

Two distinct DNA sequences recognized by transcription factors represent enthalpy and entropy optima

Ekaterina Morgunova¹, Yimeng Yin¹, Pratyush K. Das⁴, Arttu Jolma¹, Alexander Popov³,
You Xu², Lennart Nilsson² and Jussi Taipale^{1,4,5}

¹Department of Medical Biochemistry and Biophysics, Karolinska Institutet, SE 171 77
Stockholm, Sweden

²Department of Bioscience and Nutrition, Karolinska Institutet, SE 141 83 Huddinge,
Sweden

³European Synchrotron Radiation Facility, Grenoble, France

⁴Genome-Scale Biology Research Program, P.O. Box 63, FI-00014 University of
Helsinki, Finland

⁵Department of Biochemistry, University of Cambridge, United Kingdom

address correspondence to:

ajt208@cam.ac.uk

ekaterina.morgunova@ki.se

Lead Contact: Dr. Jussi Taipale, ajt208@cam.ac.uk

SUMMARY

Most transcription factors (TFs) can bind to a population of sequences closely related to a single optimal site. However, some TFs can bind to two distinct sequences that represent two local optima in the Gibbs free energy of binding (ΔG). To determine the molecular mechanism behind this effect, we solved the structures of human HOXB13 and CDX2 bound to their two optimal DNA sequences, CAATAAA and TCGTAAA. Thermodynamic analyses by isothermal titration calorimetry revealed that both sites were bound with similar ΔG . However, the interaction with the CAA sequence was driven by change in enthalpy (ΔH), whereas the TCG site was bound with similar affinity due to smaller loss of entropy (ΔS). The common presence of at least two local optima is general to all macromolecular interactions, as ΔG depends on two partially independent variables ΔH and ΔS according to the central equation of thermodynamics, $\Delta G = \Delta H - T\Delta S$.

Keywords: transcription, transcription factor, DNA recognition, HOXB13, CDX2, thermodynamics, enthalpy, entropy, epistasis.

INTRODUCTION

The binding of transcription factors (TFs) to their specific sites on genomic DNA is a key event regulating cellular processes. Analysis of structures of known TFs bound to DNA has revealed three different mechanisms of recognition of the specifically bound sequences: **1)** the “direct readout” mechanism involving the formation of specific hydrogen bonds and hydrophobic interactions between DNA bases and protein amino acids (Aggarwal, Rodgers, Drottar, Ptashne, & Harrison, 1988; J. E. Anderson, Ptashne, & Harrison, 1987; Wolberger, Dong, Ptashne, & Harrison, 1988); **2)** “indirect readout” of the DNA shape and electrostatic potential (Dror, Zhou, Mandel-Gutfreund, & Rohs, 2014; Hizver, Rozenberg, Frolow, Rabinovich, & Shakked, 2001; Joshi et al., 2007; Lavery, 2005; Rohs, Sklenar, & Shakked, 2005) by protein contacts to the DNA backbone or the minor groove, and **3)** water mediated interactions between bases and amino-acids (Bastidas & Showalter, 2013; Garner & Rau, 1995; Ladbury, Wright, Sturtevant, & Sigler, 1994; Morton & Ladbury, 1996; Patikoglou & Burley, 1997; Poon, 2012; Spolar & Record, 1994). Each of these mechanisms contributes to binding specificity of most TFs, with their relative importance varying depending on the TF and the recognized sequence.

The modes of DNA recognition differ from each other also in their thermodynamic characteristics. For example, direct hydrogen bonds can contribute strongly to enthalpy of binding, whereas indirect hydrogen bonds mediated by water are weaker, due to the loss of entropy caused by immobilization of the bridging water

molecule. In many cases, the contributions of the loss of entropy and the gain in enthalpy are similar in magnitude, leading to a phenomenon called "enthalpy-entropy compensation" (Chodera & Mobley, 2013; Jen-Jacobson, Engler, & Jacobson, 2000; Klebe, 2015; Patikoglou & Burley, 1997). This leads to binding promiscuity, allowing a TF to bind to several different but closely related sequences with a biologically relevant affinity.

Many transcription factors appear to only recognize sequences closely related to a single optimal site. Their binding to DNA can be approximated by a position weight matrix (PWM) model, which describes a single optimal site, and assumes that individual substitutions affect binding independently of each other. Thus, the combined effect of multiple mutations is predictable from the individual effects. However, some TFs have been shown to bind with high affinity to multiple different sequences, and populations of sequences that are closely related to these optimal sites (Badis et al., 2009; Johnson, Mortazavi, Myers, & Wold, 2007; Jolma et al., 2015; Morris, Bulyk, & Hughes, 2011; Zhao & Stormo, 2011). In such cases, the effect of substitution mutations is not independent, and instead the mutations display strongly epistatic behavior (D. W. Anderson, McKeown, & Thornton, 2015; Lehner, 2011), where the combined effect of two mutations can be less severe than what is predicted from the individual effects. Many cases of such multiple specificity can be explained by different spacing of homodimeric TFs (Aggarwal et al., 1988), but in some cases a single monomeric TF appears to be able to bind to two distinct sequences with similar affinities.

The molecular mechanism behind the phenomenon has not been understood. To elucidate the mechanism, we performed structural analysis of two homeodomain proteins, the posterior homeodomain protein HOXB13, and the parahox protein CDX2, each bound to two distinct high-affinity sequences. The optimal sequences for HOXB13 are CCAATAAA and CTCGTAAA that differ from each other by the three underlined base pairs, whereas CDX2 binds with high affinity to similar two sequences that begin with a G instead of a C.

This analysis, together with thermodynamic measurements of HOXB13, CDX2 and two other transcription factors, BARHL2 and MYF5 that also display multiple specificity revealed that in each case, one of the optimal sequences is bound primarily due to an optimal enthalpic contribution, whereas the other is bound due to an optimum of entropy. This result is likely to be general to most macromolecular interactions, as they commonly involve interaction of the macromolecules with a network of interconnected water-molecules, whose formation involves a trade-off between enthalpy and entropy.

RESULTS AND DISCUSSION

Modeling the binding of many TFs requires more than one PWM model

Many TFs have been reported to display multiple specificity, including many biologically important transcription factors such as the MYF family of basic helix-loop helix factors (Yin et al., 2017), the nuclear receptor HNF4A (Badis et al., 2009), and the homeodomain proteins BARHL2, CDX2 and HOXB13 (Jolma et al., 2015; Nitta et al., 2015). Analysis of enrichment of subsequences by MYF6, BARHL2, CDX2 and HOXB13 in SELEX reveals that a single PWM model cannot describe the binding affinity of these factors to DNA (**Figure 1A-D**). Each of these factors has more than one locally optimal sequence. All sequences between these optima have lower affinity and enrich less in SELEX than the optimal sequences. Therefore, more than one positionally independent position weight matrix (PWM) model is required for describing their affinity towards DNA (**Figure 1**).

Combinations of mutations affecting the optimal sites of these TFs display extremely strong epistatic effects. For example, the effect of mutating three first bases of the optimal HOXB13 motif TCGTAAAA is more than 400-fold smaller than what is expected from the individual single mutants (**Figure 1E, F**), and the generated CAATAAAA site binds to HOXB13 with almost the same affinity as the initial unmutated sequence.

Structural analysis of HOXB13 and CDX2 bound to DNA^{TCG} and DNA^{CAA}

To understand the molecular basis of the epistatic effect, we decided to solve the structure of HOXB13 and CDX2 bound to their two optimal sequences. These proteins are related, but diverged significantly in primary sequence, showing 43% identity at amino-acid level (**Figure S1 and S2**). For structural analysis, the DNA-binding domains (DBD) of HOXB13 (the 75 amino-acids Asp-209 to Pro-283) and CDX2 (the residues Arg-154 – Gln256) were expressed in *E.coli*, purified and crystallized bound to synthetic 19 or 18 bp double stranded DNA fragments containing the CTCGTAAA/GTCGTAAA (DNA^{TCG}) or CCAATAAA/GTCGTAAA (DNA^{CAA}) motifs, respectively. These core sequences were obtained by PBM (Berger et al., 2008) and HT-SELEX (Jolma et al., 2013), and validated by ChIP-seq experiments (Yin et al., 2017), and represent the two distinct binding sites of HOXB13 and CDX2 (**Figure 2**). The structures were solved using molecular replacement at resolutions 3.2 and 2.2 Å for HOXB13, and 2.57 and 2.95 Å for CDX2, respectively.

All complexes displayed a high overall similarity to HOXB13 bound to methylated DNA (Yin et al., 2017), and to the previously known DNA-bound HOX protein structures (Hovde, Abate-Shen, & Geiger, 2001; Joshi et al., 2007; LaRonde-LeBlanc & Wolberger, 2003; Passner, Ryoo, Shen, Mann, & Aggarwal, 1999; Piper, Batchelor, Chang, Cleary, & Wolberger, 1999; Zhang, Larsen, Stadler, & Ames, 2011) (**Figure 2; Figure S1**). Two parts of both HOXB13 and CDX2 DBDs interact with DNA: the recognition helix α_3 , which tightly packs into the major groove, and the N-terminal

tail interacting with the minor groove (**Figure 2A, C**). The residue Gly-84 that is affected by a coding variant that is strongly implicated in prostate cancer was not included in our construct; two other residues mutated in single prostate cancer families (Ewing et al., 2012) were predicted to destabilize the protein, or its interaction with DNA (**Figure S2**).

The core interactions between both HOXB13 and CDX2 DBDs and DNA are similar to those known from earlier structures (Hovde et al., 2001; Joshi et al., 2007; LaRonde-LeBlanc & Wolberger, 2003; Passner et al., 1999; Piper et al., 1999; Zhang et al., 2011). The TAAA sequence characteristic of the posterior homeodomains is recognized by a combination of a direct hydrogen bond to the A₁₀ base opposite of the T, A₁₀ base on the other stand and an insertion of the N-terminal basic amino-acids to the narrow minor groove induced by the stretch of four As. The overall protein structure in the four complexes is highly similar, showing only minor differences in the conformation of the N-termini (**Figure 2A, C**). The most remarkable difference between the complexes is in the conformation of DNA of the HOXB13-DNA^{TCG} complex at the position of the divergent bases (**Figure 2A, C; Figure S3**). To quantitate the shape of the DNA in the protein binding region we determined the helicoidal parameters using the program Curves+ (Lavery, Moakher, Maddocks, Petkeviciute, & Zakrzewska, 2009), and found the most prominent differences between the two complexes were in twist, shift, slide, X- and Y-displacement, minor groove width, and major groove depth at the positions of the divergent CAA and TCG sequences (**Figure S3B**). The DNA^{TCG} backbone is bent towards the major groove, facilitating contact with Arg-258 of the recognition helix with the DNA backbone. The corresponding contact (Arg-228 to backbone) is also observed in

both CDX2 structures. In contrast, the DNA^{CAA} backbone is bent towards the minor groove, leading to a contact with N-terminal Arg-217 and Lys-218 (**Figure 2B**; **Figure 3**). Instead of contacting the DNA backbone, Arg-258 assumes an alternative conformation in which it turns inside of the major groove, forming a water-mediated contact with Gln-265. The Gln-265, in turn, recognizes C_{6'} via a direct hydrogen bond. In addition, the CAA sequence is recognized by a hydrophobic interaction between Ile-262 and the T₁₁ methyl group. In CDX2 complexes the DNA bend in CDX2:DNA^{TCG} is slightly smaller due to the replacement of Thr-261 with Lys-231 which does not allow the alternative conformation of Arg-228. The other contacts in CDX2:DNA complexes are very similar to those listed for HOXB13:DNAs.

In order to understand the role of individual residues in binding of specific DNA we created 48 different single and combined mutations in DBD of HOXB13. The resulting data are presented in **Figures 2E** and **S4**. The replacement of Thr-261 either as a single mutation or in combination with any other amino-acids resulted in changing of HOXB13 specificity from Ctcg/Ccaa towards the sequence recognized by CDX2 (Gtcg/Gcaa; **Figure 2E**, left panel). No substitutions were identified that would lead to a specific loss of binding to the TCG or CAA sequences. However, several mutations affecting backbone contacts between HOXB13 amino-acids and DNA 5' of the divergent trinucleotide moderately increased the relative affinity towards the CAA sequence (**Figure 2E**, right panel).

Analysis of the mutation data together with the structures revealed that amino-acids involved in the protein-DNA interface formation cannot fully explain the specificity

preferences of HOXB13 and CDX2. The lack of direct interactions between protein and DNA in this region instead suggests that the specificity would be conferred in part by bridging water-molecules located at the protein-DNA interface.

