

On the prediction of DNA-binding preferences of C2H2-ZF domains using structural models: application on human CTCF

Alberto Meseguer^{1,†}, Filip Årman^{1,†}, Oriol Fornes², Ruben Molina-Fernández¹,
Jaume Bonet³, Narcis Fernandez-Fuentes^{4,5} and Baldo Oliva^{1,*}

¹Structural Bioinformatics Lab (GRIB-IMIM), Department of Experimental and Health Science, University Pompeu Fabra, Barcelona, Catalonia 08005, Spain, ²Centre for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute, Department of Medical Genetics, University of British Columbia, Vancouver, BC V5Z 4H4, Canada, ³Laboratory of Protein Design & Immunoengineering, School of Engineering, Ecole Polytechnique Federale de Lausanne, Lausanne 1015, Vaud, Switzerland, ⁴Department of Biosciences, U Science Tech, Universitat de Vic-Universitat Central de Catalunya, Vic, Catalonia 08500, Spain and ⁵Institute of Biological, Environmental and Rural Science, Aberystwyth University, Aberystwyth SY233EB, United Kingdom

Received January 13, 2020; Revised May 07, 2020; Editorial Decision June 04, 2020; Accepted June 10, 2020

ABSTRACT

Cis2-His2 zinc finger (C2H2-ZF) proteins are the largest family of transcription factors in human and higher metazoans. To date, the DNA-binding preferences of many members of this family remain unknown. We have developed a computational method to predict their DNA-binding preferences. We have computed theoretical position weight matrices (PWMs) of proteins composed by C2H2-ZF domains, with the only requirement of an input structure. We have predicted more than two-third of a single zinc-finger domain binding site for about 70% variants of Zif268, a classical member of this family. We have successfully matched between 60 and 90% of the binding-site motif of examples of proteins composed by three C2H2-ZF domains in JASPAR, a standard database of PWMs. The tests are used as a proof of the capacity to scan a DNA fragment and find the potential binding sites of transcription-factors formed by C2H2-ZF domains. As an example, we have tested the approach to predict the DNA-binding preferences of the human chromatin binding factor CTCF. We offer a server to model the structure of a zinc-finger protein and predict its PWM.

INTRODUCTION

Despite that physical interactions of transcription factors (TFs) with DNA do not always confer a regulatory consequence (1,2), their identification and the characterization

of binding sites is still key to understand how gene expression is regulated. Experimental techniques, such as ChIP (3), PBM (4), HT-SELEX (5), MPRA (6) or bacterial and yeast one-hybrid (7,8) have allowed the characterization of TF-binding sites at large-scale. However, experimental techniques are expensive and time consuming, and yet the binding preferences of many TFs remain unknown (9,10). Given the current limitations, the usage of computational tools to complement experimental techniques is necessary.

Cis2-His2 zinc finger (C2H2-ZF) proteins are the largest family of TFs in higher metazoans (11). They represent around the 45% of all known human TFs, being the largest TF family in humans (10). C2H2-ZF proteins are involved in a wide range of biological processes, such as development (12) or chromatin compartmentalization (13). C2H2-ZF proteins have been related to many diseases (14,15) and can be used as tools for precise gene editing (16,17). At this point, knowing the binding preferences of C2H2-ZF proteins becomes crucial, despite for many are yet unknown (10). Besides, many members of the C2H2-ZF do not have close homologs across metazoans and thus, sequence homology cannot be used to infer their binding preferences (18). Still, all members of this family have the same structure in the DNA binding domain. DNA binding domains (DBD) of C2H2-ZF proteins are composed by small domains called zinc fingers arranged in tandem (19). Each zinc finger is able to recognize DNA sequences of 3 nt (20) and, by combining adjacent zinc fingers, C2H2-ZF proteins are able to recognize long and complex DNA patterns (21). Human C2H2-ZF proteins contain an average of around 10 domains, leading to binding sites of about 30 bases (22).

*To whom correspondence should be addressed. Tel: +34 933160509; Fax: +34 933160550; Email: baldo.oliva@upf.edu

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Several computational tools have been developed to predict the binding preferences of TFs and in particular C2H2-ZF proteins. Some tools are based on combining experimental data with the structure of the interaction between proteins and DNA. Among them, some approaches use different machine learning algorithms: random forest regression (18), support vector machines (23,24), single layer perceptrons (25), hidden Markov models (26) and other statistical models (27,28); using residue–residue contacts (18,23–24,28) or context dependencies and sequence similarities (25–26,29–30). Other tools are based on the analysis of the structural patterns extracted from protein–DNA interactions (27,31–34), or from the accessibility (35) and the flexibility (36) on the DNA chain. Some of these structure-based tools use statistical potentials (33,34). Statistical potentials (also known as knowledge-based potentials) are scoring functions derived from the analysis of contacts in a set of structures. Statistical potentials have been widely used to evaluate the quality and the stability of protein folds, protein–protein interactions and protein–DNA interactions (37).

Here we offer a computational tool to predict the binding preferences of C2H2-ZF proteins. We combine experimental bacterial one-hybrid (B1H) data with structural, three-dimensional, information of TF–DNA complexes to derive C2H2-ZF-specific statistical potentials. We use the potentials to predict the position weight matrix (PWM) of the binding-site of C2H2-ZF domains (38,39). We also predict PWMs for proteins with three C2H2-ZF domains and compare them with their motifs in the JASPAR database (40). We apply statistical potentials to predict the binding preferences of human CTCF, a transcriptional repressor with a key role in genome compartmentalization (13).

MATERIALS AND METHODS

Software

The following software is used in this study: DSSP (version CMBI 2006) (41) to obtain protein structural features; X3DNA (version 2.0) (42) to analyze and generate DNA structures; *matcher* and *needle*, from the EMBOSS package (version 6.5.0) (43), to obtain local and global alignments, respectively; BLAST (version 2.2.22) (44) to search homologs of a query (target) protein sequence; MODELLER (version 9.9) (45) to construct structural models; and the program TOMTOM from the MEME suite (46) to compare two PWMs.

Databases

Atomic coordinates of protein–DNA complex structures are retrieved from the Protein Data Bank (PDB) repository (47). Protein codes and sequences are extracted from UniProt (January 2019 release) (48). We generate an internal database of protein–DNA structures with all C2H2-ZF proteins as identified in the CIS-BP database (version 1.62) (9). DNA-binding information of the C2H2-ZF family is retrieved from a bacterial one-hybrid (B1H) experiments (21). The experiment distinguishes between the inner (F2) and C-tail (F3) C2H2-ZF domains. The experiment tests the interactions of all 64 possible combinations of 3 base-pairs

(bp), which is the characteristic binding site length of individual C2H2-ZF domains (20), against multiple large protein libraries based on Zif268, with six variable amino acid positions in domains F2 and F3 (49).

Interface and triads of protein–DNA structures

We define *triads* as a type of contacts between the protein and the double-strand DNA helix. Triads consist of one amino acid and two contiguous nucleotides (i.e. dinucleotide) of the same strand. The distance associated with a triad is defined by the distance between the C_β atom of the amino acid residue and the average position of the atoms of the nitrogen-base of the 2 nt plus their complementary pairs in the opposite strand of the DNA helix (37). Each triad has an associated amino acid residue number in the protein and a dinucleotide position in the DNA. For instance, a *triad* with amino acid residue number p , dinucleotide in position q and associated distance d is represented as $(triad, d, p, q)$. Specific environmental features can be added on a *triad*, defining an *extended-triad* (*etriad*). For the amino acid, these features are: hydrophobicity, surface accessibility and secondary structure (determined with DSSP); and for the dinucleotide: types of nitrogenous bases, the closest strand, the closest groove and the closest chemical group to the amino acid.

Statistical potentials

We use the definition of **statistical potentials** described by Feliu *et al.* (50) and Fornes *et al.* (37) to define several **scoring functions** for the interaction between a protein and a DNA binding site. To calculate the statistical potentials, we use the distribution of triads at distances up to 30 Å to calculate the statistical potentials. The total potential of an interaction is calculated as the sum of the potentials of all triads, or triads with extended environmental features (*etriads*). In the case of *etriads*, the completeness of the reference dataset, basically C2H2-ZF/DNA complex structures from the PDB, is not sufficient to sample all possible combinations. We use interactions from B1H to extend the number of interacting triads (see further details in Supplementary Methods). Besides, we transform the statistical potentials into *Z*-scores (described below), to simultaneously identify the best distance associated with a triad and the best pair consisting on an amino acid and a dinucleotide.

Z-scores

The optimal condition of a statistical potential often yields a minimum. However, this minimum does not necessarily have to be negative. The variability of signs of the potentials affects the quality criterion of the scores. We define *Z*-scores to follow a criterion that incorporates signs. The goal is that the *Z*-score simultaneously identifies the best distance associated with a triad and the best pair consisting of an amino acid and a dinucleotide. Consequently, we construct a ***z-score function*** for any type of *score* using a standard normalization with respect the average of all amino acid types (see details in Supplementary Data).

