1997)

B. Quantitative models

• Suzuki *et al.* (1995)

• Margalit *et al.* (1998)

EGR family of transcription factors



Zn-finger (Cys_2His_2):

- 3 fingers, binding in a modular fashion
- target site: 4 bases long (for each finger)
- one base overlap in the target of each finger

Qualitative models...

- Usually in the form of a simple “binary” table.
- Difficult to expand to other than “one-to-one” model.
- By nature, unsuitable for quantitative predictions.

Position in triplet		5'	Middle	3'
A	Gln 6	Asn 3 Ser 3 His 3	Gln -1	
C	Ser 2	Asp 3 Thr 3 Val 3	Asp -1	
G	Arg 6 Lys 6 Asp 2 Ser 2 Phe 2	His 3 Lys 3	Arg -1	
T	Lys 2 Asp 2	Thr 3 Ala 3 Ser 3 Val 3	Leu -1 Thr -1 Asn -1	

		Position in Zn finger			
		-1	+2	+3	+6
A.	1	A			Q
	2	A		HNS DTV HK ASTV	
	3	A	Q D R LNT		
	4	A			S DFS DK
		A C G T			
Position in DNA					

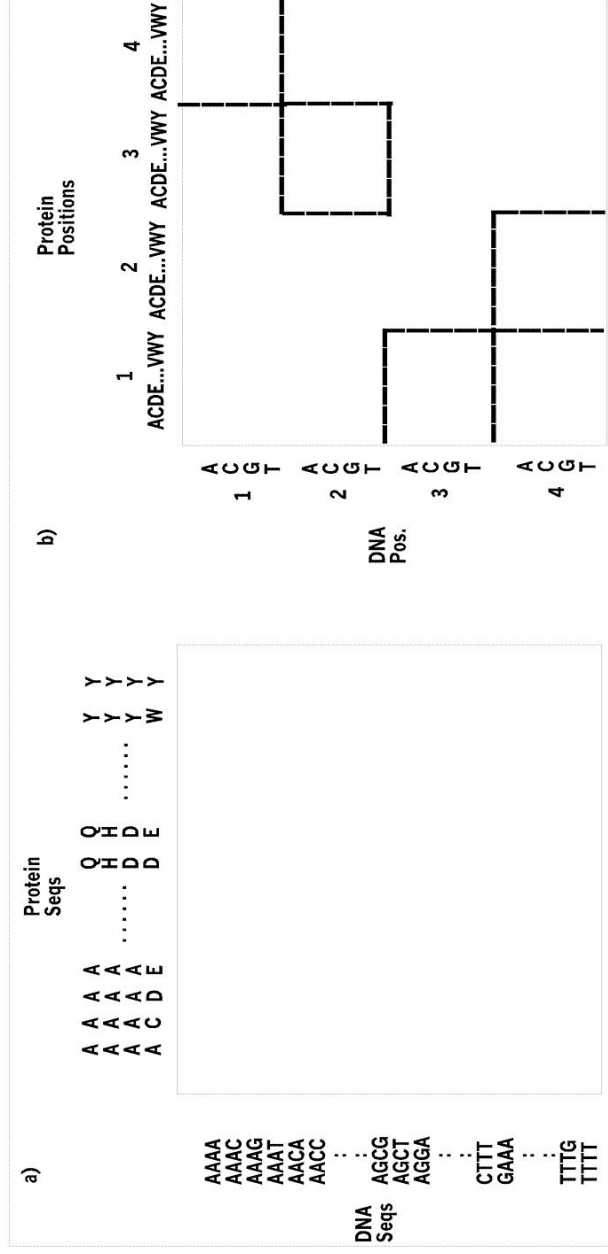
		Position in Zn finger			
		-1	+2	+3	+6
B.	1	A			
	2	A			
	3	A			
	4	A			
		A C G T			
Position in DNA					

a)

AAAA	2.63	:	:	:
AAAC	2.63	:	:	:
AAAG	2.63	:	:	:
AAAT	0.46	:	:	:
AACA	1.79	:	:	:
AACC	1.79	:	:	:
AACG	1.79	:	:	:
AACT	-0.38	:	:	:
:	:	:	:	:
:	:	:	:	:
:	:	:	:	:
AGCG	-0.25	:	:	:
AGCT	-2.42	:	:	:
AGGA	0.59	:	:	:
:	:	:	:	:
:	:	:	:	:
:	:	:	:	:
CTTT	-0.37	:	:	:
GAAA	4.01	:	:	:
:	:	:	:	:
:	:	:	:	:
:	:	:	:	:
TTTG	2.63	:	:	:
TTTT	0.46	:	:	:

b)

	1	2	3	4
A	-0.55	+1.38	+0.42	+1.38
C	0	+0.55	-0.42	+1.38
G	+0.83	-0.66	+0.42	+1.38
T	+0.83	+0.42	0	-0.79



4442-4451 Nucleic Acids Research, 2002, Vol. 30 No. 20

© 2002 Oxford University Press

Additivity in protein-DNA interactions: how good an approximation is it?