Role of water molecules in the protein-DNA interface

The main difference between the complexes with DNA^{CAA} and DNA^{TCG} is revealed by analysis of the bridging water molecules. The HOXB13-DNA^{CAA} structure (2.2 Å) contains chains of water molecules that interact with both HOXB13 amino-acids and each of the DNA bases in the CAA sequence (**Figure 4A-C**). In contrast, no water molecules are visible in the HOXB13-DNA^{TCG} structure, despite the 3.2 Å resolution that should allow identification of strongly bound water molecules as well as much fewer water molecules are found in the complex CDX2:DNA^{TCG}. A relatively large solvent channel (6.4 Å in smallest diameter) exists between the α 3 helix of HOXB13 and DNA (**Figure 4D**). The electron density in this region is low ($\sigma < 0.5$), similar to that found in the surrounding solvent, indicating that the water-molecules in this region are highly mobile. Thus, the optimal binding of HOXB13 to the CAA sequence can be rationalized by the visible interactions that contribute to the enthalpy of binding (ΔH). In contrast, no such interactions can be identified that could explain the preference of HOXB13 to the TCG trinucleotide. The absence of ordered solvent molecules, and the lower resolution of the HOXB13-DNA^{TCG} structure is consistent with the possibility that the TCG sequence

is preferred because it represents a relatively disordered, high entropy state. In complex of CDX2:DNA^{TCG} with high resolution (2.57 Å) the water molecules were well visible but they did not form the corresponding water-chains (**Figure 4E, F**) supporting the idea of entropically driven binding.

Thermodynamic features of the protein-DNA interactions

We next performed molecular dynamics simulations and free energy perturbation calculations to probe the behavior of water molecules in the protein-DNA interface for the two optimal sequences for HOXB13. The relative free energy (Hansson, Marelius, & Aqvist, 1998) estimates for the affinities of HOXB13 for the two DNA sequences obtained from the simulations indicate that both sequences are bound with similar affinities ($\Delta\Delta G = -0.1$ kcal/mol). Analysis of the mobility of water molecules at the protein-DNA interface revealed that, while there is a similar number of water molecules in both systems, the waters at the HOXB13-DNA^{TCG} interface are more mobile (**Figure S5**), consistent with a model where this complex has higher entropy than the HOXB13-DNA^{CAA} complex.

To more directly test if the two states are driven by enthalpy and entropy, we measured these thermodynamic parameters using isothermal titration calorimetry (ITC). ITC directly measures the heat of binding (ΔH) and K_d of the binding reaction. Conversion of the K_d to ΔG then allows the inference of the entropy of binding (ΔS) from the data. The measured thermodynamic parameters for the TCG site were very similar to

those we reported previously (Yin et al., 2017). Comparison of the parameters for the TCG and CAA sites revealed that consistent with SELEX (Jolma et al., 2013) and molecular modeling data, the ΔG values for both sequences were similar. However, as predicted, the CAA site displayed much higher change in enthalpy, and larger loss of entropy compared to those of the TCG site (**Figure 5A, B**). These results indicate that HOXB13 binding to one optimal site, CAA, is driven by enthalpy, whereas strong binding to the other, TCG, is due to a lower loss of entropy.

To test if the identified mechanism is general to other cases of multiple specificity, we used ITC to determine the thermodynamic parameters for CDX2 and two other TFs, the MYF family TF MYF5 and the homeodomain protein BARHL2, both of which can optimally bind to two distinct sequence populations. Analysis of the data confirmed that in both cases, the ΔG values for the two optimally bound sequences were similar, whereas the relative contributions of entropy and enthalpy to the binding were strikingly different (**Figure 5C to H**). These results suggest that the ability of some TFs to bind to two distinct sequences with high affinity can be caused by the presence of both an enthalpic and an entropic optima.

Conclusions

In drug development, multiple optimal compounds can often be found that bind to a particular target molecule (Klebe, 2015). However, biological macromolecules are composed of a small set of relatively large monomers, and thus populate the shape-space

more sparsely than synthetic small molecules, which can be modified at the level of single atoms. Therefore, the finding that TFs can bind to two distinct DNA sequences with equal affinity was unexpected (Berger et al., 2008; Jolma et al., 2013), and has been controversial in the field. Our initial hypothesis was that the two optimal states could be due to an ability of the TF to adopt two distinct conformational states, or due to a similarity of the shapes of the two distinct DNA sequences. To address these hypotheses, we solved the structure of HOXB13, a central transcription factor involved in both development (Economides & Capecchi, 2003; Krumlauf, 1994; Nolte, 2015) and tumorigenesis (Ewing et al., 2012; Huang & Cai, 2014; Pomerantz et al., 2015) bound to its two optimal DNA sequences. Surprisingly, the conformational differences between the HOXB13 proteins in the two structures were minor, and the same effect we observed in two complexes of CDX2. In addition, the shape and charge-distribution of the optimally bound DNA sequences were not similar to each other. Thus, the structural analysis failed to support either the dual protein conformation or the DNA shape similarity models. Instead, thermodynamic analyses of HOXB13, CDX2, BARHL2 and MYF5 revealed that the two optimal states were bound because of their distinct effects on enthalpy and entropy, principally caused by differential stability of the water network at the protein-DNA interface.

The mechanism by which TFs bind to two optimal DNA sequences is fundamental, and applies to all macromolecular interactions. In principle, enthalpy and entropy of binding vary partially independently as a function of the shape and charge distribution of the interacting molecules. Thus, different sequences are likely to be

optimal with respect to enthalpy and entropy, in such a way that one optimal sequence is close to the optimal enthalpy state, and another is close to the optimal entropy state (**Figure 6**). We are not aware of work that has previously described such a situation. The observed effect that commonly results in at least two distinct ΔG minima may have been missed before because it is generally strongest when there are solvent molecules at the interacting surface. Interacting through solvent molecules can increase enthalpy of binding, but also causes a large loss of entropy due to fixing of the solvent molecule (s). However, in macromolecular interactions that are driven by direct contact between residues, entropy generally has lower impact on binding than enthalpy, and thus one of the optima is at a higher ΔG than the other. Another reason for overlooking this mechanism could have been the fact that multiple local optima can also exist via other mechanisms (see for example (Klebe, 2015)), and measurements that allow inference of entropy and enthalpy separately are not commonly performed in studies of macromolecular interactions. In addition, simple additive binding models such as position weight matrices (PWMs) can hide the effect, as they can only describe a single optimal state.

The cases we studied here represent some of the strongest deviations from the PWM model, and also present two optima of very similar ΔG that are located relatively far from each other in sequence space. It is likely that many other biologically relevant examples will be identified where the sequences representing entropic and enthalpic optima are more closely related to each other. This would manifest as a "flat bottom" in the affinity landscape, where many sequences would bind with similar affinity. In

addition, situations may be found where one of the local optima is located at lower affinity than the other. This would manifest as a minor peak or a shoulder in the affinity landscape farther from the optimal sequence. In each case, the measured affinities would deviate from those predicted from a single PWM model. Our results and the underlying theory suggest that the ability of TFs to bind to distinct sequences could be widespread, and that the importance of the optimal states in determining TF-DNA binding preferences should be reinvestigated. In addition, models for TF binding that are used to identify TF sites should also be adjusted to include features that allow two or more optima. In a broader sense, our results are relevant to all macromolecular interactions, particularly in the presence of a polar solvent such as water that can contribute to bridging interactions, whose contributions to the enthalpy and entropy of binding are in the same order of magnitude. Therefore, in addition to explaining the observed epistasis in protein-DNA interactions, the presence of two optima is likely to also explain the molecular mechanisms behind other types of genetic epistasis.

Author Contributions

E.M. and J.T. designed the project; E.M. performed crystallization, solved and refined crystal structures, performed ITC experiments, analyzed the results, A.P. collected data in ESRF, Y.Y., P.K.D. and A.J. performed HT-SELEX experiments and data analysis, L.N and Y.H. performed and analyzed molecular dynamics simulations, E.M. and J.T. prepared the figures and wrote the manuscript

Acknowledgments

The authors thank Drs. Minna Taipale, Inderpreet Sur, Bernhard Schmierer, Eevi Kaasinen, Sten Linnarsson and Johan Elf for the critical review of the manuscript, Karolinska Institutet Protein Science Facility and Sandra Augsten for protein production, as well as Lijuan Hu and Anna Zetterlund for technical assistance. This work was supported by the Center for Innovative Medicine at Karolinska Institutet, Cancerfonden, the Knut and Alice Wallenberg Foundation and the Swedish Research Council.

Accession codes

The atomic coordinates and diffraction data have been deposited to Protein Data Bank with the accession codes 5EDN and 5EEA, for HOXB13-DNA^{TCG} and HOXB13-DNA^{CAA}, respectively and 6ES3 and 6ES2 for CDX2-DNA^{TCG} and CDX2-DNA^{CAA}, respectively.

Competing Financial Interests

The authors declare no competing financial interests.

Figure legends

Figure 1.

Multiple TFs prefer to bind to two optimal sequences. (A) MYF6; (B) BARHL2; (C) HOXB13; (D) CDX2. Note that single PWM models (top) fail to describe sequence specificity towards different sequences shown in the bar graphs (middle). For example, a single PWM model for HOXB13 (panel C, top) predicts near-equal affinities towards sequences TCGT and TCAT at the position of the bracket, and lower affinity towards CAAT. Analysis of the counts of the subsequences (middle), instead, reveals that the TCAT sequence is bound more weakly than the two most preferred sequences TCGT and CAAT. Counts for maxima (dark blue) and related sequences that differ from the maxima by one or more base substitutions are also shown (light blue). The bars between the maxima represent sequences that can be obtained from both maximal sequences and have the highest count between the maxima. Bottom of each figure: Two distinct models that can represent the binding specificity of the TFs, the divergent bases are indicated by shading. For clarity, the PWM for the MYF6 optima that contains both AA and AC dinucleotide flanks (middle bar in a) is not shown. **(E)** Sequences representing the highest (blue line) and lowest (red line) affinity sequences between the two optimal HOXB13 sequences. y-axis: counts for 8-mer sequences containing the indicated trinucleotide followed by TAAA. **(F)** Epistasis in HOXB13-DNA binding. The effect of individual mutations (single mutants) to the optimal sequence TCGTAAAA (top) are relatively

severe, with binding decreasing by more than 70% in all cases (observed binding). However, combinations of the mutations (double mutants) do not decrease HOXB13 binding in a multiplicative manner (compare predicted and observed binding). A multiplicative model predicts that combining all three substitutions would abolish binding, but instead the CAA site is bound more strongly than any other mutant (triple mutant).

Figure 2.

Comparison of HOXB13-DNA complexes. (A) The view of superposition of HOXB13 (wheat) bound to DNA^{TCG} and HOXB13 (red) bound to DNA^{CAA} (rmsd = 0.813 Å on 57 residues). The respective DNAs are in blue and green. The dissimilar base pairs are presented as ball-and-stick models and colored as the proteins, DNA^{TCG} is wheat and DNA^{CAA} is red. Note the different bending of the DNA backbone at these positions (orange). **(B)** Schematic representation of interactions formed between HOXB13 DBD and the two different DNAs: left panel shows the interactions between HOXB13 and the primary binding site (DNA^{TCG}) and right panel represents the interactions of HOXB13 with the secondary site (DNA^{CAA}), respectively. Dashed lines represent interaction with backbone phosphates and deoxyribose and solid lines interactions with the bases. The protein residues belonging to the HOXB13-DNA^{TCG} and HOXB13-DNA^{CAA} structures are colored wheat and red, respectively. The divergent parts of the DNA sequences are highlighted by a light green box. Note that the TCG site lacks direct contacts to the DNA

bases, whereas the CAA site is recognized by direct contacts by Gln-265 and Ile-262. Most other contacts are well conserved in both structures. The four As of the TAAAA sequence are recognized by N-terminal amino-acids interacting with the DNA backbone via the minor groove, whereas the T is recognized by a bidentate interaction formed between its complementary adenine A₁₀ and side chain of asparagine Asn-266. Two hydrogen bonds are formed between nitrogen atoms N⁶ and N⁷ from adenine base and oxygen and nitrogen atoms of the Asn-266 side chain. This adenine-specific asparagine is totally conserved in the HOX family. **(C)** Superposition of CDX2 (cyan) bound to DNA^{TCG} and CDX2 (magenta) bound to DNA^{CAA} (rmsd = 0.270 Å on 64 residues). The respective DNAs are in blue and green. The dissimilar base pairs are presented as ball-and-stick models and colored as the proteins, DNA^{TCG} is green and DNA^{CAA} is blue. Note the different bending of the DNA backbone at these positions (orange). **(D)** Schematic representation of interactions formed between CDX2 DBD and the two different DNAs. **(F)** Structural interpretation of mutations that change the specificity of HOXB13: the mutations changing Ccaa/Ctcg to Gcaa/Gtcg are shown in a small box and, as a close view, on the left panel, and mutations, which switch the preferences of HOXB13 from CTCG to CCAA, are shown in big box and, as a close view, on the right panel. The mutations are presented in structural alignment of HOXB13 (red), HOXA9 (blue, PDB entry 1PUF) and CDX2 (pink) bound to DNA. Note the unique mutation of Lys (small box), which is conserved in all known HOXes, to Thr in HOXB13 allows HOXB13 to accept any base pair in the position before TCG/CAA. The left panel is representing the close view to the interactions formed by Lys in HOXA9 and CDX2.