Structural modeling of C2H2-ZF complexes

We obtain the structure of a complex by means of homology modeling using the program MODELLER (45). The DNA binding sequence of a C2H2-ZF protein composed by three zinc-finger domains has a length of 9 bp (e.g. for Zif268). We use the same Zif268 sequences as in the B1H experiment (21). For each Zif268 sequence, its complex structure with DNA is modeled with 23 different templates (see details in the supplementary extension of methods and in http://sbi.upf.edu/C2H2ZF_repo). We complete each complex by modeling the structure of the bound DNA sequence with the program X3DNA (42). The full DNA sequence of the experiment is 29 bp long. We embed the binding site in positions 11–19. We also model several structures with the complex of Zif268 binding a non-specific DNA region to be used as non-binding examples (or background). The non-binding sequence is obtained by randomly selecting a 9 bp region of the weak promoter GAL1 (see details in extended Supplementary Methods).

Use of experimental TF–DNA interactions to calculate statistical potentials

We use a mapping function that associates the amino acids of each hexamer core sequence from the of B1H experiment with the amino acids of a template structure to derive interacting *etriads*. We use a similar mapping for the nucleotides. For each C2H2-ZF domain (F2 and F3) and combination of 3 nt, we collect all hexamer core sequences producing a significant binding signal in the B1H experiment (see details in Supplementary Data). Only triads affecting the amino acids and nucleotides under test are considered for the calculation of statistical potentials. We restrict each set of Zif268 sequences to those with the highest signal from the B1H binding experiment. Specifically, we define three thresholds based on the affinity percentile between a hexamer core sequence and a DNA: (i) higher than 90%; (ii) higher than 75%; and (iii) higher than 50%. Affinity percentiles are calculated as in the original publication (21). We impose around 500 DNA sequences per hexamer. A DNA binding sequence is repeated proportionally to the number of observations in the B1H experiment. The contacts derived from the B1H experiment are limited to short distances (the furthest contacts are around 15–20 Å). We include contacts extracted from other C2H2-ZF/DNA structures in the PDB to cover distances of up to 30 Å.

Scoring TF–DNA interactions

First, we calculate the TF–DNA interface and extract all *etriads* at distances shorter than 30 Å. Then, the score of the interaction is defined as the sum of the scores (i.e. a specific statistical potential) of all *etriads* with their associated distances. The same is applied for Z-scores. Provided that it can be modeled, we can obtain the score of a C2H2-ZF TF from its sequence alone (see details in supplementary).

Construction of PWMs using Zif268 structural models

Given the modeled structure of a C2H2-ZF/DNA complex, we obtain the PWM using the Z-score of *ES3DC_{dd}*

(*ZES3DC_{dd}*, as defined in supplementary). We extract the set of *etriads* up to a maximum of 30 Å, with their associated distances, amino acid and dinucleotide positions. We create a **test set** with all possible nucleotide sequences of the same length as the DNA molecule of the structure. We calculate the score of every sequence of the test set (see details in supplementary describing a heuristic approach for sequences longer than 9 bp), and normalize the scores as follows:

$$\text{normal}(\text{score}_{\text{seq}}) = \frac{\text{score}_{\text{seq}} - \min(\{\text{score}_{\text{seq}}\})}{\max(\{\text{score}_{\text{seq}}\}) - \min(\{\text{score}_{\text{seq}}\})}$$

Where $\text{score}_{\text{seq}} = -ZES3DC_{dd}$ and $\{\text{score}_{\text{seq}}\}$ is the set of all scores in the ‘test set’. Normalized scores range between 0 and 1. Then, we rank the normalized scores and select the top scoring DNA sequences over a **cut-off threshold** (i.e. 0.95). The selected sequences are used to build an ungapped MSA, which is used to calculate the **theoretical PWM** of the TF.

Construction of the experimental PWM

The experimental PWM of a Zif268 sequence from the B1H experiment is calculated based on its affinities for different binding sites. The DNA strand of a binding site is formed by trinucleotides flanked by two fixed nucleotides (G and A for F2, and two A for F3). All binding sites targeted by a specific hexamer-fragment with affinity higher than a threshold are stored and gapless aligned without gaps to construct the PWM (e.g. the top 20% threshold uses all binding sites with affinity percentile higher than 80%, while for a threshold of 100% we use all detected sites with an affinity percentile that is not null). We construct experimental PWMs for top 10%, top 25%, top 50% and top 100% binding sites. These experimental PWMs are also named **hexamer-specific PWMs**, to distinguish them from PWMs obtained with other experiments or with a different approach.

RESULTS

Outline of the method

The main objective of this work is to predict the PWM of a TF of the C2H2-ZF family from a complex structure of the TF/DNA interaction. Briefly, we score the interaction with statistical potentials, a mathematical formulation of the observed frequency of contacts between amino acids and dinucleotides (i.e. named triads). We use these potentials to rank DNA sequences potentially bound by the TF and use the top scoring sequences to calculate its PWM.

To calculate the statistical potentials we require a large population of contacts, otherwise the potentials may be too sparse. To overcome this limitation, we have developed a computational approach that increases the amount of available contacts with non-structural experimental information gathered for the C2H2-ZF family. Figure 1 shows a flowchart of the method, from the calculation of statistical potentials (at the top) to the prediction of the PWM (at the bottom). The first step (1) is to calculate the contacts between amino acids and dinucleotides of all non-redundant TF–DNA complex structures of the C2H2-ZF family in PDB, as in Fornes *et al.* (37). Contacts and their

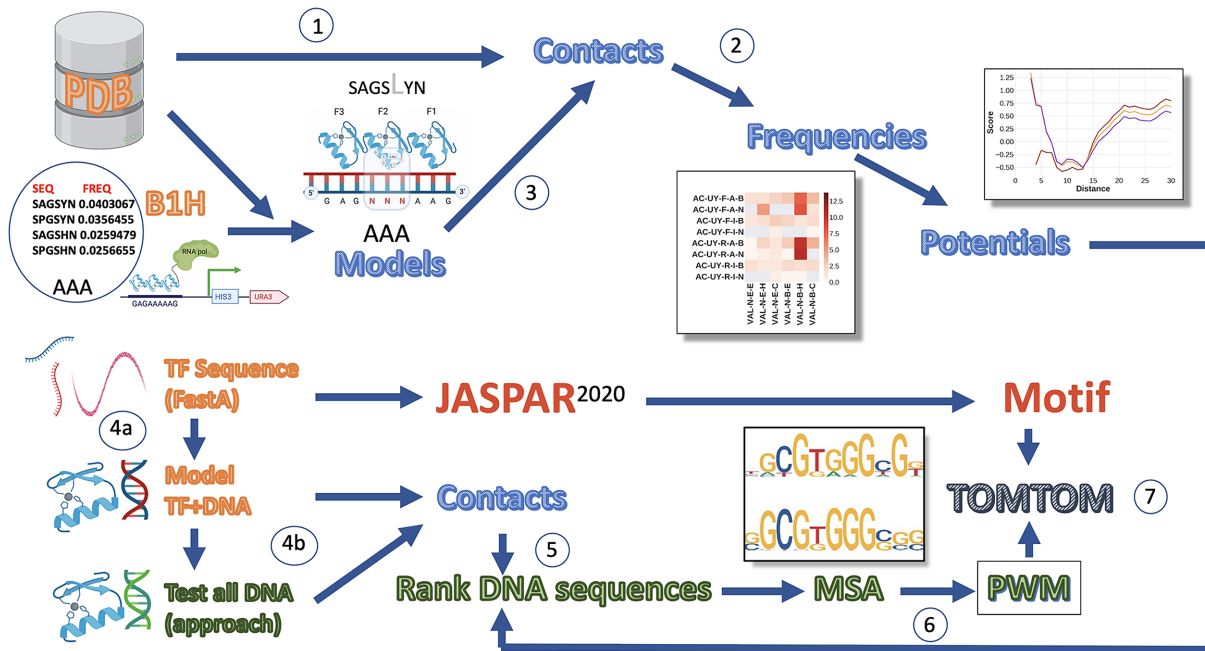


Figure 1. Flowchart of the approach. The flowchart follows three main tasks: i) steps 1–3 to generate the statistical potential using information from Protein Data Bank (PDB) and experiments of bacteria-one-hybrid (B1H); ii) steps 4–6 to build the PWM of a sequence or a structure of a TF of the C2H2-ZF family; and iii) step 7 to compare the PWM of a TF of the C2H2-ZF family with its motif in JASPAR. Steps are numbered and encircled, following the description in the main text. Input data for the approach is colored in orange, scripts for the construction of the statistical potential in blue and scripts for the construction of the PWM in green.

specific features are stored as triads and their frequencies are used to calculate the statistical potentials (step 2). To increase the amount of available contacts, we transform the bacterial-one-hybrid (B1H) data from Persikov *et al.* (21) as follows (step 3): first we select the hexamer core amino acid sequences of Zif268 C2H2-ZF domains 2 and 3 (F2 and F3, respectively) that bind a specific 3 bp of DNA with sufficient affinity. Second, for each hexamer, we model a balanced set of contacts by substitution of the hexamer amino acids and its DNA bound sequence(s) in structural models of Zif268 (see details in ‘Materials and Methods’ section and Supplementary Data). Finally, the resulting contacts are used to recalculate the frequencies and statistical potentials from step 2.