Panayiotis V. Benos^{1,2}, Martha L. Bulyk³ and Gary D. Stormo^{1,*}

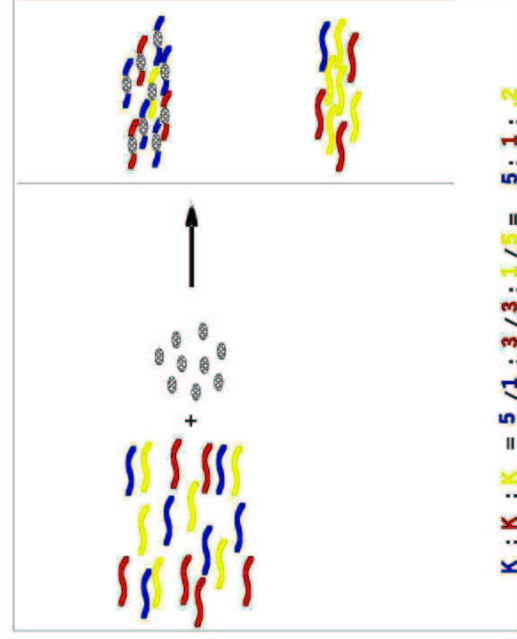
¹Department of Genetics, Campus Box 8232, Washington University, School of Medicine, St Louis, MO 63110, USA, ²Department of Human Genetics and Center for Computational Biology and Bioinformatics and Cancer Institute, University of Pittsburgh, PA 15261, USA and ³Division of Genetics, Department of Medicine and Department of Pathology, Brigham and Women's Hospital and Harvard Medical School, and Harvard-MIT Division of Health Sciences and Technology, Boston, MA 02115, USA

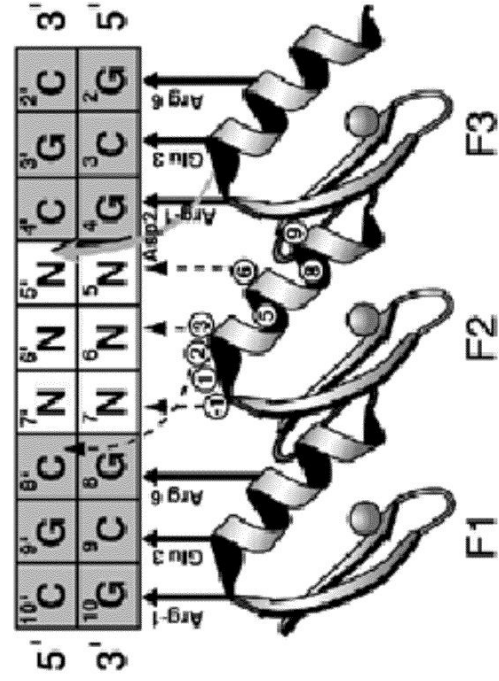
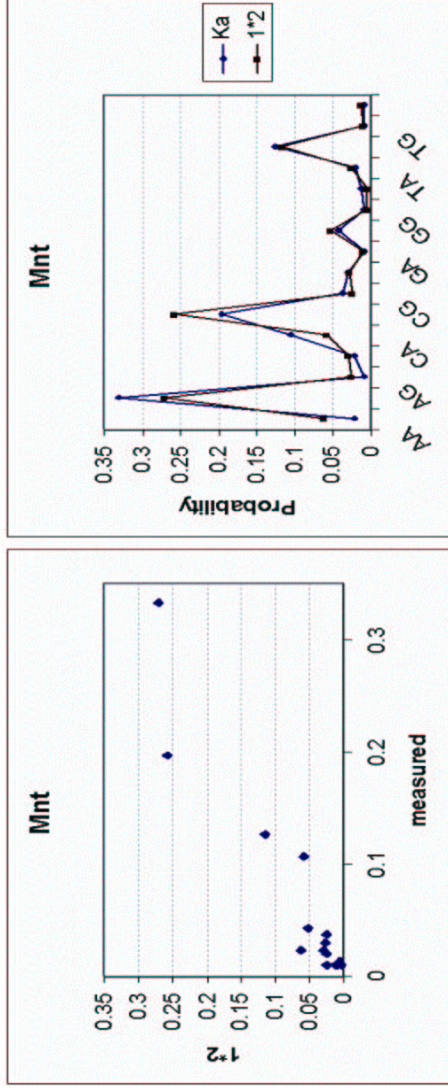
Nucleotides of transcription factor binding sites exert interdependent effects on the binding affinities of transcription factors