Long aliphatic chain of Lys increases the hydrophobicity of this part of protein-DNA interface, pushing out the water molecules. Dashed line indicates water-mediated interaction between the ϵ -Amino group of Lys and the N7 and O6 of the guanine base at the Gtcg sequence. The right panel is representing the close view of triple mutation in the loop connecting helix 1 and helix 2: Lys-239/Met, Phe-240/Tyr and Ile-241/Leu; and single mutation of Lys-272/Arg. Those mutations change the hydrogen bond network inside the protein and between protein and DNA and lead to a preference towards the more rigid, more B-shaped DNA^{CAA}.

Figure 3.

Close view of the protein-DNA interactions. (A) HOXB13-DNA^{TCG} and (B) HOXB13-DNA^{CAA} complexes. The 2mFo-Fc maps contoured with 1.5σ are shown around the key residues. The residues and base pairs involved in protein to DNA contacts are also labeled. (C, D) Surface representation of the major groove in HOXB13-DNA^{TCG} and HOXB13-DNA^{CAA} complexes, respectively. The divergent bases are colored to indicate electrostatic charges of the atoms: neutral carbon atoms are green, oxygen atoms (negative) are red and nitrogen atoms (positive) are blue. Note the larger solvent-accessible space between amino-acids and bases in the TCG structure (C) and the difference in distribution of the positively and negatively charged spots on the surface that can contribute to differences in distribution of water molecules on the surface. (E) CDX2-DNA^{TCG} and (F) CDX2-DNA^{CAA} complexes. The 2mFo-Fc maps contoured with

1.5 σ are shown around the key residues. The residues and base pairs involved in protein to DNA contacts are also labeled.

Figure 4.

Close view of the role of water molecules in HOXB13-DNA interaction. (A) Schematic representation of water-mediated interactions between amino-acids (red typeface) of HOXB13 and DNA bases in the HOXB13-DNA^{CAA} structure. Different water chains are indicated with different shades of blue. Thick dashed lines represent interactions formed between water molecules and bases or amino acids; thin dashed lines represent contacts formed between water molecules, and solid blue line indicates the direct interaction between A₁₀ and Asn-266. Note that all of the base positions in the CAA sequence (boxes) are recognized via direct or water-mediated hydrogen bonds. (B) Structural representation of the network of interactions schematically presented in (A). Note the three water chains colored by slightly varied blue color. The amino acids and bases involved in interactions are presented as stick models. (C) Close view to the different conformations of amino-acids observed in HOXB13-DNA^{TCG} and HOXB13-DNA^{CAA} structures. Note that the conformations of the key amino-acids Gln-265 and Arg-258 that interact with the water network in HOXB13-DNA^{CAA} (amino-acids in red, DNA carbons in green) are not suitable for interacting with the network in HOXB13-DNA^{TCG} (amino-acids and DNA carbons in wheat). (D) Surface representation of

protein-DNA interface of HOXB13-DNA^{TCG} complex. Relatively large channel between the protein and DNA that goes through the protein-DNA interface (white) lend support to the presence of mobile water molecules in this region. TCG-bases are colored by atoms: carbon atoms are yellow; oxygen atoms are red and nitrogen atoms are blue. **(E)** Schematic representation of water-mediated interactions between amino-acids (cyan typeface) of CDX2 and DNA bases in the CDX2-DNA^{TCG} structure. Different water chains are indicated with different shades of red. Thick dashed lines represent interactions formed between water molecules and bases or amino acids; thin dashed lines represent contacts formed between water molecules, and solid red line indicates the direct interaction between A₁₀ and Asn-236. Note that only the position of the GC pair is recognized (boxes) via water-mediated hydrogen bonds. **(F)** Structural representation of the network of interactions schematically presented in **(E)**. Note the three water chains colored by varied red-pink color. The amino acids and bases involved in interactions are presented as stick models.

Figure 5.

Calorimetric titration data reveals that two optimal DNA sequences recognized by HOXB13 (A, B), CDX2 (C, D), BARHL2 (E, F) and MYF5 (G, H) represent enthalpy and entropy optima. The optimal sequences with higher enthalpic contribution to binding are presented on the left side (A, C, E, G) and the reactions with higher entropic contribution are presented on the right side (B, D, F, H). Note that for each protein both DNAs are bound with similar ΔG . The top panels of the ITC figures

represent raw data; the bottom panels show the integrated heat of the binding reaction. The red line represents the best fit to the data, according to the model that assumes a single set of identical sites. The determined changes of enthalpy and calculated losses of entropy are shown on the bottom panel. The changes of Gibbs free energy, $\Delta G = \Delta H - T\Delta S$, are also calculated and presented on the bottom panel of each isotherm.

Figure 6.

The two optimal sites bound by HOXB13 represent enthalpy and entropy-driven optima. (A-B) Schematic illustrations of the binding mechanism driven by the low enthalpy (A) and by high entropy (B) are presented in the left panels. The DNA bases are presented as pyrimidine and purine rings, protein is represented as ellipsoid, N-terminus is shown bound to the minor groove created by A-stretch, and water molecules are shown schematically and colored blue. The dashed lines represent hydrogen bonds observed in the low enthalpy state; solid line represents direct interactions between amino acids and bases. The blurred water molecules indicate the high entropy state. Hydrogen bonds that are common to both complexes are omitted for clarity. Graphs on the right show schematic illustration of the variance of enthalpy (ΔH , top), entropy ($-T\Delta S$, middle) and Gibbs free energy (ΔG) (bottom) as a function of an idealized one-dimensional continuous variable representing the high-dimensional variables of shape, electrostatic charge and vibration of DNA that vary as a function of the DNA sequence. As DNA is composed of only four bases, only discrete positions in this axis are possible (indicated by dots). Example models of shape and charge distribution of different DNA sequences

(from Figure 1C) are shown as surface representation above the scheme. The surfaces are colored according to the charge distribution: positively charged atoms are in blue, negatively charged are in red and neutral atoms are in green. Note that enthalpy and entropy are partially negatively correlated, leading to binding promiscuity (wider optima in ΔG compared to ΔH and ΔS). The remaining uncorrelated component leads to the presence of two optima for ΔG (bottom). Shaded boxes on the right show simplified dinucleotide binding models that illustrate how this leads to two distinct locally optimal sequences.

Table 1. Data collection and refinement statistics.

	HOXB13-DNA ^{TCG}	HOXB13-DNA ^{CAA}	CDX2-DNA ^{TCG}	CDX2-DNA ^{CAA}
Data collection				
Wavelength (Å)	0.9724	0.9724	0.9724	0.9724
Resolution range (Å)	45.95 - 2.19 (2.27 - 2.19)	46.29 - 3.2 (2.97 - 2.87)	43.23 - 2.57 (2.66 - 2.57)	55.96 - 2.95 (3.13 - 2.95)
Space group	P 1 2 1	P 2 2 2 ₁	C 1 2 1	I 1 2 1
Unit cell (Å, °)	77.35 57.92 101.28; 90 101.57 90	52.62 52.52 389.33; 90 90 90	127.95 46.49 68.89; 90 113.27 90	70.25 46.69 128.63; 90 101.40 90
Total reflections	241614 (21747)	86877 (3476)	19575 (1958)	27018 (4003)
Unique reflections	44125 (3912)	20590 (1049)	12095 (1197)	8802 (1264)
Multiplicity	5.5 (5.6)	4.2 (3.3)		3.2 (3.2)
Completeness (%)	97.42 (87.37)	97.5 (90.4)	99.5 (100)	96.6 (90.5)
Mean I/sigma(I)	8.11 (1.10)	7.91 (0.10)	8.47 (2.77)	7.5 (1.1)
R-merge	0.12 (1.21)	0.085 (4.59)	0.13 (5.49)	0.071 (7.24)
R-meas	0.13	0.09	0.08	0.09
CC1/2	0.99 (0.71)	0.99 (0.72)	0.99 (0.80)	0.99 (0.61)
Refinement				
R-work	0.25 (0.37)	0.21	0.23	0.19
R-free	0.29 (0.35)	0.28	0.29	0.25

Number of non-hydrogen atoms	5591	5197	2841	2783
macromolecule	5072	5172	2748	2717
water	519	17	93	66
Protein residues	274	242	144	141
RMS (bonds)	0.011	0.025	0.018	0.012
RMS (angles)	1.26	2.03	2.11	1.83
Ramachandran favored (%)	97	92	97.8	97.1
Ramachandran outliers (%)	0.41	1.7	1.43	0.73
Clashscore	5.31	10.51	4.42	6.43
Average B-factor	41.70	124.40	30.54	74.75
macromolecule	42.10	124.70	29.30	74.41

Statistics for the highest-resolution shell are shown in parentheses.

MATERIAL AND METHODS

Protein expression, purification and crystallization

Expression and purification of the DNA-binding domain fragment of human HOXB13 (residues 209-283) as well as CDX2 (residues 184-256) were performed as described in Refs. (Savitsky et al., 2010) and (Yin et al., 2017). The DNA fragments used in crystallization were obtained as single strand oligos (Eurofins), and annealed in 20 mM HEPES (pH 7.5) containing 150 mM NaCl and 0.5 mM *Tris* (2-carboxyethyl) phosphine (TCEP) and 5% glycerol. For each complex, the purified and concentrated protein was first mixed with a solution of annealed DNA duplex at a molar ratio 1:1.2 and after one hour on ice subjected to the crystallization trials. The crystallization conditions for all complexes were optimized using an in house developed crystal screening kit of different PEGs. Complexes were crystallized in sitting drops by vapor diffusion technique from solution containing 50 mM Tris buffer (pH 8.0), 100 mM MgCl₂, 150 mM KCl, 8% of PEG (400) and different concentrations of various PEGs. PEG (3350) was used at 14 % for HOXB13-DNA^{TCG} and 21 - 27 % (w/v) of polyethylene glycol monomethyl ether (PEGmme (5000)) was used in crystallizations of HOXB13 with DNA^{CAA} and CDX2 with both DNAs. The data sets were collected at ESRF from a single crystal on beam-lines ID29 (HOXB13:DNA^{TCG}) and ID23-1 (HOXB13:DNA^{CAA}, and both CDX2 complexes), at 100 K using the reservoir solution as cryo-protectant. The data collection strategy was optimized with the program BEST (Bourenkov & Popov, 2006). Data were integrated with the program XDS (Kabsch, 2010) and scaled with SCALA (Murshudov et al., 2011; Winn et al., 2011). Statistics of data collection are presented in **Table 1**.

Structure determination and refinement

All structures were solved by molecular replacement using program Phaser (McCoy et al., 2007) as implemented in Phenix (Adams et al., 2010) and CCP4 (Winn et al., 2011) with the structure of HOXA9 (pdb entry 1PUF) as a search model for HOXB13 and structure of CDX2:DNA^{meth} (pdb entry 5LTY) as a search model for CDX. After the positioning of protein, the density of DNA was clear and the molecule was built manually using COOT (Emsley, Lohkamp, Scott, & Cowtan, 2010). The rigid body refinement with REFMAC5 was followed by restrain refinement with REFMAC5, as implemented in CCP4 (Winn et al., 2011) and Phenix.refine (Afonine et al., 2012). The manual rebuilding of the model was done using COOT. The refinement statistics are presented in **Table 1**. The first 7 amino acid from N-termini and the last 7 from C-termini were found disordered and were not built in the maps. The end base pairs of the DNA in HOXB13-DNA^{TCG} structure were also found slightly disordered but it was possible to build them to the maps. Figures showing structural representations were prepared using PyMOL (Schrödinger, 2015).

