



**Motivation:**

Interactions between proteins and DNA molecules lie at the core of fundamental cellular processes such as transcriptional regulation. Some of these interactions have been experimentally described at atomic scale, but the molecular details of many others remain to be discovered. This work investigates ways to exploit the current knowledge about protein-DNA interfaces contained in the Protein Data Bank, with two aims: 1) modelling similar interfaces related by homology and 2) identifying genomic regulatory motifs.

**Results:**

First, we describe the structural and evolutionary conservation of protein-DNA interfaces, and the limits they impose on modelling accuracy. Second, we estimate the error rate associated to modelling interface amino acid side chains, those likely to be contacting nucleotide N-bases. Third, a simplistic protocol is implemented and different parameters are benchmarked on a set of 85 regulators from *Escherichia coli*. Finally, this approach is made available via the TFmodeller web server.

**protein-DNA interface conservation and modelling limitations**

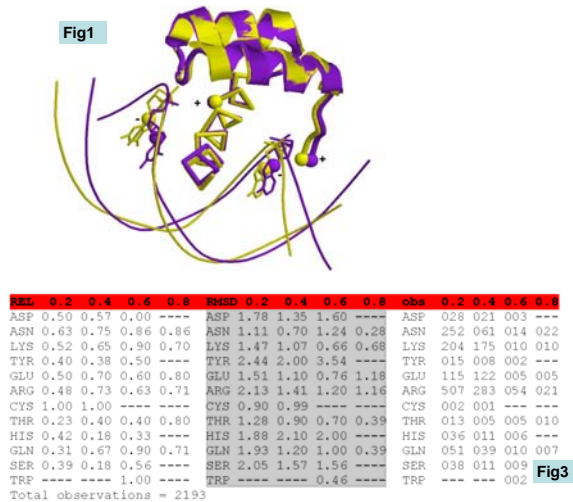
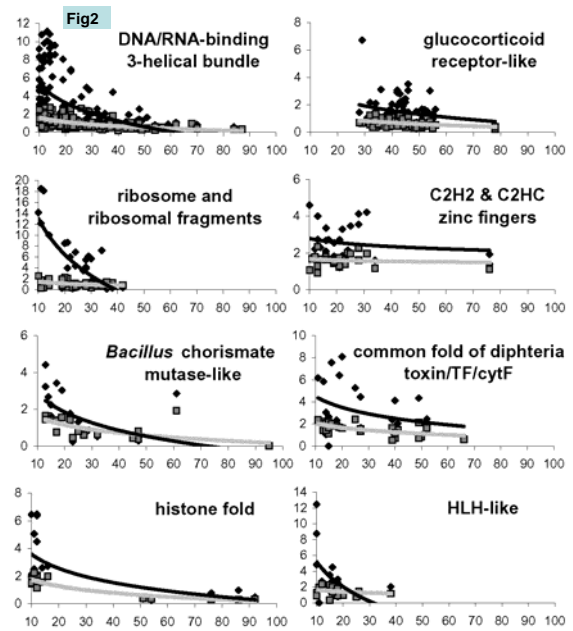


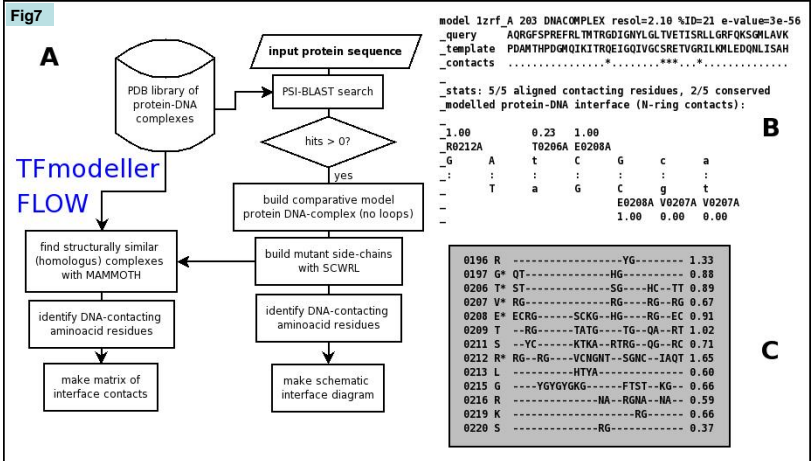
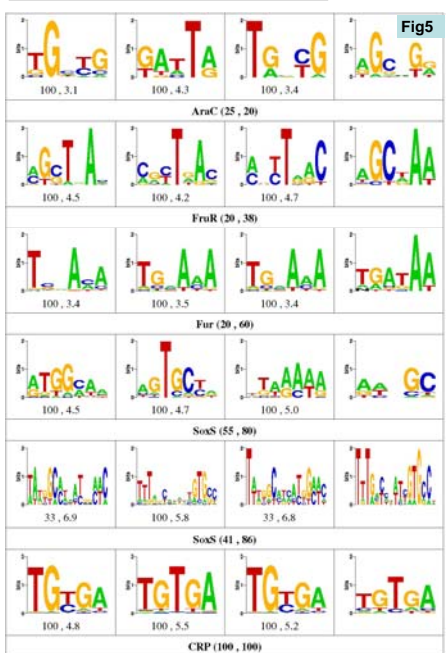
Figure 1 shows two superposed *Drosophila* TFs with protein sequences that are 40% identical, highlighting their structural conservation. Figure 2 shows the conservation trends observed with 442 superposed DNA-binding domains annotated in SCOP, with grey data corresponding to the conservation of the protein interfaces (+) and the black points to the DNA interface (-). This data is derived from 442 superposed protein-DNA complexes, that were also analyzed to generate Table 3, that shows how accurately different H-bonding residues are modelled by SCWRL. Accuracy is measured in terms of RMSD and in terms of coincident DNA contacts (REL).



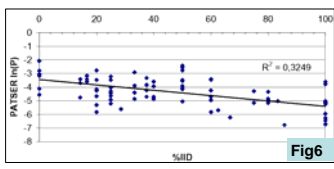
**Fig4**

	G	A	T	C
Gly	-3.93	-3.93	-3.93	-3.93
Ala	-3.93	-3.93	0.66	-3.72
Val	-3.93	-3.93	-0.17	-3.57

**footprint benchmark and TFmodeller web server (www.ccg.unam.mx/tfmodeller)**



By using a simplistic scoring scheme (Table 4) it is possible to estimate the preferred DNA sequences discriminated by a set of 85 modelled *E. coli* TFs. Results from selected examples are shown in Figure 5, including TFs modelled with remote templates, such as AraC, and perfect templates, such as CRP. Templates are evaluated with two measures (interface, and overall sequence identity). Three different predictions are shown for each TF (left), next to the sequence logo derived from annotated binding sites in RegulonDB (right). Figure 6 shows a significant correlation between the accuracy of the predicted DNA motifs and the interface identity of the template used to build the complex. Figure 7 summarizes how this approach was implemented in order to set up the TFmodeller server, and how the output looks like.



**References**

Contreras-Moreira, B., Branger, P.A. & Collado-Vides, J. (2007). TFmodeller: comparative modelling of protein-DNA complexes. *Bioinformatics*, doi: 10.1093/bioinformatics/btm148.  
 Contreras-Moreira, B. & Collado-Vides, J. (2006). Comparative footprinting of DNA-binding proteins. *Bioinformatics*, 22(14): e74.