To predict the PWM of a C2H2-ZF protein we require either its sequence or its complex structure with DNA. If starting from the sequence (step 4a), we model the complex structure of its interaction with DNA. Depending on the number of templates used, this step may produce multiple structural models. Each generated model follows step 4b. If starting from the structure (step 4b), we replace the DNA of the structure by all possible DNA sequences of the same length (we use a heuristic approach for long sequences to reduce the computational cost, see supplementary) and apply the statistical potentials to score the contacts of the TF with each replaced DNA. Scores are normalized between 0 (for the worst binding) and 1 (for the best binding). In step 5, we rank all the replaced DNA sequences by the normalized scores and select the top ones (with normalized score higher than 0.95) to build an ungapped multiple sequence alignment (MSA). In step 6 the MSA is converted into a

PWM (in MEME format). To validate the approach, we perform the comparison with experimental PWMs derived experimentally. Specifically, out of 40 TFs from JASPAR (40) with three C2H2-ZF domains, we model 29 with multiple templates (step 4a). Then, we predict the theoretical PWM of each structure (step 4b). We use TOMTOM (46) to compare the JASPAR PWM of a TF with the theoretical PWMs obtained for the different models of that TF (step 7).

Analysis of the statistical potentials

We have constructed several statistical potentials to describe the interaction between the C2H2-ZF domains (F2 and F3 of Zif268) and the DNA. We have applied a Z-score modification (see ‘Materials and Methods’ section and further details in Supplementary Data) on top of the classical definition of potential (51) (named PAIR). As an example of statistical potential, we have selected the interactions between asparagine (Asn) with dinucleotide guanine-cytosine (GC) and arginine (Arg) with dinucleotides adenine-guanine (AG) and cytosine-thymine (CT). We selected these two residues because Arg is a classical amino acid positively charged found involved in unspecific protein DNA contacts, while Asn is a polar residue with specificity for some nitrogenous bases. Figure 2A–D shows the PAIR and ZPAIR potentials between Asn and the dinucleotide with bases GC in finger domains F2 and F3. This example shows that the Z-score function preserves the optimum shortest distance, but different between domains F2 and F3. Figure 2E–H shows the ZPAIR potential of Arg interacting with dinucleotides with bases AG and CT. Supplementary fig-

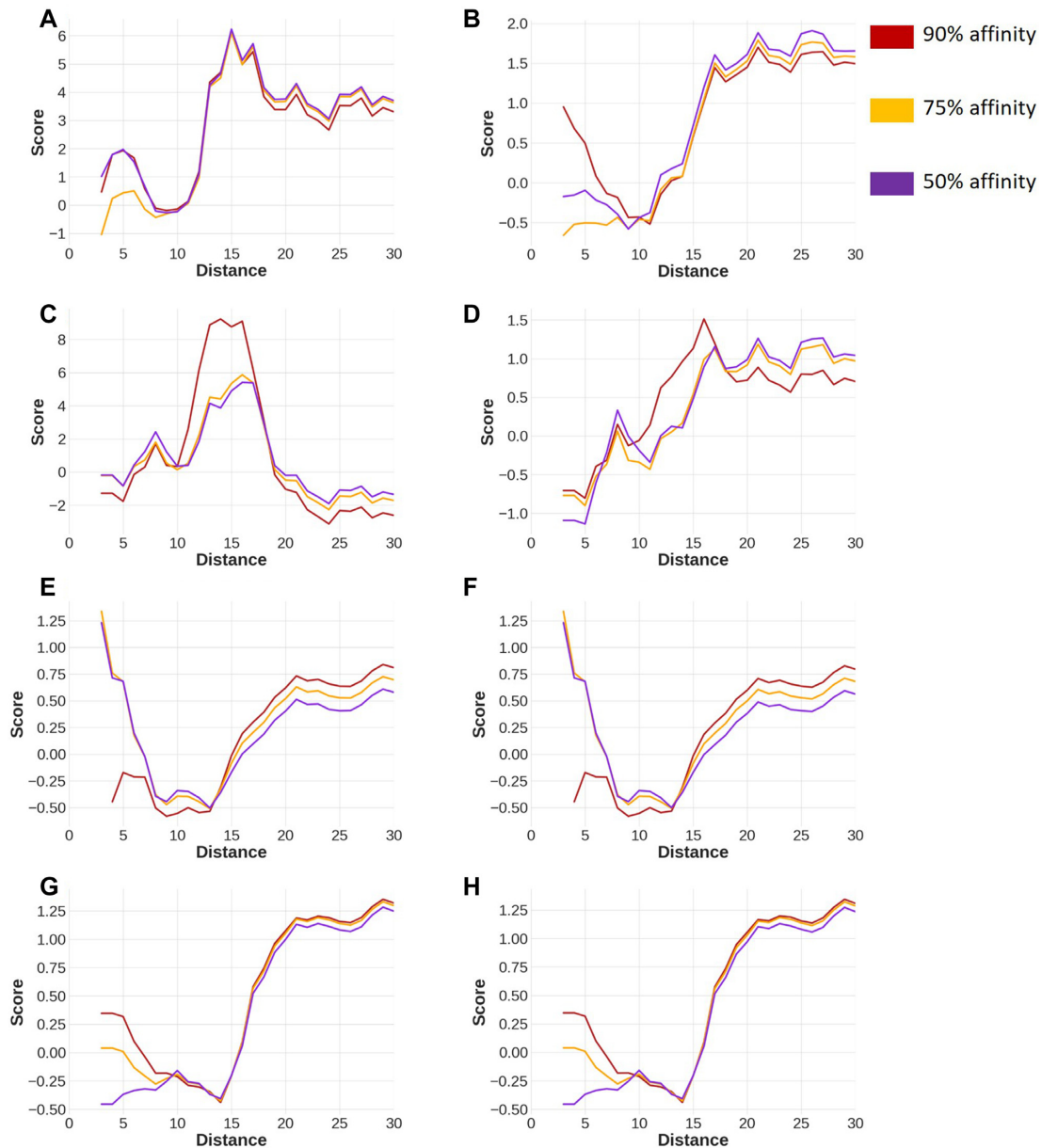


Figure 2. Statistical energy profiles PAIR and ZPAIR obtained with F2 and F3 domains. (A) Profile of Asn–GC PAIR score using F2. (B) Profile of Asn–GC ZPAIR score using F2. (C) Profile of Asn–GC PAIR score using F3. (D) Profile of Asn–GC ZPAIR score using F3. (E) Profile of Arg–AG ZPAIR using F2. (F) Profile of Arg–CT ZPAIR using F2. (G) Profile of Arg–AG ZPAIR using F3. (H) Profile of Arg–CT ZPAIR using F3.

ures showing the potential PAIR and ZPAIR for all amino acids and dinucleotides can be accessed in http://sbi.upf.edu/C2H2ZF_repo. We use a set of hexamer sequence-fragments yielding affinity percentiles higher than 90, 75 or 50% to construct the potentials. We observe a relatively small difference using either 90, 75 or 50% affinities for distances longer than 10Å. We also observe that the potential is symmetric for the reversed dinucleotide (i.e. the potential resulting for the interaction of Arg with AG in Figure 2E and G is the same with CT in Figure 2F and H). However, the finger-domain has the ability to distinguish forward and reverse dinucleotides depending on structural and topological features of the DNA helix. In previous works we already developed a topological-dependent potential named ES3DC

(see details in Supplementary Methods and in Fornes *et al.* (37)). The limitation of such specific potential is the completeness of the dataset, as the large number of combinations to be sampled is very high and thus requiring a large number of observations. The use of experimental data from B1H is a good opportunity to populate many triads in close distance (shorter than 20 Å) between the finger domain and the DNA binding site (21,49). Figure 3 shows the increase of different types of contacts produced with the help of B1H data with respect to those obtained only with structures of the C2H2-ZF family in PDB (47). Only some topological features of both DNA and protein conformation highlight the increase, as they are specific of the C2H2-ZF family. As an example, Figure 3 shows the granular-