HT-SELEX and motif analysis

HOXB13, CDX2, MYF6 and BARHL2 HT-SELEX experiments were performed essentially as described in Yin et al. (Ref. (Yin et al., 2017)). The PWM models were generated from cycles 4, 4, 3

and 3 of HOXB13 (from Ref. (Yin et al., 2017)), CDX2, MYF6 and BARHL2 HT-SELEX reads, respectively, using the multinomial (setting=1) method (Jolma et al., 2010) with the following seeds: HOXB13 single PWM: NCYMRTAAAAN, TCG: NCTCGTAAAAN, CAA: NCCAATAAAAAN; MYF6 single PWM: NRWCAGCTGWYN, AA...TT flank: NAACAGCTGTTN, GT...AC flank: NGTCAGCTGACN; BARHL2 single PWM: NSYTAAWYGNYN, TT: NSYTAATTGNYN, AC: NSYTAAACGKYN.

Molecular Dynamics

Molecular dynamics simulations were performed for HOXB13B complexed with either DNA^{TCG} or DNA^{CAA}; the DNA sequence used in the simulations contained nucleotides G5 – C18 from the crystal structure. The CHARMM 36 forcefield (Best et al., 2012; Foloppe & MacKerell, 2000; Hart et al., 2012; A. D. MacKerell et al., 1998; A. D. MacKerell, Jr., Feig, & Brooks, 2004) and CHARMM program (Brooks et al., 2009), with the CHARMM interface to OpenMM (Friedrichs et al., 2009) to allow the use of NVIDIA graphical processing units (GPUs), were used for all simulations. The starting structure was placed in a cubic solvent box with 8 nm side length containing water (Jorgensen, Chandrasekhar, Madura, Impey, & Klein, 1983) and 0.15M NaCl; Na⁺ ions were then added to neutralize the system. After energy minimization to relax initial strain the systems were heated from 100K to 300K over 0.1 ns followed by 0.3 ns simulation at constant pressure (1 bar) and constant temperature (300K), with soft harmonic positional restraints on the protein and DNA atoms. For each complex 3 x 0.8 μ s production runs were performed using the GPU, with the pressure and temperature maintained at 1 bar and 300K, respectively, and without the positional restraints. Particle mesh Ewald

summation was used to treat the long range electrostatic interactions, using a 6th order cubic spline interpolation for the charge distribution on the 0.1nm spaced grid points, $\kappa=0.34$. The same 0.9 nm cutoff was used for both the direct space part of the PME and for the van der Waals interactions, which were switched to zero from 0.8 nm to 0.9 nm, and the non-bond list was generated with a 1.1 nm cutoff. SHAKE (Ryckaert, 1977) was used to keep the lengths of all covalent X-H bonds fixed, allowing a time-step of 2 fs.

In the free energy perturbation calculations (Zwanzig, 1954) we changed the three base pairs in the TCG sequence into those of the CAA sequence using a total of 43 intermediate states, where the order of change was: turn off charges, change Lennard-Jones parameters, turn on charges. In each state, a 10 ns equilibration was followed by 10 ns production. The free energies were calculated using the Bennett Acceptance Ratio method (Bennett, 1976).

Isothermal titration calorimetry.

The ITC experiments were carried similarly to described in Ref. (Yin et al., 2017). Briefly, an ITC200 microcalorimeter (MicroCal Inc., Northampton, Massachusetts, USA) in PSF (Protein Science Facility at Karolinska Institute, Sweden) was used to measure binding isotherms of DNAs by direct titration of protein to the cell containing DNA. The measurements were taken at 25 C°. Both protein and DNA were prepared in a buffer containing 20 mM HEPES pH 7.5, 300 mM NaCl, 10% glycerol and 2 mM Tris (2-carboxyethyl)phosphine (TCEP). To measure binding affinity, a solution of 0.15mM protein was titrated to 0.012–0.016 mM solution of DNA. A total of 23 injections were made with 240 s between injections. Each experiment was

repeated three times for the reliability of the results. All data were evaluated using the OriginPro 7.0 software package (Microcal) supplied with the calorimeter. The apparent binding constant K_b , binding enthalpy ΔH and stoichiometry n , together with their corresponding standard deviation (s.d.), were determined by a nonlinear least square fit of the data to standard equations for the binding using a model for one set of independent and identical binding sites as implemented in the package. The entropy and free energy of binding were obtained from the relation $\Delta G = -RT\ln K_d = \Delta H - T\Delta S$.

REFERENCES

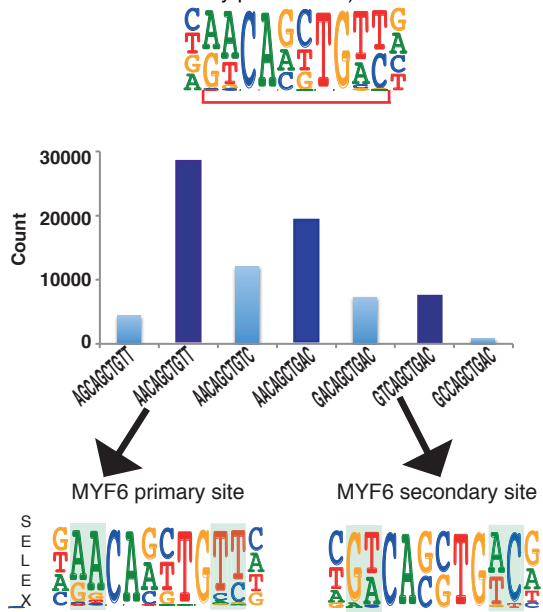
- Adams, P. D., Afonine, P. V., Bunkoczi, G., Chen, V. B., Davis, I. W., Echols, N., . . . Zwart, P. H. (2010). PHENIX: a comprehensive Python-based system for macromolecular structure solution. *Acta crystallographica. Section D, Biological crystallography*, 66(Pt 2), 213-221. doi:10.1107/S0907444909052925
- Afonine, P. V., Grosse-Kunstleve, R. W., Echols, N., Headd, J. J., Moriarty, N. W., Mustyakimov, M., . . . Adams, P. D. (2012). Towards automated crystallographic structure refinement with phenix.refine. *Acta crystallographica. Section D, Biological crystallography*, 68(Pt 4), 352-367. doi:10.1107/S0907444912001308
- Aggarwal, A. K., Rodgers, D. W., Drottar, M., Ptashne, M., & Harrison, S. C. (1988). Recognition of a DNA operator by the repressor of phage 434: a view at high resolution. *Science*, 242(4880), 899-907.
- Anderson, D. W., McKeown, A. N., & Thornton, J. W. (2015). Intermolecular epistasis shaped the function and evolution of an ancient transcription factor and its DNA binding sites. *eLife*, 4, e07864. doi:10.7554/eLife.07864
- Anderson, J. E., Ptashne, M., & Harrison, S. C. (1987). Structure of the repressor-operator complex of bacteriophage 434. *Nature*, 326(6116), 846-852. doi:10.1038/326846a0
- Badis, G., Berger, M. F., Philippakis, A. A., Talukder, S., Gehrke, A. R., Jaeger, S. A., . . . Bulyk, M. L. (2009). Diversity and complexity in DNA recognition by transcription factors. *Science*, 324(5935), 1720-1723. doi:10.1126/science.1162327
- Bastidas, M., & Showalter, S. A. (2013). Thermodynamic and structural determinants of differential Pdx1 binding to elements from the insulin and IAPP promoters. *Journal of molecular biology*, 425(18), 3360-3377. doi:10.1016/j.jmb.2013.06.011
- Bennett, C. H. (1976). Efficient Estimation of Free-Energy Differences from Monte-Carlo Data. *Journal of Computational Physics*, 22(2), 245-268. doi:10.1016/0021-9991(76)90078-4
- Berger, M. F., Badis, G., Gehrke, A. R., Talukder, S., Philippakis, A. A., Pena-Castillo, L., . . . Hughes, T. R. (2008). Variation in homeodomain DNA binding revealed by high-resolution analysis of sequence preferences. *Cell*, 133(7), 1266-1276. doi:10.1016/j.cell.2008.05.024
- Best, R. B., Zhu, X., Shim, J., Lopes, P. E., Mittal, J., Feig, M., & Mackerell, A. D., Jr. (2012). Optimization of the additive CHARMM all-atom protein force field targeting improved sampling of the backbone phi, psi and side-chain chi(1) and chi(2) dihedral angles. *Journal of chemical theory and computation*, 8(9), 3257-3273. doi:10.1021/ct300400x
- Bourenkov, G. P., & Popov, A. N. (2006). A quantitative approach to data-collection strategies. *Acta crystallographica. Section D, Biological crystallography*, 62(Pt 1), 58-64. doi:10.1107/S0907444905033998
- Brooks, B. R., Brooks, C. L., 3rd, Mackerell, A. D., Jr., Nilsson, L., Petrella, R. J., Roux, B., . . . Karplus, M. (2009). CHARMM: the biomolecular simulation program. *Journal of computational chemistry*, 30(10), 1545-1614. doi:10.1002/jcc.21287
- Chodera, J. D., & Mobley, D. L. (2013). Entropy-enthalpy compensation: role and ramifications in biomolecular ligand recognition and design. *Annual review of biophysics*, 42, 121-142. doi:10.1146/annurev-biophys-083012-130318

- Dror, I., Zhou, T., Mandel-Gutfreund, Y., & Rohs, R. (2014). Covariation between homeodomain transcription factors and the shape of their DNA binding sites. *Nucleic acids research*, 42(1), 430-441. doi:10.1093/nar/gkt862
- Economides, K. D., & Capecchi, M. R. (2003). Hoxb13 is required for normal differentiation and secretory function of the ventral prostate. *Development*, 130(10), 2061-2069. doi:10.1242/dev.00432
- Emsley, P., Lohkamp, B., Scott, W. G., & Cowtan, K. (2010). Features and development of Coot. *Acta crystallographica. Section D, Biological crystallography*, 66(Pt 4), 486-501. doi:10.1107/S0907444910007493
- Ewing, C. M., Ray, A. M., Lange, E. M., Zuhlke, K. A., Robbins, C. M., Tembe, W. D., . . . Cooney, K. A. (2012). Germline Mutations in HOXB13 and Prostate-Cancer Risk. *New England Journal of Medicine*, 366(2), 141-149.
- Foloppe, N., & Mackerell, A. D. (2000). All-atom empirical force field for nucleic acids: I. Parameter optimization based on small molecule and condensed phase macromolecular target data. *Journal of computational chemistry*, 21(2), 86-104. doi:10.1002/(Sici)1096-987x(20000130)21:2<86::Aid-Jcc2>3.0.Co;2-G
- Friedrichs, M. S., Eastman, P., Vaidyanathan, V., Houston, M., Legrand, S., Beberg, A. L., . . . Pande, V. S. (2009). Accelerating molecular dynamic simulation on graphics processing units. *Journal of computational chemistry*, 30(6), 864-872. doi:10.1002/jcc.21209
- Garner, M. M., & Rau, D. C. (1995). Water release associated with specific binding of gal repressor. *The EMBO journal*, 14(6), 1257-1263.
- Hansson, T., Marelus, J., & Aqvist, J. (1998). Ligand binding affinity prediction by linear interaction energy methods. *Journal of computer-aided molecular design*, 12(1), 27-35.
- Hart, K., Foloppe, N., Baker, C. M., Denning, E. J., Nilsson, L., & Mackerell, A. D., Jr. (2012). Optimization of the CHARMM additive force field for DNA: Improved treatment of the BI/BII conformational equilibrium. *Journal of chemical theory and computation*, 8(1), 348-362. doi:10.1021/ct200723y
- Hizver, J., Rozenberg, H., Frolow, F., Rabinovich, D., & Shakked, Z. (2001). DNA bending by an adenine-thymine tract and its role in gene regulation. *Proceedings of the National Academy of Sciences of the United States of America*, 98(15), 8490-8495. doi:10.1073/pnas.151247298
- Hovde, S., Abate-Shen, C., & Geiger, J. H. (2001). Crystal structure of the Msx-1 homeodomain/DNA complex. *Biochemistry*, 40(40), 12013-12021.
- Huang, H., & Cai, B. (2014). G84E mutation in HOXB13 is firmly associated with prostate cancer risk: a meta-analysis. *Tumor Biology*, 35(2), 1177-1182. doi:10.1007/s13277-013-1157-5
- Jen-Jacobson, L., Engler, L. E., & Jacobson, L. A. (2000). Structural and thermodynamic strategies for site-specific DNA binding proteins. *Structure*, 8(10), 1015-1023.
- Johnson, D. S., Mortazavi, A., Myers, R. M., & Wold, B. (2007). Genome-wide mapping of in vivo protein-DNA interactions. *Science*, 316(5830), 1497-1502. doi:10.1126/science.1141319
- Jolma, A., Kivioja, T., Toivonen, J., Cheng, L., Wei, G., Enge, M., . . . Taipale, J. (2010). Multiplexed massively parallel SELEX for characterization of human transcription factor binding specificities. *Genome research*, 20(6), 861-873. doi:10.1101/gr.100552.109
- Jolma, A., Yan, J., Whittington, T., Toivonen, J., Nitta, K. R., Rastas, P., . . . Taipale, J. (2013). DNA-Binding Specificities of Human Transcription Factors. *Cell*, 152(1-2), 327-339. doi:10.1016/j.cell.2012.12.009
- Jolma, A., Yin, Y., Nitta, K. R., Dave, K., Popov, A., Taipale, M., . . . Taipale, J. (2015). DNA-dependent formation of transcription factor pairs alters their binding specificity. *Nature*, 527(7578), 384-388. doi:10.1038/nature15518