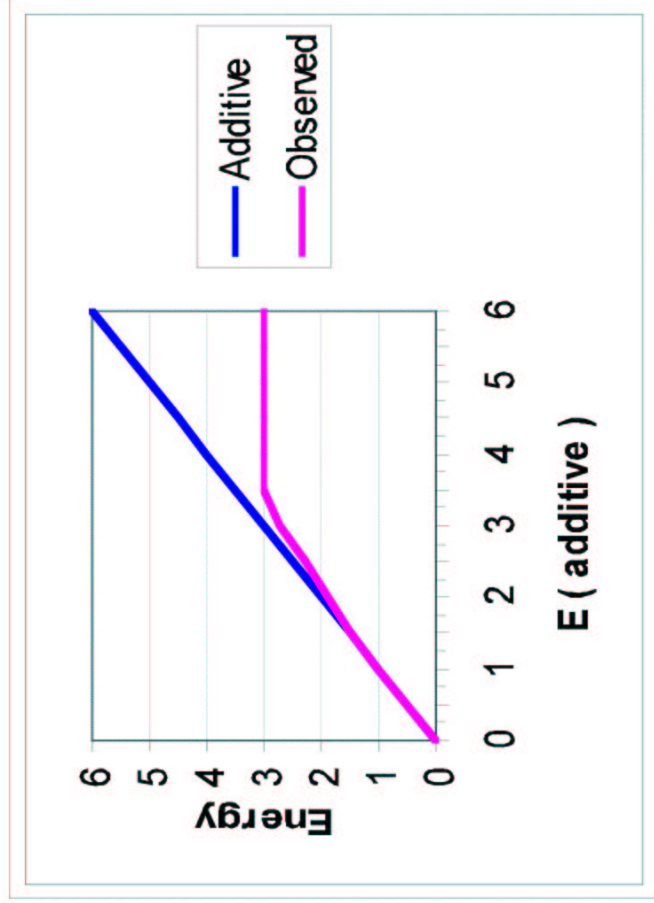
Martha L. Bulyk , Philip L. F. Johnson and George M. Church
Nucleic Acids Research, 2002, Vol. 30, 1255–1261

Non-independence of Mnt repressor-operator Interaction determine by a new quantitative multiple Fluorescence relative affinity (QuMFRA) assay

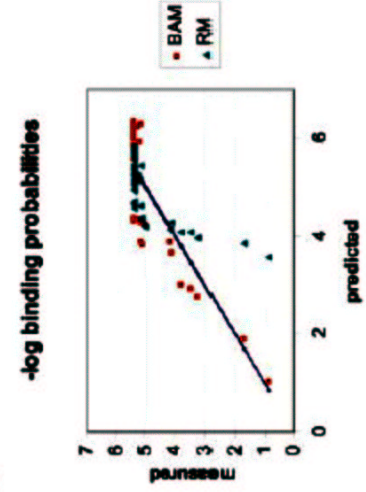
Tsz-Kwong Man and Gary D. Stormo
Nucleic Acids Research, 2001, Vol. 29, 2471-2478