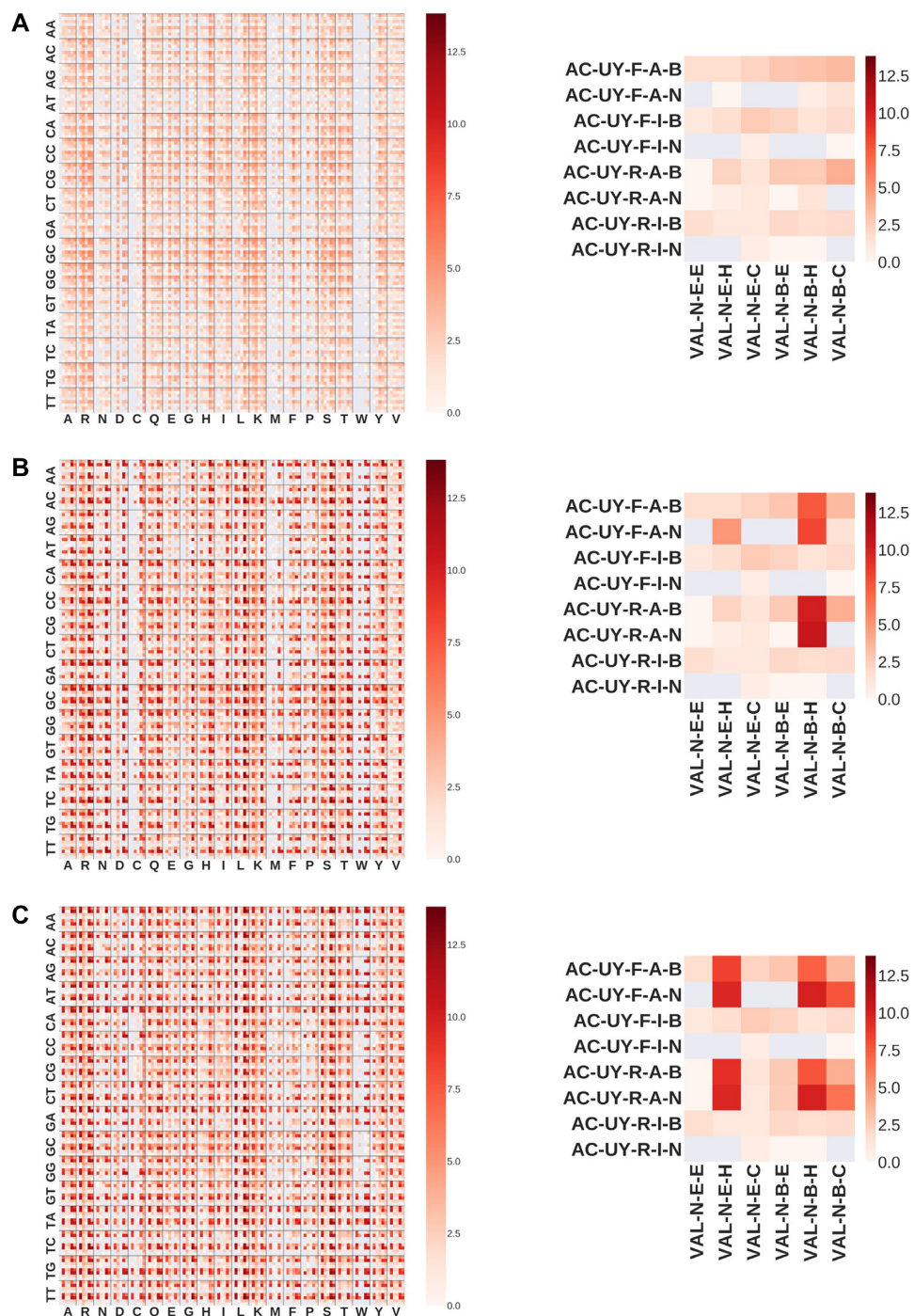


Figure 3. Heatmap plots of the number of amino acid—dinucleotides and their environments (etriads) at distance shorter than 30A in a logarithmic scale. Detailed view of a cell of the heatmap is shown in the right side of each heatmap. Each square inside the cell shows the number of extended-triads (in logarithmic scale) for a specific amino acid—dinucleotide (the example uses valine, Val and adenosine-cytosine, AC) and their environments. Amino acid environments are: hydrophobicity (P as polar, N non polar), surface accessibility (E if exposed, B if buried) and secondary structure (E for β -strand, H for helix and C for coil). Dinucleotide environments are: type of nitrogenous bases (U for purine, I for pyrimidine), closest DNA strand (F for forward, R for reverse), closest DNA groove (A for major, I for minor) and closest chemical group (B if phospho-ribose backbone atoms, N if nucleobase). (A) Extended-triads obtained from PDB structures. (B) Extended-triads obtained from PDB structures and B1H experiments of the F2 domain. (C) Extended-triads obtained from PDB structures and B1H experiments of the F3 domain. A detailed analysis of panels A, B and C is available in the web repository (http://sbi.upf.edu/C2H2ZF_repo).

ity of features covering the interaction between valine and adenosine-cytosine. The conclusions are similar for other interactions: (i) the B1H data help to populate certain specific features; and (ii) the specificity is derived mostly from B1H data rather than from PDB structures (this is caused by the large amount of information extracted from the experiment in comparison to the amount of structures available).

Prediction of PWMs in domains F2 and F3 of Zif268

To evaluate the quality of the theoretical PWMs we compare them with the results of B1H experiments (21). We construct two types of PWMs: hexamer-specific PWMs and trinucleotide-specific PWMs. Hexamer-specific PWMs are the experimental PWMs as defined in ‘Materials and Methods’ section, obtained by aligning all DNA trinucleotides targeted by the same amino acid hexamer. Trinucleotide-specific PWMs are artificial PWMs containing the combination of 3 nt (64 in total) with 100% weight in each specific position. These trinucleotides are flanked by nucleotides specific of each assayed zinc finger (G and A for F2 and two A for F3).

We create the hexamer-specific and theoretical PWMs with the DNA binding sites of the hexamer sequences tested by B1H. For each hexamer fragment of a finger domain, we select the binding site with highest affinity in the B1H experiment and assume that these are the three bases binding-specific of the hexamer. Theoretical PWMs are obtained by combining homology modeling and Z-scores $ZES3DC_{dd}$. Structural models of the variants of Zif268 are constructed using 23 different templates (see Supplementary Data). Hence, there are 23 theoretical PWMs for each hexamer sequence of amino acids.

When comparing the theoretical PWM with the trinucleotide-specific PWMs we check the ranking position of the correct trinucleotide-specific PWM. When comparing with the experimental PWMs we use the number of nucleotide-matches to evaluate the quality. We define the number of nucleotide-matches as the number of positions where two compared PWMs share the same nucleotides with highest frequencies. Since we are focused on the trinucleotides of the binding site, we are only interested in the number of central nucleotide matches, which ranges between 0 and 3.

Using an affinity threshold of 90%, we find 131 hexamers for F2 and 82 for F3 with at least one theoretical PWM ranking the correct binding site on the top. This represents at least one hexamer sequence in 28 (for F2 domain) and 32 (for F3 domain) trinucleotide combinations out of 64. Table 1 shows the number of hexamers with three, 2 or 1-nt-matches with the hexamer-specific PWM for each trinucleotide in F2 domain. Supplementary Table S1 shows the results for domain F3 and additional details. Considering 3 or 2-nt matches, we are able to match at least one theoretical PWM for 71% of the hexamer variants in F2 and 74% in F3. This proves that for most hexamers we are able to find a theoretical PWM with an almost perfect match with the corresponding hexamer-specific PWM. However, the selection of the template (or templates) is crucial to predict the specific binding site of a single domain. The comparisons of PWMs of all hexamer sequences tested in F2 and F3 do-

main with affinity percentile higher than 90% is shown in Supplementary Table S2.

The results obtained using affinity percentile higher than 75 and 50% are similar to using 90%. Around 65–75% hexamers of F2 domain (and 75–85% for F3 domain) have 3 or 2 nt-matches between the theoretical and the hexamer-specific PWM. Details are shown in Supplementary Tables S1 and 2. The comparison of all PWMs are shown in http://sbi.upf.edu/C2H2ZF_repo. The fact that the quality of the theoretical PWMs is the same for the range of affinity percentiles analyzed (from 50 to 90%) suggests that our method is not able to distinguish high from low affinity binding sites; but it allows to identify the TF binding site regardless of the affinity (see Annexure 1 in Supplementary Material and Figure S1).

In Figure 4 we show the comparison of some examples of hexamer-specific and theoretical PWMs (all theoretical and experimental PWMs and their comparisons can be retrieved from http://sbi.upf.edu/C2H2ZF_repo). Among these examples we observe some theoretical PWMs that, although different than their expected binding sites, share common trends of the nucleotide frequencies of the experimental PWM. For example, for the binding site ATG in domain F2 by the SQSGCN hexamer (top left PWM in Figure 3A), we observe similar nucleotides underlying lower frequencies between theoretical and experimental PWMs. Similarly, other examples are shown in Figure 4 with combinations of nucleotides of binding sites displaying nucleotide matches with underlying lower frequencies.

Examples of binding site prediction of C2H2-ZF transcription factors

We compare the theoretical PWMs with the PWMs retrieved from JASPAR (40) for several TFs. We use some members of the C2H2-ZF family, composed by three finger-domains, with a known PWM (coded as motif) in JASPAR (40), for which the structure of the complex with DNA is known or it can be modeled, to obtain the theoretical PWM (see Supplementary Table S3 and other details in Supplementary Material). We obtain two PWMs using statistical potentials $ZES3DC_{dd}$ calculated with variant sequences in F2 domain (**ZES3DCF2**) and in F3 (**ZES3DCF3**) of the B1H experiment. We compare the theoretical PWMs of each TF using all contacts under 30Å, then we repeat the comparison by decreasing this threshold down to 15Å.

Figure 5 shows the JASPAR PWMs of some selected TFs, compared with the theoretical PWMs calculated with a distance threshold of 30 Å (see Supplementary Table S4 for more details). All theoretical PWMs and structural models can be downloaded from http://sbi.upf.edu/C2H2ZF_repo. We are able to find at least one PWM significantly similar to its motif in JASPAR (P -value < 0.05 with TOMTOM) for almost all TFs (27 out of 29). The PWMs of some TFs are compared with more than one possible motif in JASPAR, often associated by some relationship in evolution (i.e. among orthologs and paralogs of different species, see details in Supplementary Data).