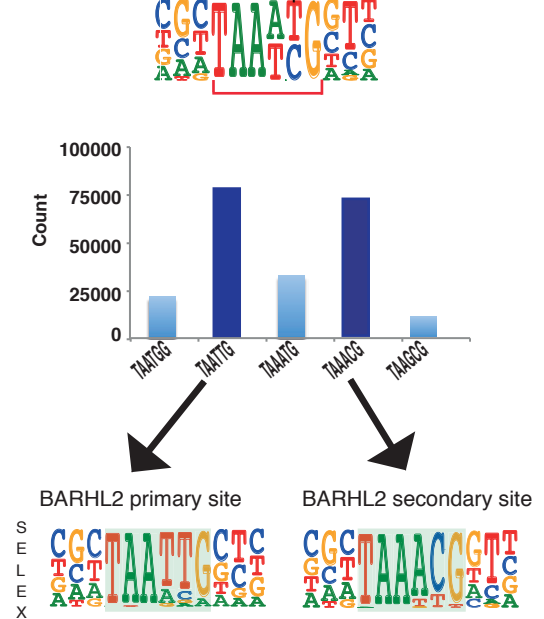
- Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. W., & Klein, M. L. (1983). Comparison of Simple Potential Functions for Simulating Liquid Water. *Journal of Chemical Physics*, 79(2), 926-935. doi:10.1063/1.445869
- Joshi, R., Passner, J. M., Rohs, R., Jain, R., Sosinsky, A., Crickmore, M. A., . . . Mann, R. S. (2007). Functional specificity of a Hox protein mediated by the recognition of minor groove structure. *Cell*, 131(3), 530-543. doi:10.1016/j.cell.2007.09.024
- Kabsch, W. (2010). Xds. *Acta crystallographica. Section D, Biological crystallography*, 66(Pt 2), 125-132. doi:10.1107/S0907444909047337
- Klebe, G. (2015). Applying thermodynamic profiling in lead finding and optimization. *Nature reviews. Drug discovery*, 14(2), 95-110. doi:10.1038/nrd4486
- Krumlauf, R. (1994). Hox Genes in Vertebrate Development. *Cell*, 78(2), 191-201. doi:10.1016/0092-8674(94)90290-9
- Ladbury, J. E., Wright, J. G., Sturtevant, J. M., & Sigler, P. B. (1994). A thermodynamic study of the trp repressor-operator interaction. *Journal of molecular biology*, 238(5), 669-681. doi:10.1006/jmbi.1994.1328
- LaRonde-LeBlanc, N. A., & Wolberger, C. (2003). Structure of HoxA9 and Pbx1 bound to DNA: Hox hexapeptide and DNA recognition anterior to posterior. *Genes & development*, 17(16), 2060-2072. doi:10.1101/gad.1103303
- Lavery, R. (2005). Recognizing DNA. *Quarterly reviews of biophysics*, 38(4), 339-344. doi:10.1017/S0033583505004105
- Lavery, R., Moakher, M., Maddocks, J. H., Petkeviciute, D., & Zakrzewska, K. (2009). Conformational analysis of nucleic acids revisited: Curves+. *Nucleic acids research*, 37(17), 5917-5929. doi:10.1093/nar/gkp608
- Lehner, B. (2011). Molecular mechanisms of epistasis within and between genes. *Trends in genetics : TIG*, 27(8), 323-331. doi:10.1016/j.tig.2011.05.007
- MacKerell, A. D., Bashford, D., Bellott, M., Dunbrack, R. L., Evanseck, J. D., Field, M. J., . . . Karplus, M. (1998). All-atom empirical potential for molecular modeling and dynamics studies of proteins. *The journal of physical chemistry. B*, 102(18), 3586-3616. doi:10.1021/jp973084f
- MacKerell, A. D., Jr., Feig, M., & Brooks, C. L., 3rd. (2004). Improved treatment of the protein backbone in empirical force fields. *Journal of the American Chemical Society*, 126(3), 698-699. doi:10.1021/ja036959e
- McCoy, A. J., Grosse-Kunstleve, R. W., Adams, P. D., Winn, M. D., Storoni, L. C., & Read, R. J. (2007). Phaser crystallographic software. *Journal of Applied Crystallography*, 40(Pt 4), 658-674. doi:10.1107/S0021889807021206
- Morris, Q., Bulyk, M. L., & Hughes, T. R. (2011). Jury remains out on simple models of transcription factor specificity. *Nature biotechnology*, 29(6), 483-484. doi:10.1038/nbt.1892
- Morton, C. J., & Ladbury, J. E. (1996). Water-mediated protein-DNA interactions: the relationship of thermodynamics to structural detail. *Protein science : a publication of the Protein Society*, 5(10), 2115-2118. doi:10.1002/pro.5560051018
- Murshudov, G. N., Skubak, P., Lebedev, A. A., Pannu, N. S., Steiner, R. A., Nicholls, R. A., . . . Vagin, A. A. (2011). REFMAC5 for the refinement of macromolecular crystal structures. *Acta crystallographica. Section D, Biological crystallography*, 67(Pt 4), 355-367. doi:10.1107/S0907444911001314
- Nitta, K. R., Jolma, A., Yin, Y., Morgunova, E., Kivioja, T., Akhtar, J., . . . Taipale, J. (2015). Conservation of transcription factor binding specificities across 600 million years of bilateria evolution. *eLife*, 4. doi:10.7554/eLife.04837
- Nolte, C. A., Tara B; and Krumlauf, Robb. (2015). Mammalian Embryo: Hox Genes. *eLS*, April. Retrieved from doi:10.1002/9780470015902.a0000740.pu3

- Passner, J. M., Ryoo, H. D., Shen, L., Mann, R. S., & Aggarwal, A. K. (1999). Structure of a DNA-bound Ultrabithorax-Extradenticle homeodomain complex. *Nature*, 397(6721), 714-719. doi:10.1038/17833
- Patikoglou, G., & Burley, S. K. (1997). Eukaryotic transcription factor-DNA complexes. *Annual review of biophysics and biomolecular structure*, 26, 289-325. doi:10.1146/annurev.biophys.26.1.289
- Piper, D. E., Batchelor, A. H., Chang, C. P., Cleary, M. L., & Wolberger, C. (1999). Structure of a HoxB1-Pbx1 heterodimer bound to DNA: role of the hexapeptide and a fourth homeodomain helix in complex formation. *Cell*, 96(4), 587-597.
- Pomerantz, M. M., Li, F., Takeda, D. Y., Lenci, R., Chonkar, A., Chabot, M., . . . Freedman, M. L. (2015). The androgen receptor cistrome is extensively reprogrammed in human prostate tumorigenesis. *Nature genetics*, 47(11), 1346-1351. doi:10.1038/ng.3419
- Poon, G. M. (2012). Sequence discrimination by DNA-binding domain of ETS family transcription factor PU.1 is linked to specific hydration of protein-DNA interface. *The Journal of biological chemistry*, 287(22), 18297-18307. doi:10.1074/jbc.M112.342345
- Rohs, R., Sklenar, H., & Shakked, Z. (2005). Structural and energetic origins of sequence-specific DNA bending: Monte Carlo simulations of papillomavirus E2-DNA binding sites. *Structure*, 13(10), 1499-1509. doi:10.1016/j.str.2005.07.005
- Ryckaert, J. P., G. Ciccotti, G., Berendsen H.J.C. (1977). Numerical Integration of the Cartesian Equations of Motion of a System with Constraints: Molecular Dynamics of n-Alkanes. *J. Comput. Phys.*, 23, 15.
- Savitsky, P., Bray, J., Cooper, C. D., Marsden, B. D., Mahajan, P., Burgess-Brown, N. A., & Gileadi, O. (2010). High-throughput production of human proteins for crystallization: the SGC experience. *Journal of structural biology*, 172(1), 3-13. doi:10.1016/j.jsb.2010.06.008
- Schrödinger, L. L. C. (2015). The PyMOL Molecular Graphics System. Version 1.8.
- Spolar, R. S., & Record, M. T., Jr. (1994). Coupling of local folding to site-specific binding of proteins to DNA. *Science*, 263(5148), 777-784.
- Winn, M. D., Ballard, C. C., Cowtan, K. D., Dodson, E. J., Emsley, P., Evans, P. R., . . . Wilson, K. S. (2011). Overview of the CCP4 suite and current developments. *Acta Crystallographica Section D-Biological Crystallography*, 67, 235-242. doi:10.1107/S09074444910045749
- Wolberger, C., Dong, Y. C., Ptashne, M., & Harrison, S. C. (1988). Structure of a phage 434 Cro/DNA complex. *Nature*, 335(6193), 789-795. doi:10.1038/335789a0
- Yin, Y., Morgunova, E., Jolma, A., Kaasinen, E., Sahu, B., Khund-Sayeed, S., . . . Taipale, J. (2017). Impact of cytosine methylation on DNA binding specificities of human transcription factors. *Science*, 356(6337). doi:10.1126/science.aaj2239
- Zhang, Y., Larsen, C. A., Stadler, H. S., & Ames, J. B. (2011). Structural basis for sequence specific DNA binding and protein dimerization of HOXA13. *Plos One*, 6(8), e23069. doi:10.1371/journal.pone.0023069
- Zhao, Y., & Stormo, G. D. (2011). Quantitative analysis demonstrates most transcription factors require only simple models of specificity. *Nature biotechnology*, 29(6), 480-483. doi:10.1038/nbt.1893
- Zwanzig, R. W. (1954). High-Temperature Equation of State by a Perturbation Method .1. Nonpolar Gases. *Journal of Chemical Physics*, 22(8), 1420-1426.

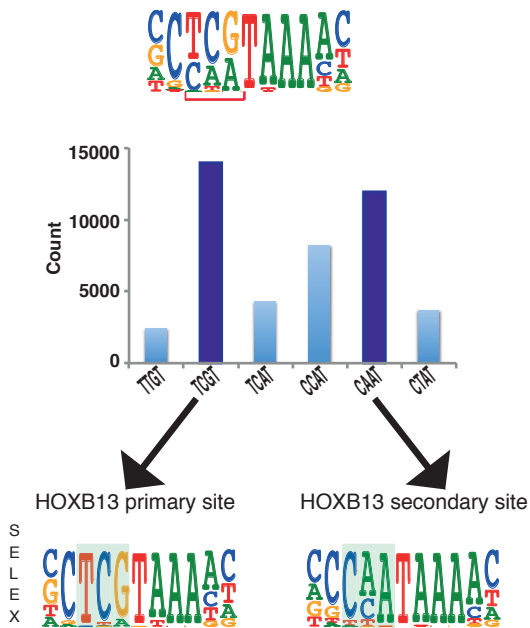
A Single PWM model for MYF6
 bioRxiv preprint doi: <https://doi.org/10.1101/205906>; this version posted November 7, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



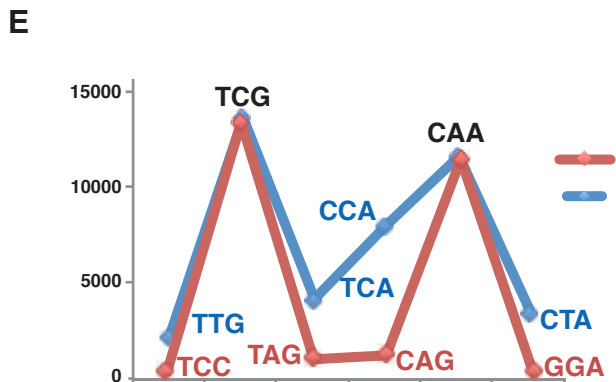
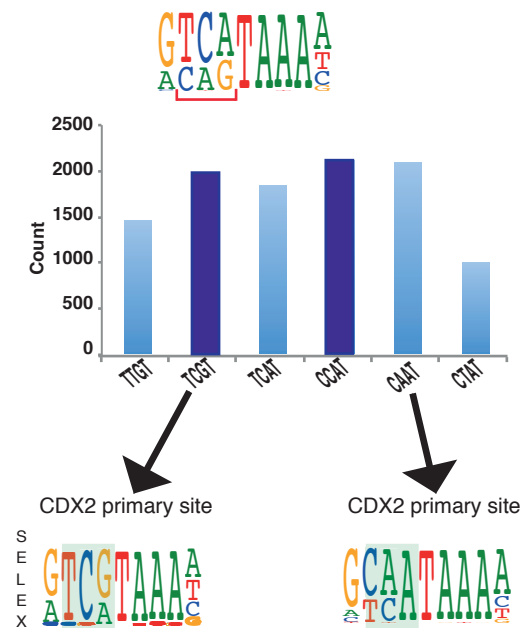
B Single PWM model for BARHL2