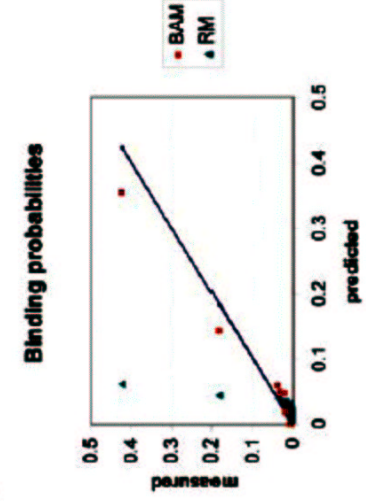


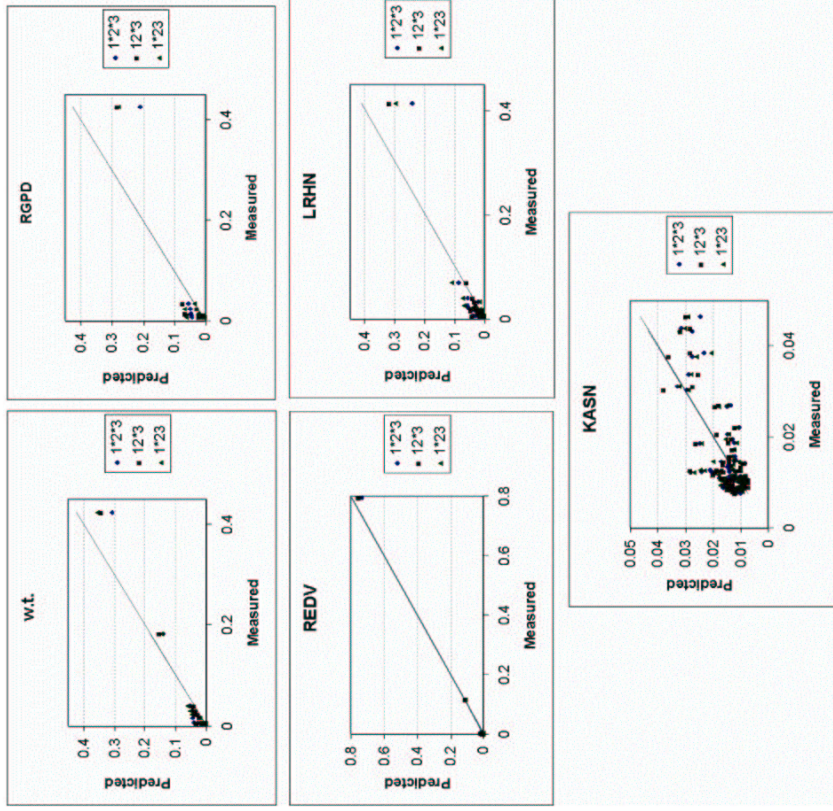


A



B





Correlation coefficients				
Zif268 variant	1*2*3	12*3	1*23	1*23
w.t.	0.973	0.986	0.987	
RGPD	0.883	0.942	0.941	
REDV	0.999	0.999	0.999	
LRHN	0.927	0.978	0.956	
KASN	0.695	0.791	0.718	

doi:10.1016/S0022-2836(02)00917-8 available online at <http://www.idealibrary.com on>  J. Mol. Biol. (2002) 323, 701–727

JMB



**Probabilistic Code for DNA Recognition by Proteins of
the EGR Family**

Panayiotis V. Benos¹, Alan S. Lapedes² and Gary D. Stormo^{1*}

**Statistical
Algorithm for
Modeling
Interaction
Energies**

Statistical Algorithm for Modeling Interaction Energies

Assuming equilibrium, Boltzmann leads us...

$$P(\underline{c}) = \frac{e^{-E(\underline{c})}}{\sum_{\underline{x}} e^{-E(\underline{x})}}$$



Two kinds of *in vitro* selection experiments

- SELEX
 - Fixed protein sequence selects preferred DNA sites from a random pool; usually returns more than one DNA sequence for each protein
- Phage display
 - Fixed DNA sequence selects preferred protein from a random pool (only “critical residues” are randomized); usually returns more than one protein sequence

SAMIE

In a SELEX experiment:

$$P(\underline{\mathbf{n}} | \underline{\mathbf{A}}) = \frac{P_{\text{ref}}(\underline{\mathbf{n}}) e^{-E(\underline{\mathbf{n}}, \underline{\mathbf{A}})}}{\sum_{\underline{\mathbf{x}}} P_{\text{ref}}(\underline{\mathbf{x}}) e^{-E(\underline{\mathbf{x}}, \underline{\mathbf{A}})}}$$