We further test if the similarity of the TFs with the sequence of Zif268, from which the statistical potentials are derived, affects the quality of the results. We calculate the

Table 1. Results of the prediction of PWMs in domain F2

BS	3M	2M	1M	#HEXAMER	TEMPLATE	HEXAMER	<M>	RANK	TOP
AAA	0	3	1	4	5ke9_A	SPGSHN	0.65	2	0
AAC	1	7	24	32	2wbu_A	WHSSVH	0.66	1	1
AAG	0	5	10	15	5ke9_A	RSDYTM	0.55	4	0
AAT	4	3	1	8	1g2d_C	FQSNVS	0.39	1	3
ACA	1	8	16	25	1ali_A	QQSTSR	0.68	4	0
ACC	0	8	0	8	5ke8_A	HPSTSH	1.04	4	0
ACG	4	21	0	25	1alj_A	WASSSN	0.80	1	10
ACT	3	37	10	50	1p47_A	FSSSSA	0.79	1	3
AGA	0	2	0	2	1zaa_C	SSGSWN	0.70	1	1
AGC	7	4	0	11	5ke8_A	WHSSIH	0.74	1	6
AGG	0	22	7	29	5ke6_A	RKDHTN	0.82	3	0
AGT	0	8	14	22	1ali_A	YHSNLS	0.31	2	0
ATA	0	1	0	1	1ali_A	NAHNCL	0.37	9	0
ATC	1	6	5	12	5ke7_A	SSSGLH	0.68	1	1
ATG	0	1	11	12	4r2c_A	WHSGLN	0.40	9	0
ATT	0	7	11	18	5keb_A	FQSGSS	0.29	1	1
CAA	0	1	1	2	1zaa_C	TKGNTQ	0.57	3	0
CAC	0	11	0	11	5keb_A	DPSNRS	0.94	2	0
CAG	0	15	16	31	1p47_A	TKWNTS	0.75	2	0
CAT	3	4	5	12	5ke7_A	AQSNSS	0.57	1	3
CCA	0	3	0	3	5keb_A	QLSTNY	0.80	4	0
CCC	3	0	0	3	1zaa_C	TRRDRR	2.42	1	3
CCG	1	1	0	2	1zaa_C	RKDTRD	1.78	1	1
CCT	0	7	0	7	1zaa_C	RKQDSR	1.25	1	1
CGA	2	1	0	3	2wbu_A	QYGHST	0.77	1	2
CGC	0	3	0	3	5keb_A	SRPNLG	1.59	2	0
CGG	22	8	1	31	5keb_A	RASHSD	1.43	1	23
CGT	0	20	3	23	5keb_A	MSHHRD	0.99	2	0
CTA	1	2	0	3	4r2c_A	SQSGCQ	0.78	1	1
CTC	3	0	0	3	1p47_A	SRSGCH	1.04	1	3
CTG	0	4	1	5	1p47_A	RKFIIE	0.79	2	0
CTT	0	1	7	8	5ke9_A	YRHVSD	0.76	1	1
GAA	1	0	9	10	1jk1_A	TKGNTR	0.56	2	0
GAC	1	17	1	19	1zaa_C	WASSSR	0.95	3	0
GAG	0	35	3	38	5keb_A	TRFNLR	0.78	1	1
GAT	0	8	1	9	4r2a_A	FASNRR	0.65	2	0
GCA	1	2	0	3	5ke9_A	QLATNR	0.97	3	0
GCC	7	0	0	7	5keb_A	WLTNRR	2.21	1	7
GCG	9	15	0	24	5ke9_A	RRDTAN	1.41	1	20
GCT	11	4	0	15	5ke8_A	FRSTSR	0.97	1	11
GGA	0	6	2	8	1ali_A	QLSTKY	0.74	4	0
GGC	4	1	0	5	2kmk_A	WQSSIK	1.10	1	1
GGG	19	27	0	46	5ke7_A	RNAHLN	1.32	1	20
GGT	0	8	3	11	5ke7_A	FQSNLR	0.84	1	1
GTA	0	2	2	4	1alh_A	TKGSTR	0.71	7	0
GTC	0	7	0	7	2kmk_A	HASSSR	0.84	5	0
GTG	0	38	4	42	5ke9_A	RKAITD	0.87	4	0
GTT	0	5	1	6	5kea_A	FLSSSR	0.72	1	2
TAA	1	0	1	2	1zaa_C	MYIDYY	0.91	1	1
TAC	0	2	3	5	5ke7_A	LKGNTK	0.71	7	0
TAG	2	1	5	8	5ke9_A	RKWTDL	0.53	1	2
TAT	0	6	4	10	5ke9_A	WLTSNV	0.26	9	0
TCA	0	0	2	2	5ke7_A	HNIYHH	0.37	23	0
TCC	0	6	0	6	5ke8_A	TKASTP	1.26	6	0
TCG	1	3	0	4	1alh_A	RKESVI	1.31	5	0
TCT	0	7	4	11	5keb_A	WSSSAI	0.81	3	0
TGA	1	1	1	3	5ke9_A	WASSHY	0.49	1	1
TGC	0	2	0	2	2kmk_A	WPNSKA	0.78	2	0
TGG	0	37	11	48	1alk_A	RNAHSE	0.81	3	0
TGT	0	17	23	40	5ke9_A	WASSSS	0.27	5	0
TTA	0	1	0	1	5keb_A	CIHYNN	0.35	17	0
TTC	0	2	0	2	5ke7_A	SASGSH	0.62	3	0
TTG	0	5	6	11	5keb_A	RKWTML	0.69	2	0
TTT	0	0	3	3	5ke9_A	YRWIRD	0.36	4	0

BS is the trinucleotide combination of the DNA binding site. 3M, 2M and 1M show the number of hexamers with at least one theoretical PWM having 3, 2 or 1 nt matches with the experimental PWM, respectively. #HEXAMER is the total number of hexamers having as main binding the trinucleotide of the row. HEXAMER and TEMPLATE show the hexamer sequence and the code of the structure used as template for the theoretical PWM, this combination yields the highest match of nucleotides of the corresponding trinucleotide in the same row. <M> shows the average ratio of nucleotide matches of all theoretical PWMs and hexamer sequences with the same binding site of the row. RANK shows the best ranking position of the correct trinucleotide-specific PWM among all hexamers with the same binding site of the row. TOP shows the number of hexamers with at least one theoretical PWM ranking on the top the correct trinucleotide-specific PWM of the row.

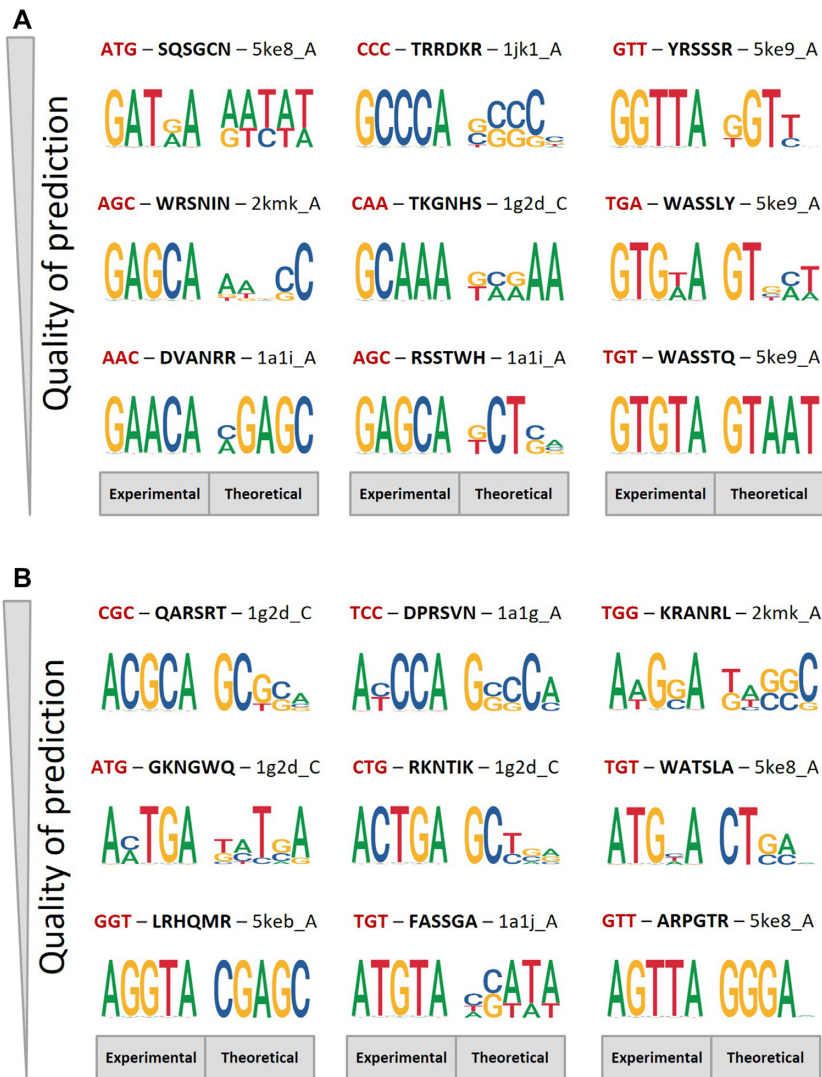


Figure 4. Comparison of PWMs. We compare theoretical PWMs with experimental PWMs of the same hexamer sequence variants. For each comparison we show the amino acid hexamer sequence (highlighted in bold) used to calculate the experimental PWM, the DNA binding site with highest affinity of the hexamer sequence (highlighted in red) and the PDB code of the structure used as template to obtain the theoretical PWM. (A) Comparison of PWMs for domain F2. (B) Comparison of PWMs for domain F3.

similarity as the percentage of identical residues aligned (%id) between the sequence of the TF and the sequence of Zif268. Certainly, the theoretical PWMs of TFs very similar to Zif268 are significantly similar to their motif in JASPAR. However, we also obtain theoretical PWMs for sequences with low similarity with Zif268 that significantly match with their corresponding motif in JASPAR (see Supplementary Data). Similar conclusions are obtained when comparing the sequences of the TFs and the templates used to construct their PWMs. The bias on the statistical potential, caused by structures of close homologs to each TF query is studied in the supplementary. The main conclusion is that the bias is avoided by using contacts shorter than 18Å to construct the theoretical PWMs. We test for each TF the capacity to predict a motif in JASPAR without biases using a modification of the statistical potentials. We generate specific statistical potentials for each TF by removing the

contacts of close homologs (%id > 50). After avoiding the bias, we are still able to find at least one PWM significantly similar to its motif in JASPAR for almost all TFs. Also, between 11 and 14 TFs have more than 50% of the theoretical PWMs significantly similar with their motif in JASPAR, being most of them the same TFs whether close homologs are removed or not in the statistical potential (see Supplementary Table S4).