C Single PWM model for HOXB13

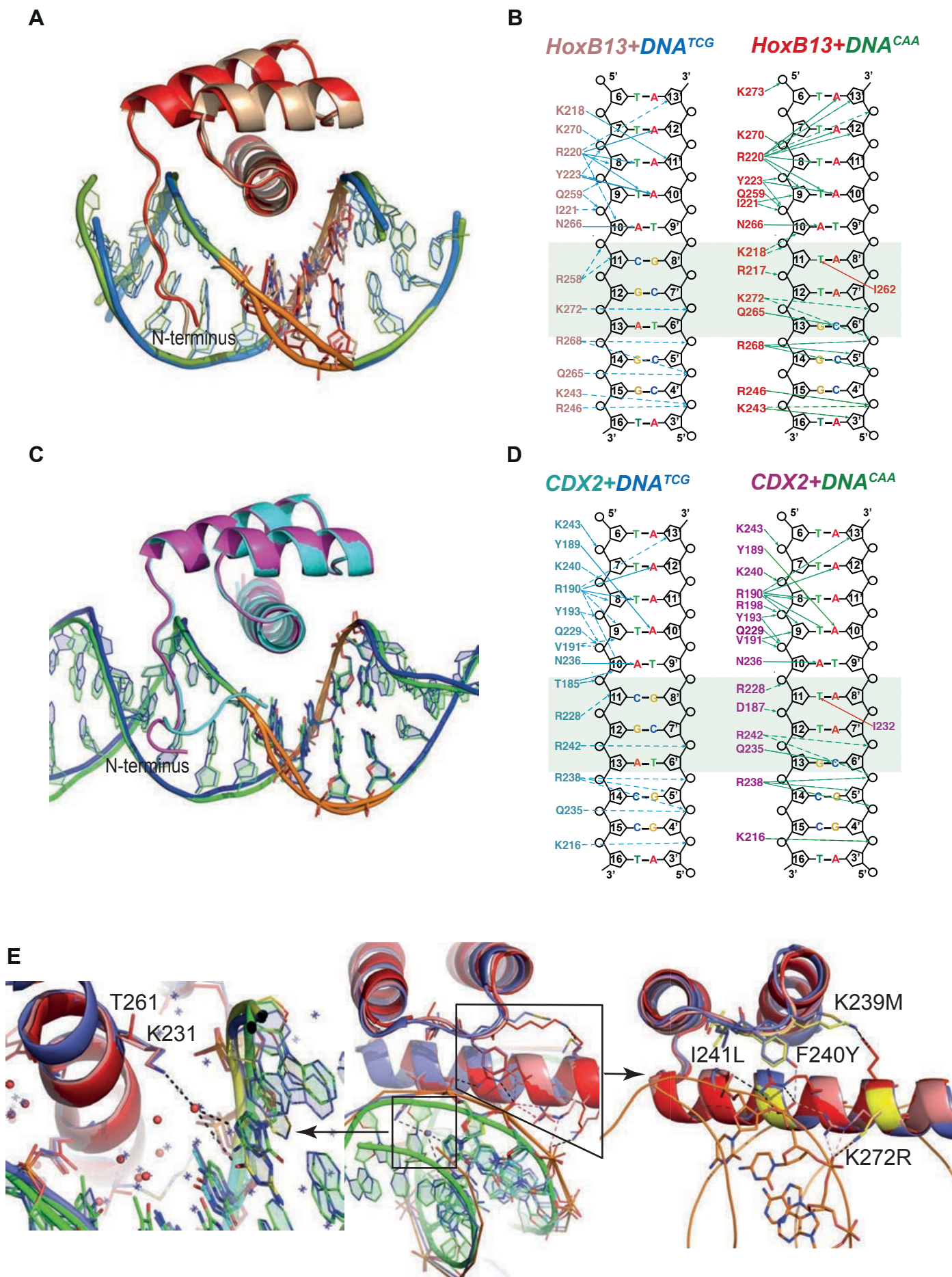


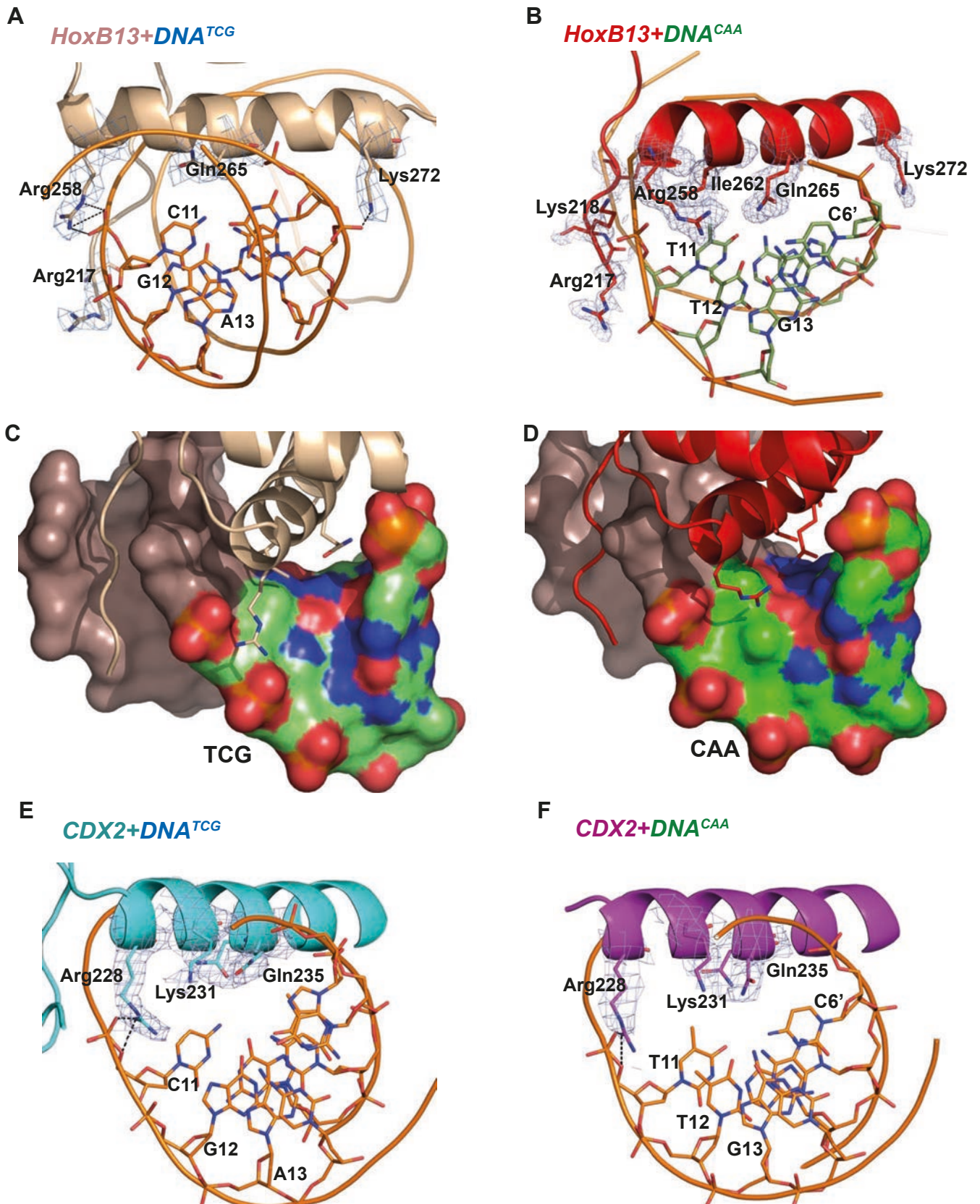
D Single PWM model for CDX2

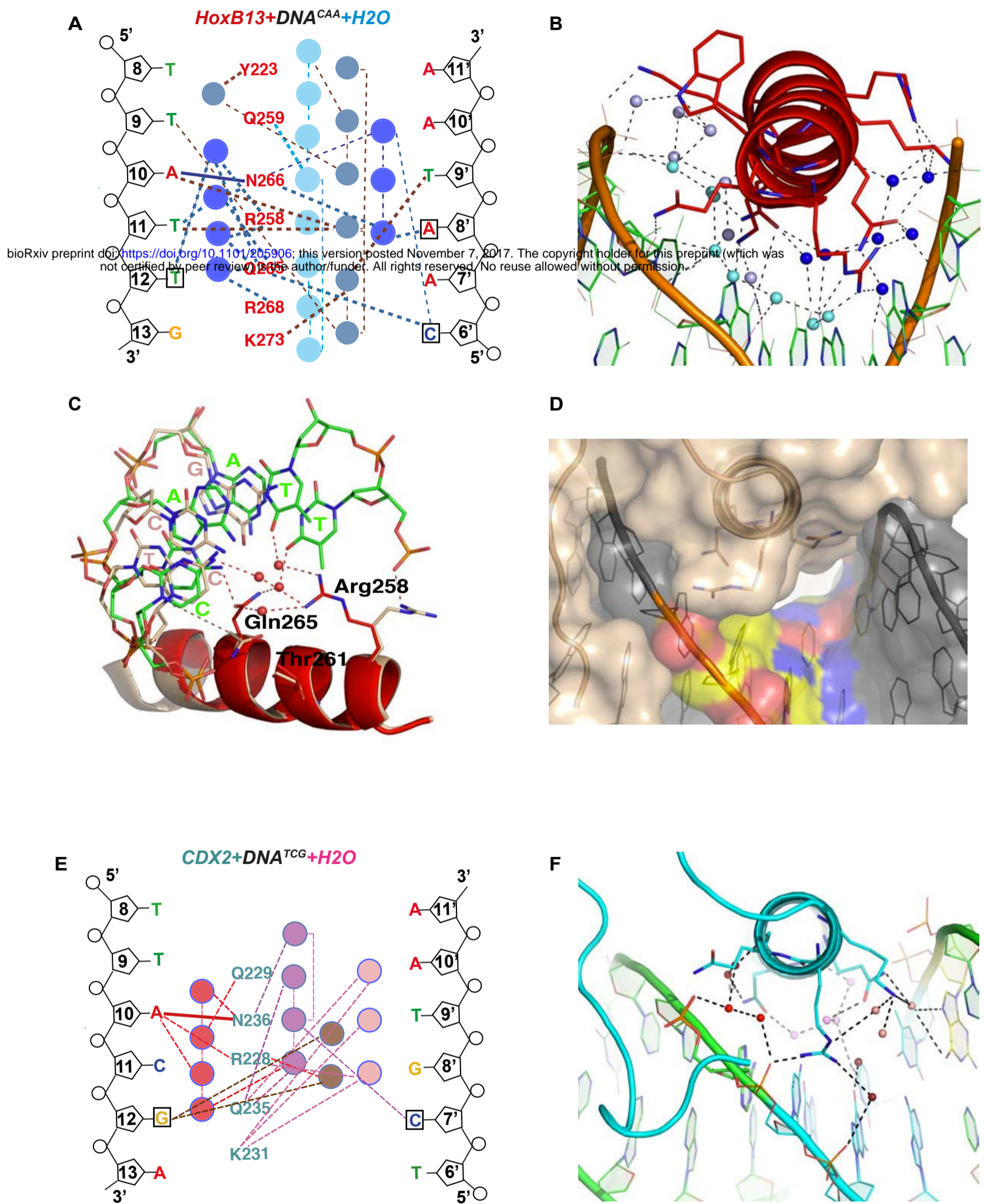


F

EPISTASIS	TCG(TAAAA)		
	100%		
Single mutants	CCG	TAG	TCA
observed binding	12.6%	5.3%	29.6%
Double mutants	CAG	CCA	TAA
observed binding	7.1%	58.2%	18.4%
predicted binding	0.7%	3.7%	1.6%
fold difference	10.5	15.7	11.6
Triple mutant	CAA		
observed binding	85.3%		
predicted binding	0.2%		
fold difference	429.2		