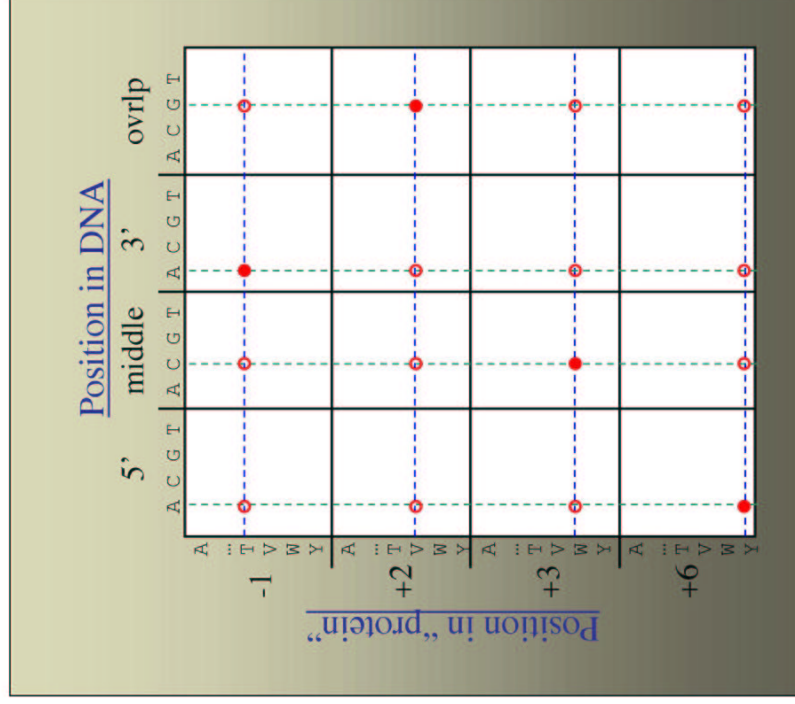
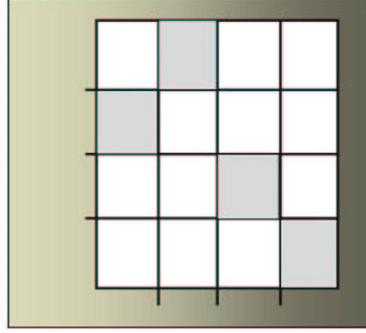
where energy:

$$E(\underline{\mathbf{n}}, \underline{\mathbf{A}}) = \sum_{ijab} C_{ij} n_i^a T_{ij}^{ab} A_j^b$$

$$E(\underline{\mathbf{n}}, \underline{\mathbf{A}}) = \sum_{ijab} C_{ij} n_i^a T_{ij}^{ab} A_j^b$$

TVWY selects **ACAG**

“Connectivity matrix” “C”



SAMIE

Training algorithm.

- Calculate T_{ij}^{ab} that maximises log-likelihood of the data.
- *Steepest ascents* method.

$$P(\{\text{data}\}) = \prod_{\{\text{data}\} \ni \underline{n}} P(\underline{n} | \underline{A})$$

$$\log P(\{\text{data}\}) = \sum_{\{\text{data}\} \ni \underline{n}} \log P(\underline{n} | \underline{A})$$

$$P(\underline{n} | \underline{A}) = \frac{P_{\text{ref}}(\underline{n}) e^{-E(\underline{n}, \underline{A})}}{\sum_{\underline{x}} P_{\text{ref}}(\underline{x}) e^{-E(\underline{x}, \underline{A})}}$$

The EGR data

Binding examples (*three fingers*).

- 1033 examples of EGR protein family (919 non-redundant)
- 322(304) SELEX and 431(399) phage display.

Randomisation experiments (*single fingers*).

- 366 distinct SELEX results
- 70 tetra-nucleotides were selected by 180 tetra-peptides.
- 444 distinct phage display results
- 361 tetra-peptides were selected by 59 tetra- nucleotides.

Total of 2009 base-amino acids interacting pairs

Table 2. The energy matrix as it was calculated by SAMEE when trained on the COMBINED_6 dataset