Not all models produce PWMs significantly similar with their corresponding motifs. For some TFs this can be explained by the low number of models produced: only one model is constructed with the length of three–four finger domains for Q86T24 and Q8GYC1. A detailed analysis shows that many theoretical PWMs, not significantly similar with their motif, are still able to match more than 50% nucleotide-matches with their JASPAR motif. The average ratio of identical nucleotides using all models varies between

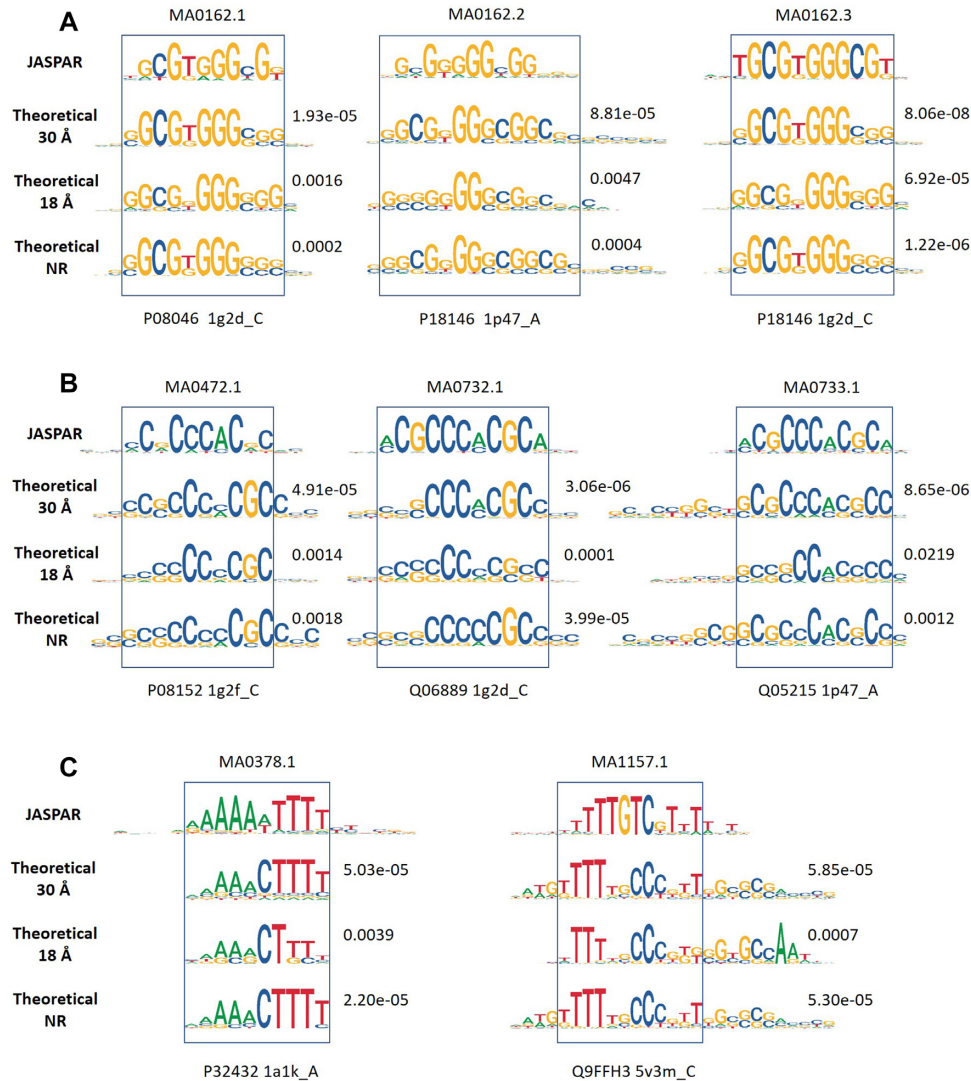


Figure 5. Comparison between theoretical PWMs of some members of the C2H2-ZF family and their motifs in JASPAR database. We use ZES3DCF2 statistical potentials for each TF with contacts under 30 Å, using all PDB structures of the C2H2-ZF family or avoiding those of its close homologs. (A) and (B) show examples rich on G and C nucleotides, while in (C) are shown examples rich on A and T nucleotides. JASPAR motifs are shown at the top of each comparison. PDB codes of the templates used to construct the theoretical PWMs are also indicated. As a quality criterion of the comparison, the right top margin of each predicted logo includes the *P*-value calculated with TOMTOM. (A) Examples of TFs P08046 and P18146. (B) Examples of TFs Q06889, P11161, Q05215 and P08152. (C) Examples of P32432, Q9FFH3 and Q8H1F5.

60 and 88% for the majority of TFs (see Supplementary Table S4), and it is only slightly reduced after avoiding the bias.

Application to CTCF

We apply statistical potentials to predict the binding preferences of human CTCF. The DNA binding domain of human CTCF is formed by 11 zinc-finger domains of the C2H2 family (residues 266–577). Different DNA binding motifs have been proposed for this domain: one for the central part, flanked by one upstream and one downstream motifs (52). The central part of human CTCF binding domain has a well-defined motif in JASPAR (MA0139.1). The structure of the complete sequence of human CTCF is unknown. However, several structures have been obtained

by crystallography of different fragments bound to DNA (structures with PDB codes 5K5H, 5K5I, 5K5J, 5T00, 5T0U, 5YEL, 5YEH, 5YEF, 5UND and 5KKQ). We construct a structural model of the almost complete sequence of the binding domain of human CTCF (see Supplementary Figure S4). The model is constructed by superimposition of the structures 5T0U (zinc-finger domains 2–7) and 5YEL (zinc-finger domains 6–11), using the overlapping fragment of fingers 6 and 7 and removing the redundant amino acids and nucleotides from 5YEL (amino acid fragment from 455 to 512, highlighted in red in Supplementary Figure S4). Finger domains are shown in the protein sequence alignment and in the alignments of DNA sequences taken from the PDB structures. The C-terminal domains, taken from 5YEL, bind on the 5' region, while the N-terminal domains

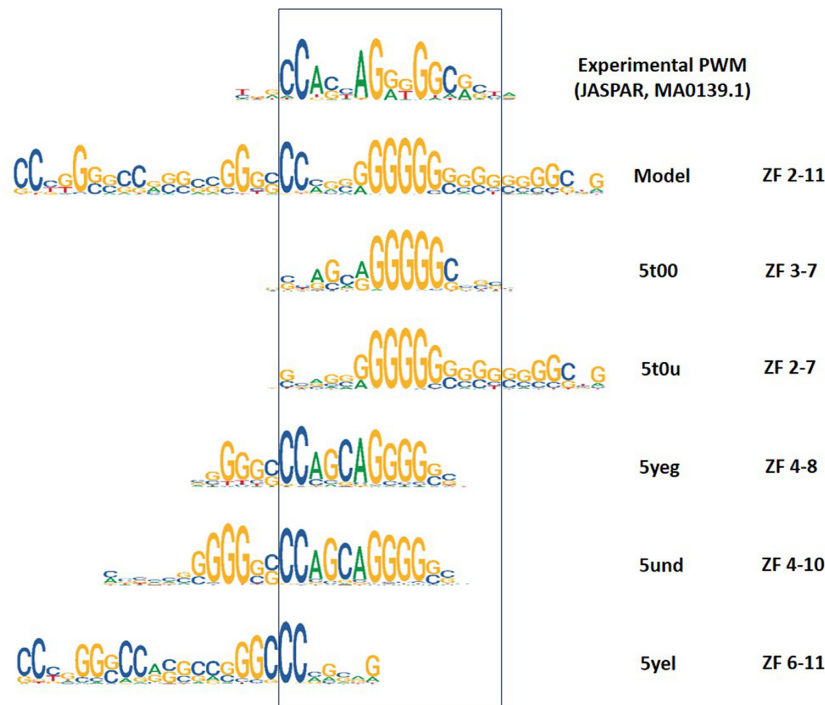


Figure 6. Comparison between experimental and theoretical PWMs of CTCF. The PWM on top of the figure is the experimental PWM, retrieved from the JASPAR database (MA0139.1). The rest of logos show the theoretical PWMs obtained with the model or the structures selected from PDB. The PDB code and the numbers of the zinc-finger domains corresponding to the human CTCF DNA-binding domain are shown on the right side.

from 5T0U bind on the 3' site, thus the alignment of the DNA fragment is shown in reverse orientation with respect to the finger-domains.