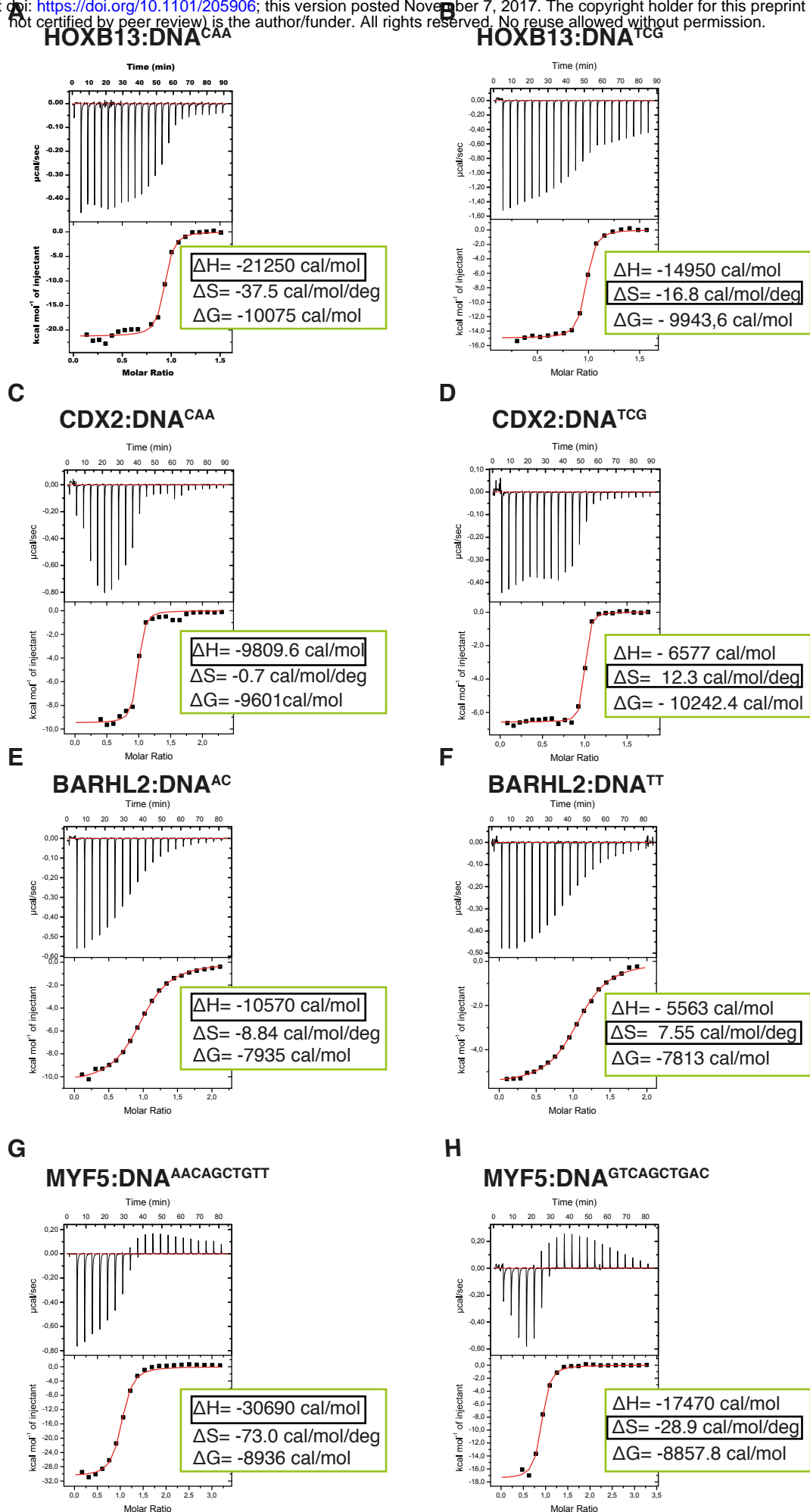






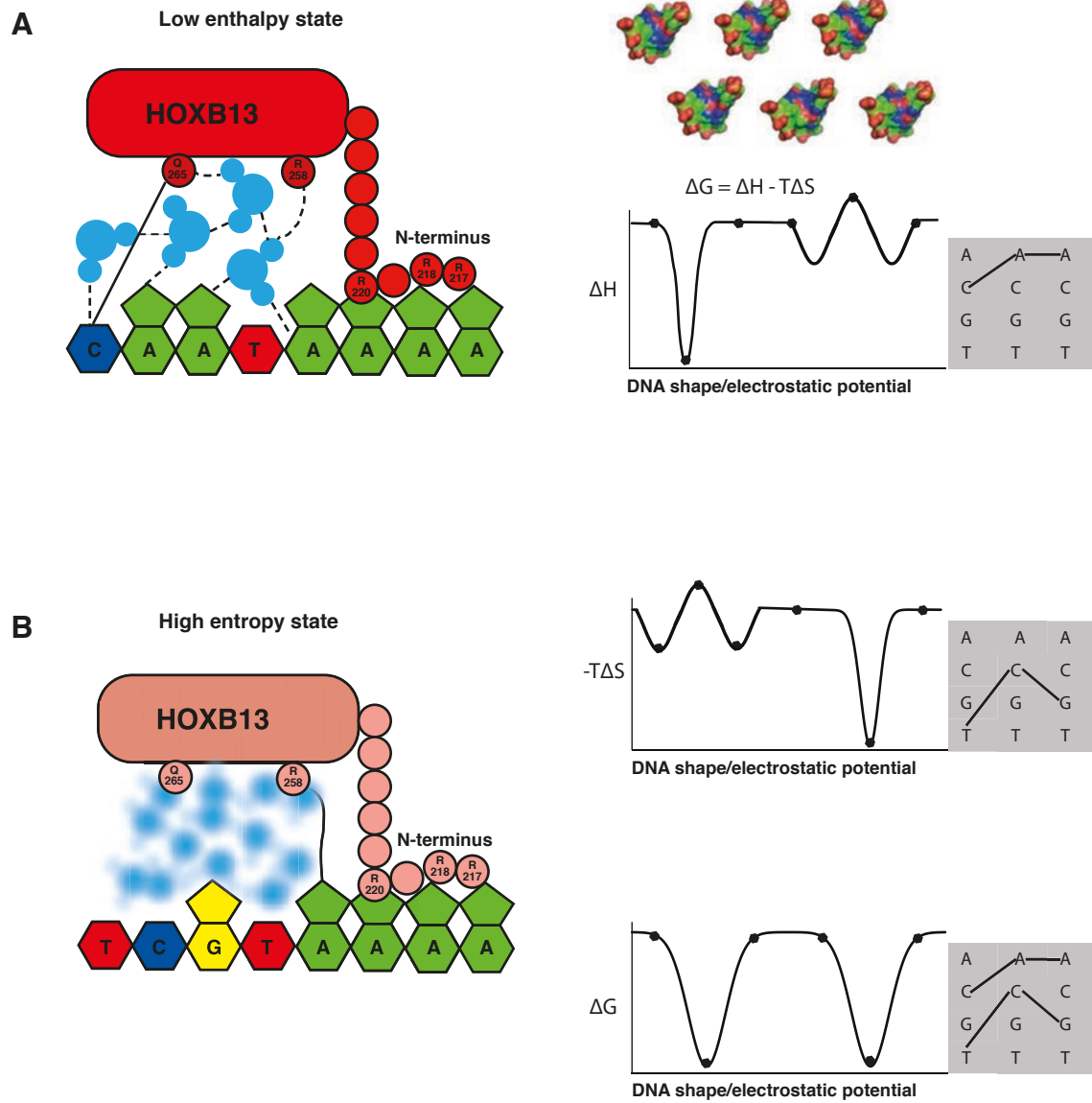
Morgunova et al., 2017 Figure 5

bioRxiv preprint doi: <https://doi.org/10.1101/205906>; this version posted November 7, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



Morgunova et al., 2017 Figure 6

bioRxiv preprint doi: <https://doi.org/10.1101/205906>; this version posted November 7, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



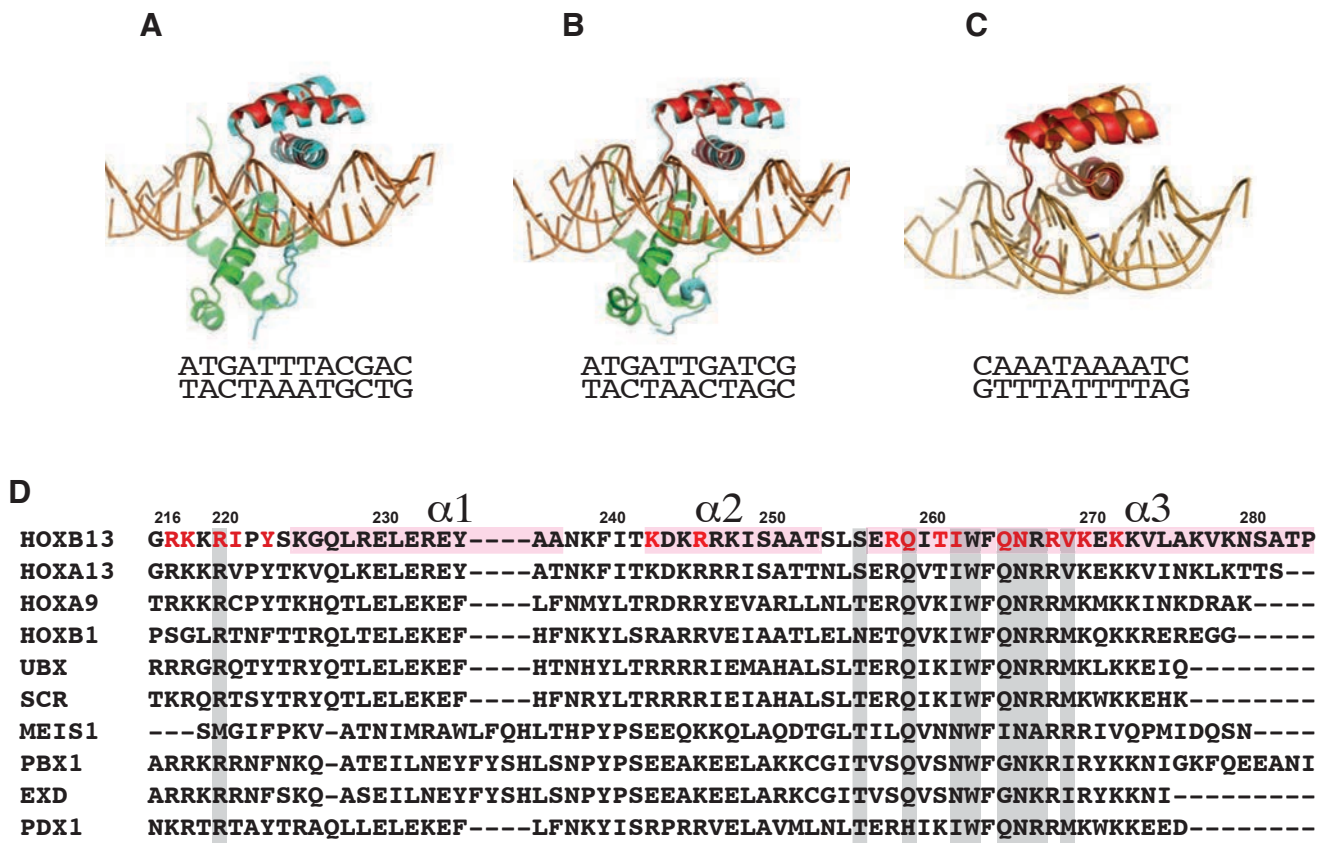
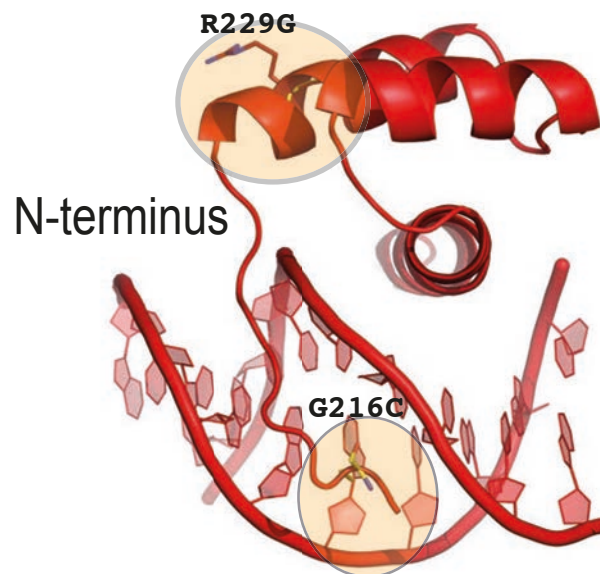


Figure S1. The comparison of HOXB13 structure with HOXB1 and HOXA9, Related to Figure 2

The superposition of HOXB13-DNA^{CAA} complex (red) with: **(A)** HOXA9:PBX1-DNA complex (HOXA9 is in cyan, PBX1 is green, PDB entry 1PUF); **(B)** with HOXB1:PBX1-DNA complex (HOXB1 is in cyan, PBX1 is in green, PDB entry 1B72); **(C)** with HOXA13-DNA complex (HOXA13 is in orange, PDB entry 2LD5). The corresponding DNA sequences are presented under pictures.