	Sligoer position = -3; base position = 3				Sligoer position = -2; base position = 4				Sligoer position = -1; base position = 5				Sligoer position = 0; base position = 6				
	A	C	G	T	A	C	G	T	A	C	G	T	A	C	G	T	
A	3.29	0.98	3.32	6.43	0.25	-0.87	3.32	6.32	0.54	3.12	3.27	-1.34	-0.35	3.36	1.36	-0.35	0.72
C	3.21	4.71	6.29	4.51	3.56	-0.95	6.32	4.08	4.16	1.89	6.27	-0.06	4.68	5.07	5.31	5.39	0.72
D	3.59	-2.36	0.39	-0.68	-0.68	2.27	0.42	-0.44	-0.29	5.82	7.20	6.32	1.75	0.04	1.37	-0.21	-0.21
E	0.94	-1.06	-0.05	-0.96	6.34	-1.72	1.96	0.29	1.50	0.43	1.96	1.16	0.09	1.08	1.26	-0.06	-0.06
F	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
G	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
H	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
I	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
J	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
K	0.53	-0.97	-1.14	-1.14	6.59	-0.84	3.70	6.47	5.54	7.47	1.04	0.07	1.04	0.87	6.09	-1.03	-1.36
L	3.18	0.39	7.89	-0.46	7.05	6.63	3.70	1.34	6.93	1.54	8.16	0.89	1.09	2.44	0.39	7.04	0.74
M	6.32	4.99	6.32	-1.55	6.33	5.55	2.41	0.28	0.26	7.32	2.86	-0.29	6.91	5.71	5.86	6.27	0.74
N	0.94	0.06	0.06	-1.35	1.14	-1.72	0.25	1.21	-2.73	0.86	-0.005	-0.80	0.81	0.81	-1.03	1.75	0.74
O	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
P	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
Q	-2.50	6.32	0.32	0.32	-1.32	1.32	1.32	1.32	0.32	0.32	0.32	0.32	0.32	0.32	0.32	0.32	0.32
R	1.01	3.26	-1.36	0.46	1.66	-1.65	2.17	7.37	6.59	3.54	2.56	2.14	0.36	0.36	0.36	1.39	0.39
S	3.06	0.35	0.41	-0.57	2.29	-1.45	0.26	-0.81	6.56	1.15	1.69	-1.65	-0.21	0.18	-0.40	-0.39	-0.39
T	0.27	1.32	0.39	-1.32	-0.18	-1.32	1.39	1.39	7.50	-0.04	3.36	0.36	-1.19	0.20	-0.09	-0.46	0.76
U	2.59	4.70	4.39	0.92	7.01	0.17	3.76	0.77	8.01	0.79	8.32	0.80	1.13	1.38	0.46	0.76	0.76
V	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
W	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
X	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
Y	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72
Z	3.21	4.71	1.39	4.51	3.86	-3.55	3.36	4.08	4.16	6.92	6.27	5.13	4.68	0.42	1.21	1.39	0.72

The COMBINED_6 dataset contains both COMPLEX and phage display data, according to the state-of-the-art model of Stormo. The energy values presented here have been renormalized to that the average binding constant for each contact is 1.0. This renormalization does not affect the calculation of probability, as described in the text. The renormalized matrix shows that the energetic potential of a base-motivator contact varies, depending on its position in the DNA target and the protein.