Following our previous approaches, we obtain the theoretical PWM with each structure using contacts up to 30Å and the potentials ZES3DCF2 and ZES3DCF3. The PWM based on experimental data is retrieved from JASPAR, with profile motif MA0139.1 and compared with the theoretical PWMs (see details in Supplementary Table S5). We show in Figure 6 the logos of the JASPAR profile aligned with the logos of the theoretical PWMs. We observe a profile pattern preserved for many theoretical PWMs in which we recognize a central common region that corresponds with the JASPAR profile MA0139.1. This is the core profile of CTCF, located between nucleotides 8 and 22 in the forward chain (oriented from 5'-3') of the modeled structure. Because structures, such as codes 5YEL or 5UND are formed mostly by zinc-finger domains at the C-terminal region, the theoretical PWMs constructed with them match incompletely the core profile. Consequently, it is more difficult to align these PWMs, resulting in lower scores (higher *P*-values of TOMTOM) and small overlap. Interestingly, the profiles obtained with PDB codes 5UND, 5YEL and the 5' site of the profile of our model, identify a pattern that has some similarities with the upstream motif mentioned by Nakahashi *et al.* (52). Motifs of flanking sites recognized by finger-domains 1–2 and 8–11 are not well-defined. Despite they are important for the recognition of the binding sites of CTCF, the current tools for motif discovery have not unraveled both downstream and upstream profiles. Consequently, we cannot test the quality of the alignments between the

theoretical PWMs and many of the proposed motifs of the flanking regions, because there is no consensus.

DISCUSSION

We have developed a method to predict the binding preferences of C2H2-ZF proteins using their structures to obtain one (or several) PWMs. We offer the use of this approach with a server. The method requires the structure of the C2H2-ZF protein or the structures of the templates to model it. The number of models depends on the number of templates. Consequently, the number of theoretical PWMs is larger for sequences with many templates than for those with few and this affects the capacity of the prediction.

Our analyses show that the percentage of nucleotide matches in binding sites of single-domains between theoretical and experimental PWMs is independent of the experimental affinity percentile. Therefore, although we can roughly distinguish binding from non-binding sites, we cannot distinguish intermediate degrees of affinity. This is relevant on the prediction of the effect of mutations affecting the binding strength of zinc-finger domains.

Given that our method provides several theoretical PWMs for the same TF, it entails an additional problem: selecting the correct or best PWM. However, rather than finding the best PWM from a set of theoretical PWMs, we bring the opportunity to select one among many potential solutions and help finding the binding site of a TF in a DNA sequence. We proof that, for a relevant number of TFs that can be modeled, the number of PWMs significantly similar to an experimental PWM is larger than 50% (and the

proof is valid after removing biases due to the similarity between the query sequence and the dataset used to construct the prediction). Therefore, by scanning with several theoretical PWMs of a TF, the majority of regions detected and predicted to bind will hit around the right location of the binding site.

Furthermore, our approach also suggests that perhaps the same TF recognizes more than one binding site depending on its conformation. When constructing theoretical PWMs with different structures, each structure is a snapshot of the interaction between the TF and the DNA. Therefore, by using many structures we introduce the dynamic nature of proteins as an additional feature. This can be useful for C2H2-ZF proteins that may interact with DNA using different conformations or different arrangements of zinc fingers, such as CTCF. It is known that CTCF has a central binding motif plus two flanking motifs, one in downstream and another upstream (52). CTCF binding sites display different combinations of downstream-central-upstream motifs that can be spaced by a variable number of nucleotides. Therefore, searching binding sites with many PWMs obtained from different conformations of CTCF–DNA complexes may be more informative of the whole conformational space of CTCF than a single model.

To sum up, we have developed a computational tool to predict the DNA binding preferences of C2H2-ZF proteins. With the help of homology modeling, we are able to predict PWMs for TFs for which we only know their amino acid sequence. We have tested our method by comparing theoretical PWMs with their motifs in JASPAR. We have used our approach to test PWM predictions of different regions of human CTCF and predicted a PWM to cover domains 2 to 11 of the DNA binding domain of CTCF (from downstream to upstream motifs). We offer a repository with the results and a server to calculate the PWM using the structure of a TF as input (see http://sbi.upf.edu/C2H2ZF_repo). We also offer a server to model a TF–DNA complex structure with an input sequence (see details in Supplementary Data). In a thorough analysis of the potential number of C2H2-ZF TFs from Uniprot eukaryotic reference proteomes (one protein per gene) that would benefit from this approach, we have detected 134 399 proteins (about 72%). Thus, our method could potentially be applied to more than 250 000 out of 353 167 UniProt proteins with a predicted C2H2-ZF domain by PROSITE (53) (see more details in Supplementary Data). We think our approach may also be applied to study the potential effect of mutations in the DNA binding sequence. However, because of the lack of specificity on the prediction of binding affinities, further research is still needed on this goal. In the near future, we plan to extend this approach to other TFs with additional experimental information from B1H or other similar experiments collected in Cis-BP database (54).

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

ACKNOWLEDGEMENTS

We thank the help of Drs Persikov and Noyes to access and understand the data on B1H we have used in this

work. B.O., N.F.F. and O.F. wish to thank the Council of the Catalan Republic for their spirit and devotion to catalan science.

FUNDING

Spanish Ministry of Economy (MICINN) [BIO2017-85329-R, RYC2015-17519, MDM2014-0370] and European Regional Development Fund (FEDER) [BIO2017-85329-R, RYC-2015-17519, MDM-2014-0370]; Erasmus+ Fellowship 2019 by EU (to F.Å.); Research Formation of ‘Generalitat de Catalunya’ (FI) Fellowship (to A.M).
Conflict of interest statement. None declared.

REFERENCES

1. Fuxman Bass, J.I., Pons, C., Kozłowski, L., Reece-Hoyes, J.S., Shrestha, S., Holdorf, A.D., Mori, A., Myers, C.L. and Walhout, A.J. (2016) A gene-centered *C. elegans* protein–DNA interaction network provides a framework for functional predictions. *Mol. Syst. Biol.*, **12**, 884.
2. Kemmeren, P., Sameith, K., van de Pasch, L.A., Benschop, J.J., Lenstra, T.L., Margaritis, T., O’Duibhir, E., Apweiler, E., van Wageningen, S., Ko, C.W. *et al.* (2014) Large-scale genetic perturbations reveal regulatory networks and an abundance of gene-specific repressors. *Cell*, **157**, 740–752.
3. Johnson, D.S., Mortazavi, A., Myers, R.M. and Wold, B. (2007) Genome-wide mapping of in vivo protein–DNA interactions. *Science*, **316**, 1497–1502.
4. Berger, M.F., Philippakis, A.A., Qureshi, A.M., He, F.S., Estep, P.W., Lenstra, T.L., Bulyk, M.L. (2006) Compact, universal DNA microarrays to comprehensively determine transcription-factor binding site specificities. *Nat. Biotechnol.*, **24**, 1429–1435.
5. Hallikas, O. and Taipale, J. (2006) High-throughput assay for determining specificity and affinity of protein–DNA binding interactions. *Nat. Protoc.*, **1**, 215–222.
6. Melnikov, A., Murugan, A., Zhang, X., Tesileanu, T., Wang, L., Rogov, P., Feizi, S., Gnirke, A., Callan, C.G. Jr, Kinney, J.B. *et al.* (2012) Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat. Biotechnol.*, **30**, 271–277.
7. Meng, X. and Wolfe, S.A. (2006) Identifying DNA sequences recognized by a transcription factor using a bacterial one-hybrid system. *Nat. Protoc.*, **1**, 30–45.
8. Deplancke, B., Dupuy, D., Vidal, M. and Walhout, A.J. (2004) A gateway-compatible yeast one-hybrid system. *Genome Res.*, **14**, 2093–2101.
9. Weirauch, M.T., Yang, A., Albu, M., Cote, A.G., Montenegro-Montero, A., Drewe, P., Najafabadi, H.S., Lambert, S.A., Mann, I., Cook, K. *et al.* (2014) Determination and inference of eukaryotic transcription factor sequence specificity. *Cell*, **158**, 1431–1443.
10. Lambert, S.A., Jolma, A., Campitelli, L.F., Das, P.K., Yin, Y., Albu, M., Chen, X., Taipale, J., Hughes, T.R. and Weirauch, M.T. (2018) The human transcription factors. *Cell*, **172**, 650–665.
11. Vaquerizas, J.M., Kummerfeld, S.K., Teichmann, S.A. and Luscombe, N.M. (2009) A census of human transcription factors: function, expression and evolution. *Nat. Rev. Genet.*, **10**, 252–263.
12. Sommer, R.J., Retzlaff, M., Goerlich, K., Sander, K. and Tautz, D. (1992) Evolutionary conservation pattern of zinc-finger domains of *Drosophila* segmentation genes. *Proc. Natl. Acad. Sci. U.S.A.*, **89**, 10782–10786.
13. Phillips, J.E. and Corces, V.G. (2009) CTCF: master weaver of the genome. *Cell*, **137**, 1194–1211.
14. Ladomery, M. and Delleire, G. (2002) Multifunctional zinc finger proteins in development and disease. *Ann. Hum. Genet.*, **66**, 331–342.
15. Fuxman Bass, J.I., Sahni, N., Shrestha, S., Garcia-Gonzalez, A., Mori, A., Bhat, N., Yi, S., Hill, D.E., Vidal, M. and Walhout, A.J.M. (2015) Human gene-centered transcription factor networks for enhancers and disease variants. *Cell*, **161**, 661–673.