(D) The sequence alignment of Hox proteins with known structures. The numbering corresponds to HOXB13. Three helices are labeled on the top and highlighted with light pink. The residues involved in interactions are highlighted in grey. The residues involved in interaction in HOXB13 are colored red.

A



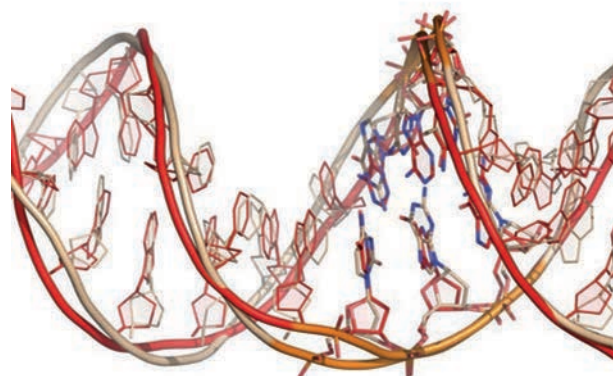
B

	216	220	230	240	250	260	270	280
HOXB13	G	R	K	R	I	P	Y	S
HOXA13	G	R	K	R	V	P	Y	T
HOXC13	G	R	K	R	V	P	Y	T
HOXD13	G	R	K	R	V	P	Y	T
HOXC12	S	R	K	R	K	P	Y	S
HOXD12	A	R	K	R	K	P	Y	T
HOXA11	T	R	K	R	C	P	Y	T
HOXC11	T	R	K	R	C	P	Y	T
HOSD11	S	R	K	R	C	P	Y	T
HOXA10	G	R	K	R	C	P	Y	T
HOXC10	G	R	K	R	C	P	Y	T
HOXD10	G	R	K	R	C	P	Y	T
HOXA9	T	R	K	R	C	P	Y	T
HOXB9	S	R	K	R	C	P	Y	T
HOXC9	T	R	K	R	C	P	Y	T
HOXD9	T	R	K	R	C	P	Y	T
CDX1	K	D	K	Y	R	V	V	T
CDX2	K	D	K	Y	R	V	V	T

Figure S2. HOXB13 prostate cancer mutation

(A) Structural representation of two of three residues found mutated in single prostate families, Gly-216-Cys and Arg-229-Gly. The mutated residues are presented in ball-and-stick style and highlighted with orange rings. Note that the first mutation Gly-216-Cys belonging to the N-termini of HOXB13 DBD can affect the interactions forming by protein in narrow minor groove. The other mutation Arg-229-Gly is located at the beginning of helix 1 and because glycine residues is known as “helix-breaker” the mutation can effect the interaction between N-termini with DNA as well as the interaction between two helixes. (B) Sequence alignment of posterior members of HOX family. The cancer mutations found in HOXB13 are colored red. Light red columns highlight the residues involved in interactions with DNA. The numbering corresponds to HOXB13.

A



DNA^{TCG} + DNA^{CAA}
rmsd = 1.128

B

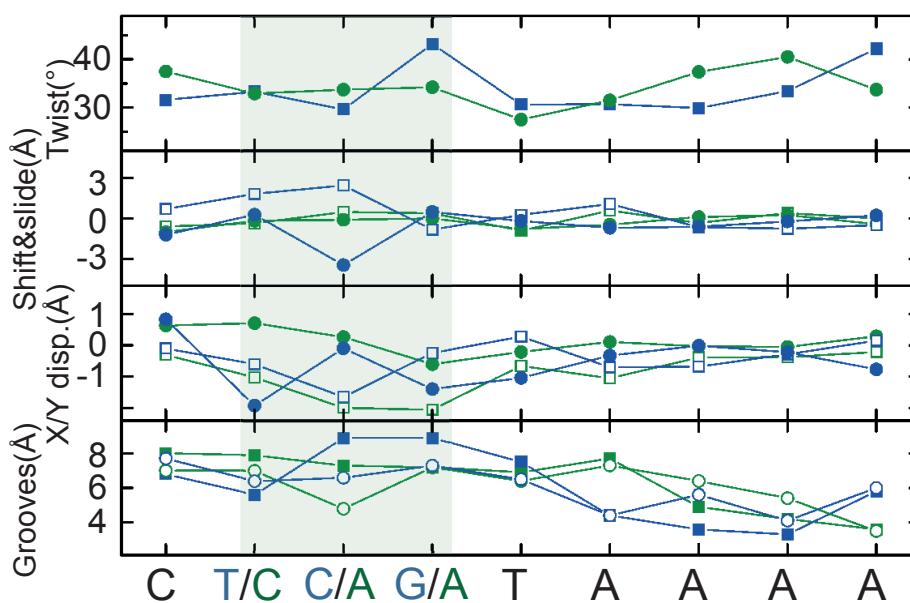


Figure S3. Pairwise comparison of two DNA molecules, Related to Figure 2

(A) Pairwise comparisons of DNA^{TCG} (wheat) and DNA^{CAA} (red); (B) Helicoidal parameters for HOXB13-DNA^{TCG} (blue) and HOXB13-DNA^{CAA} (green). Top: Helical twist; Middle top: shift (squares) and slide (circles); Middle bottom: X- (squares) and Y-displacements (circles); Bottom: Minor groove width (squares) and major groove depth (circles). The most pronounced differences are found for the TCGT and CAAT positions.

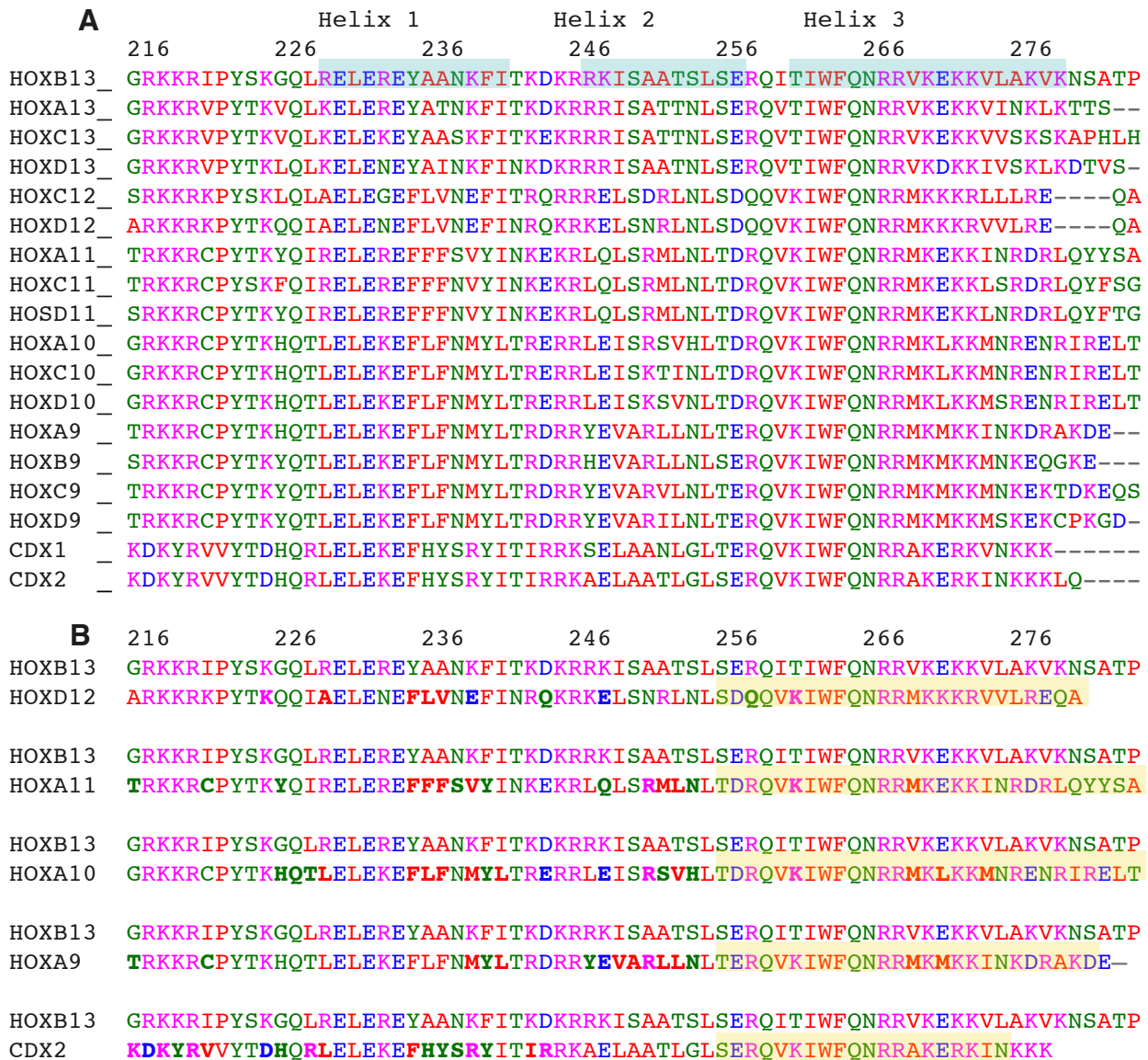


Figure S4. HOXB13 - HOXes/CDX mutations

(A) Sequence alignment of HOXB13 with other HOXes. Secondary structure (alpha-helices) of HOXB13 is highlighted in cyan. (B) The pairwise alignment; note that in addition to single mutations there are combined mutations and replacement of Helix 3 (DNA-binding helix) to corresponding helix of other HOXes (highlighted in yellow). The numbering on the top of the sequences is HOXB13 numbering.

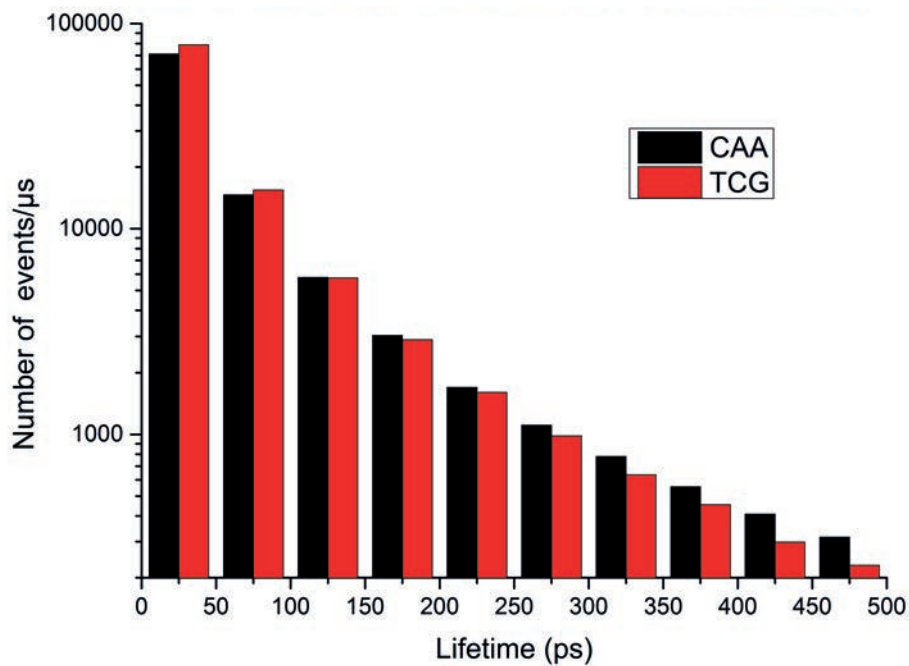


Figure S5. Distribution of water-bridge lifetimes in HOXB13:DNA complexes

Histogram showing the distribution of lifetimes of water bridges between the protein and the DNA for the HOXB13-DNA^{TCG} and HOXB13-DNA^{CAA} systems. The histogram is constructed by calculating the duration of each water bridge with 50 ps resolution from the molecular dynamics simulations; a water bridge is considered to exist when a water molecule is simultaneously hydrogen-bonded to one of the protein residues 255-272 and one of the DNA base pairs 5'-T(6)TTTACGAG(14)-3'.