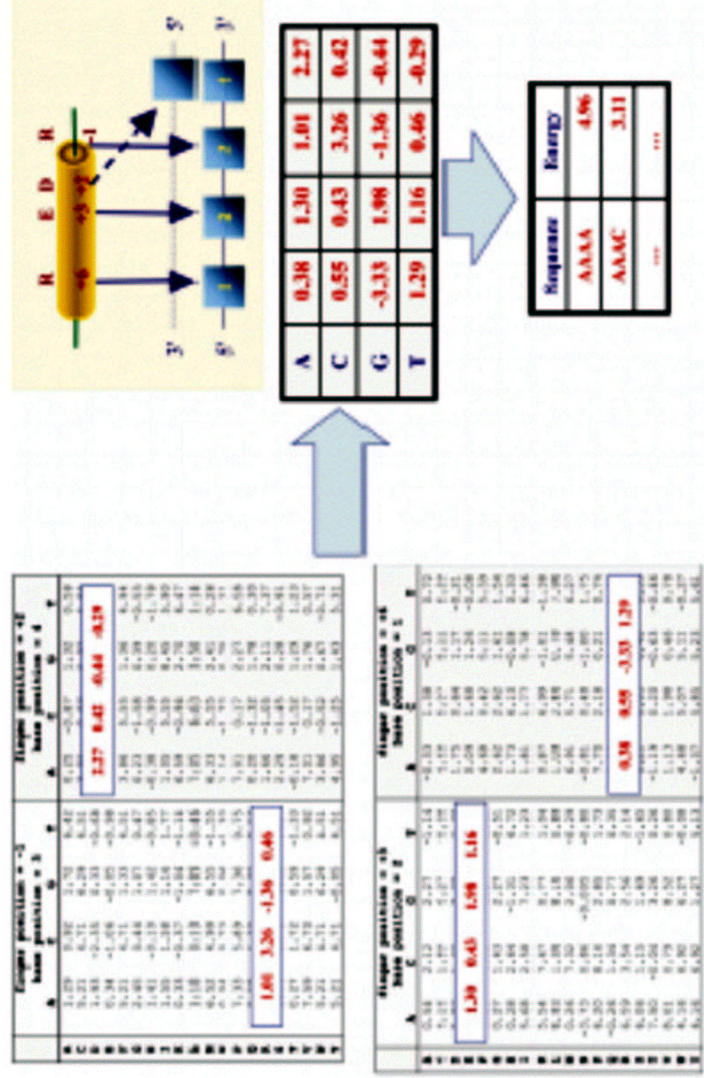


Figure 8 (Signed appendix)

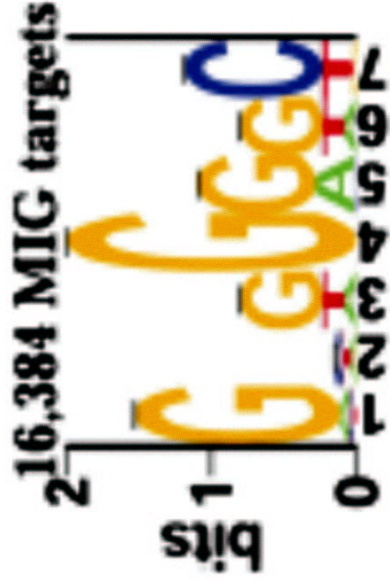


Figure 4. The weighted LOGO of the motif at six finger protein binding sites. This LOGO plots all possible nucleotide targets, weighted by the probability of occurrence, on its horizontal axis. The vertical axis formal binding site with the highest affinity to the binding protein is 5'-GGGGGA-3'.

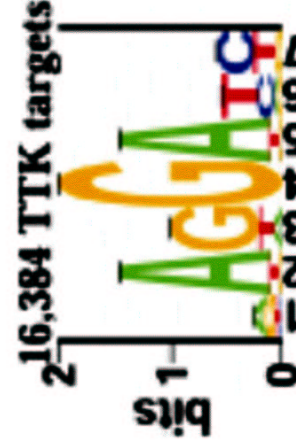
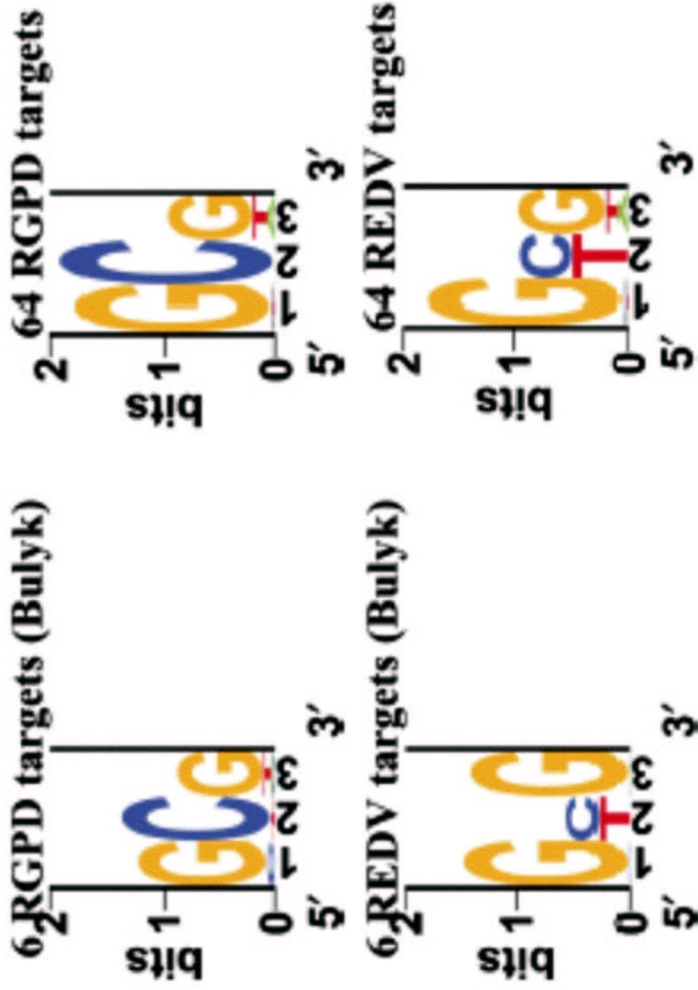


Figure 5. The weighted LOGO of the Drosophila six finger protein binding site. This LOGO plots all possible nucleotide targets, weighted by the probability of occurrence, on its horizontal axis. The vertical axis TTK binding site with the highest affinity to the binding protein is 5'-AGGATC-3'.