16. Jabalameli, H.R., Zahednasab, H., Karimi-Moghaddam, A. and Jabalameli, M.R. (2015) Zinc finger nuclease technology: advances and obstacles in modelling and treating genetic disorders. *Gene*, **558**, 1–5.
17. Oakes, B.L., Xia, D.F., Rowland, E.F., Xu, D.J., Ankoudinova, I., Borhardt, J.S., Zhang, L., Li, P., Miller, J.C., Rebar, E.J. *et al.* (2016) Multi-reporter selection for the design of active and more specific zinc-finger nucleases for genome editing. *Nat. Commun.*, **7**, 10194.
18. Gupta, A., Christensen, R.G., Bell, H.A., Goodwin, M., Patel, R.Y., Pandey, M., Enuameh, M.S., Rayla, A.L., Zhu, C., Thibodeau-Beganny, S. *et al.* (2014) An improved predictive recognition model for Cys(2)-His(2) zinc finger proteins. *Nucleic Acids Res.*, **42**, 4800–4812.
19. Wolfe, S.A., Nekludova, L. and Pabo, C.O. (2000) DNA recognition by Cys2His2 zinc finger proteins. *Annu. Rev. Biophys. Biomol. Struct.*, **29**, 183–212.
20. Pavletich, N.P. and Pabo, C.O. (1991) Zinc finger-DNA recognition: crystal structure of a Zif268-DNA complex at 2.1 Å. *Science*, **252**, 809–817.
21. Persikov, A.V., Wetzel, J.L., Rowland, E.F., Oakes, B.L., Xu, D.J., Singh, M. and Noyes, M.B. (2015) A systematic survey of the Cys2His2 zinc finger DNA-binding landscape. *Nucleic Acids Res.*, **43**, 1965–1984.
22. Najafabadi, H.S., Mnaimeh, S., Schmitges, F.W., Garton, M., Lam, K.N., Yang, A., Albu, M., Weirauch, M.T., Radovani, E., Kim, P.M. *et al.* (2015) C2H2 zinc finger proteins greatly expand the human regulatory lexicon. *Nat. Biotechnol.*, **33**, 555–562.
23. Persikov, A.V., Osada, R. and Singh, M. (2009) Predicting DNA recognition by Cys2His2 zinc finger proteins. *Bioinformatics*, **25**, 22–29.
24. Persikov, A.V. and Singh, M. (2014) De novo prediction of DNA-binding specificities for Cys2His2 zinc finger proteins. *Nucleic Acids Res.*, **42**, 97–108.
25. Liu, J. and Stormo, G.D. (2008) Context-dependent DNA recognition code for C2H2 zinc-finger transcription factors. *Bioinformatics*, **24**, 1850–1857.
26. Cho, S.Y., Chung, M., Park, M., Park, S. and Lee, Y.S. (2008) ZIFIBI: prediction of DNA binding sites for zinc finger proteins. *Biochem. Biophys. Res. Commun.*, **369**, 845–848.
27. Kaplan, T., Friedman, N. and Margalit, H. (2005) Ab initio prediction of transcription factor targets using structural knowledge. *PLoS Comput. Biol.*, **1**, e1.
28. Benos, P.V., Lapedes, A.S. and Stormo, G.D. (2002) Probabilistic code for DNA recognition by proteins of the EGR family. *J. Mol. Biol.*, **323**, 701–727.
29. Najafabadi, H.S., Albu, M. and Hughes, T.R. (2015) Identification of C2H2-ZF binding preferences from ChIP-seq data using RCADE. *Bioinformatics*, **31**, 2879–2881.
30. Lambert, S.A., Yang, A.W.H., Sasse, A., Cowley, G., Albu, M., Caddick, M.X., Morris, Q.D., Weirauch, M.T. and Hughes, T.R. (2019) Similarity regression predicts evolution of transcription factor sequence specificity. *Nat. Genet.*, **51**, 981–989.
31. Mandel-Gutfreund, Y. and Margalit, H. (1998) Quantitative parameters for amino acid-base interaction: implications for prediction of protein-DNA binding sites. *Nucleic Acids Res.*, **26**, 2306–2312.
32. Contreras-Moreira, B. (2010) 3D-footprint: a database for the structural analysis of protein-DNA complexes. *Nucleic Acids Res.*, **38**, D91–D97.
33. Lin, C.K. and Chen, C.Y. (2013) PiDNA: Predicting protein-DNA interactions with structural models. *Nucleic Acids Res.*, **41**, W523–W530.
34. Alamanova, D., Stegmaier, P. and Kel, A. (2010) Creating PWMs of transcription factors using 3D structure-based computation of protein-DNA free binding energies. *BMC Bioinformatics*, **11**, 225.
35. Pique-Regi, R., Degner, J.F., Pai, A.A., Gaffney, D.J., Gilad, Y. and Pritchard, J.K. (2011) Accurate inference of transcription factor binding from DNA sequence and chromatin accessibility data. *Genome Res.*, **21**, 447–455.
36. Battistini, F., Hospital, A., Buitrago, D., Gallego, D., Dans, P.D., Gelpi, J.L. and Orozco, M. (2019) How B-DNA dynamics decipher sequence-selective protein recognition. *J. Mol. Biol.*, **431**, 3845–3859.
37. Fornes, O., Garcia-Garcia, J., Bonet, J. and Oliva, B. (2014) On the use of knowledge-based potentials for the evaluation of models of protein-protein, protein-DNA, and protein-RNA interactions. *Adv. Protein Chem. Struct. Biol.*, **94**, 77–120.
38. Wasserman, W.W. and Sandelin, A. (2004) Applied bioinformatics for the identification of regulatory elements. *Nat. Rev. Genet.*, **5**, 276–287.
39. Stormo, G.D. (2000) DNA binding sites: representation and discovery. *Bioinformatics*, **16**, 16–23.
40. Khan, A., Fornes, O., Stigliani, A., Gheorghe, M., Castro-Mondragon, J.A., van der Lee, R., Bessy, A., Cheneby, J., Kulkarni, S.R., Tan, G. *et al.* (2018) JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res.*, **46**, D260–D266.
41. Kabsch, W. and Sander, C. (1983) Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers*, **22**, 2577–2637.
42. Lu, X.J. and Olson, W.K. (2008) 3DNA: a versatile, integrated software system for the analysis, rebuilding and visualization of three-dimensional nucleic-acid structures. *Nat. Protoc.*, **3**, 1213–1227.
43. Rice, P., Longden, I. and Bleasby, A. (2000) EMBOSS: the European Molecular Biology Open Software Suite. *Trends Genet.*, **16**, 276–277.
44. Altschul, S.F., Gertz, E.M., Agarwala, R., Schaffer, A.A. and Yu, Y.K. (2009) PSI-BLAST pseudocounts and the minimum description length principle. *Nucleic Acids Res.*, **37**, 815–824.
45. Webb, B. and Sali, A. (2016) Comparative protein structure modeling using MODELLER. *Curr. Protoc. Bioinformatics*, **54**, doi:10.1002/0471250953.bi050615.
46. Bailey, T.L., Johnson, J., Grant, C.E. and Noble, W.S. (2015) The MEME Suite. *Nucleic Acids Res.*, **43**, W39–W49.
47. Rose, P.W., Pric, A., Altunkaya, A., Bi, C., Bradley, A.R., Christie, C.H., Costanzo, L.D., Duarte, J.M., Dutta, S., Feng, Z. *et al.* (2017) The RCSB protein data bank: integrative view of protein, gene and 3D structural information. *Nucleic Acids Res.*, **45**, D271–D281.
48. UniProt, C. (2019) UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res.*, **47**, D506–D515.
49. Persikov, A.V., Rowland, E.F., Oakes, B.L., Singh, M. and Noyes, M.B. (2014) Deep sequencing of large library selections allows computational discovery of diverse sets of zinc fingers that bind common targets. *Nucleic Acids Res.*, **42**, 1497–1508.
50. Feliu, E., Aloy, P. and Oliva, B. (2011) On the analysis of protein-protein interactions via knowledge-based potentials for the prediction of protein-protein docking. *Protein Sci.*, **20**, 529–541.
51. Wiederstein, M. and Sippl, M.J. (2007) ProSA-web: interactive web service for the recognition of errors in three-dimensional structures of proteins. *Nucleic Acids Res.*, **35**, W407–W410.
52. Nakahashi, H., Kieffer Kwon, K.R., Resch, W., Vian, L., Dose, M., Stavreva, D., Hakim, O., Pruett, N., Nelson, S., Yamane, A. *et al.* (2013) A genome-wide map of CTCF multivalency redefines the CTCF code. *Cell Rep.*, **3**, 1678–1689.
53. Sigrist, C.J., Cerutti, L., de Castro, E., Langendijk-Genevaux, P.S., Bulliard, V., Bairoch, A. and Hulo, N. (2010) PROSITE, a protein domain database for functional characterization and annotation. *Nucleic Acids Res.*, **38**, D161–D166.
54. Weirauch, M.T., Yang, A., Albu, M., Cote, A.G., Montenegro-Montero, A., Drewe, P., Najafabadi, H.S., Lambert, S.A., Mann, I., Cook, K. *et al.* (2014) Determination and inference of eukaryotic transcription factor sequence specificity. *Cell*, **158**, 1431–1443.