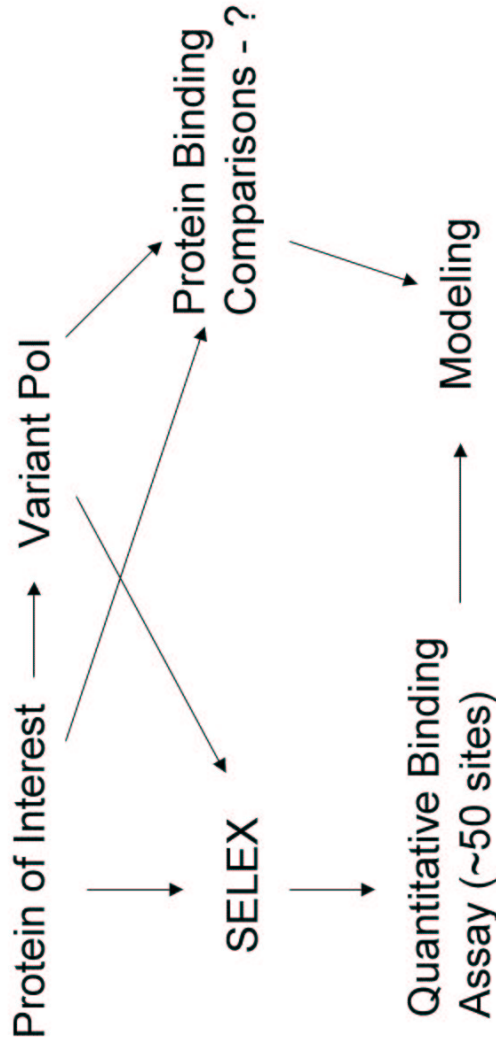


Quantitative Predictions

<u>Source/protein</u>	<u>Correlation Coefficient</u>
Hamilton	
EGR1	0.75
WT1	0.76
Segal (several variants)	0.82
Miller	
EGR1	0.93
D20A	0.81
Bulyk	
EGR1	0.61
RGP	0.99
REDV	0.73
LRHN	0.56
KASN	0.16

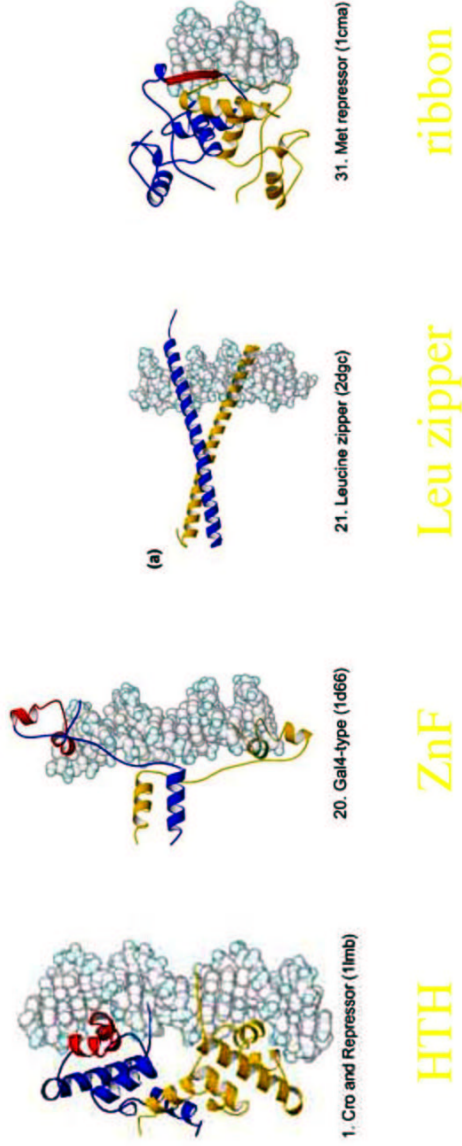


Next: need more quantitative data
for protein families of interest



Being done by Post-Docs: JJ Liu, Yong Yin

Examples of DNA-binding Proteins



from Luscombe *et al.* (2000)
<http://genomebiology.com/2000/1/1/reviews/001.1>