

Local DNA Topography Correlates with Functional Noncoding Regions of the Human Genome

Stephen C. J. Parker,¹ Loren Hansen,^{1,2} Hatice Ozel Abaan,³ Thomas D. Tullius,^{1,4*} Elliott H. Margulies^{3*}

¹Bioinformatics Program, Boston University, Boston, MA 02215, USA. ²National Center for Biotechnology Information, National Institutes of Health, Bethesda, MD, USA. ³Genome Technology Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ⁴Department of Chemistry, Boston University, Boston, MA 02215, USA.

*To whom correspondence should be addressed. E-mail: elliott@nhgri.nih.gov (E.H.M.); tullius@bu.edu (T.D.T.)

The three-dimensional molecular structure of DNA, specifically the shape of the backbone and grooves of genomic DNA, can be dramatically affected by nucleotide changes, which can cause differences in protein binding affinity and phenotype. We developed an algorithm to measure constraint on the basis of similarity of DNA topography among multiple species using hydroxyl radical cleavage patterns to interrogate the solvent accessible surface area of DNA. This algorithm found that 12% of bases in the human genome are evolutionarily constrained—double the number detected by nucleotide sequence-based algorithms. Topography-informed constrained regions correlated with functional non-coding elements, including enhancers, better than regions identified solely on the basis of nucleotide sequence. These results support that the molecular shape of DNA is under selection and can identify evolutionary history.

Genomic sequences that code for proteins are relatively well understood, but comprise only ~2% of the human genome (1). Many functions are encoded in the remaining ~98% non-coding portion of the genome, but little is known about how functional non-coding information is specified (2). It has been hypothesized that functional regions are likely to be evolutionarily constrained because of their importance to the organism (3–11). Nonetheless, evolutionary sequence-constraint algorithms fail to identify many non-coding functional elements (12–15). This may be because these methods only analyze the primary nucleotide sequence of a genome (i.e., the order of A's, C's, G's, and T's). However, DNA is a molecule with a three-dimensional structure that varies according to the nucleotide sequence (16–18).

We employed a previously developed method on the basis of the hydroxyl radical cleavage pattern of DNA (19) to quantitatively evaluate how the structure of DNA varies throughout a genome. This approach can be used to predict the shape of the DNA backbone and grooves of genomic DNA at single-nucleotide resolution (20). We call this pattern the structural profile of a DNA region. Structural profiles

reveal that the relationship of DNA structure to the corresponding DNA sequence is not always simple. While similar sequences often adopt similar structures (Fig. 1A) divergent nucleotide sequences can have similar local structures (Fig. 1B) (20). Conversely (albeit less frequently), similar sequences can adopt very different structures (Fig. 1C). These observations indicate that DNA regions that differ on the basis of the order of nucleotides may be similar in structure, which suggests they may perform similar biological functions.

We quantified the effect of single-base substitutions on DNA structure by computing the structural profiles of all possible 11mer sequences (4,194,304 in total), measuring the similarity between profiles for all pairs of 11mers that differ only by a single substitution at the central nucleotide (21). A histogram of structural differences for all 11mers reveals a range of effects from minor to drastic (Fig. 1, D to G). We note that in some cases, the one-base change in an 11mer sequence has a minor effect, but the 11mer may have a dramatically different structure if another base is substituted. Since current sequence-constraint algorithms do not compute the effect of a base change on DNA structure, we developed a computer program, Chai, which incorporates structural information. This method works in a manner analogous to the DNA sequence-based binomial conservation (binCons) algorithm (8) but instead of computing the binomial probability of observed base substitutions between species, Chai calculates the difference between DNA structural profiles as a measure of similarity (21). We used the binCons and Chai algorithms to analyze 30 Mb of high-quality, comparative sequence data for 36 different species [the ENCODE pilot project regions (15, 22)]. We defined a false discovery rate (FDR) on the basis of a neutral (or null) alignment (10, 22) in an identical manner for assessing both sequence- and structure-informed conservation, so that each type of conserved region was identified with equal statistical confidence. We found that at any given FDR, the structure-informed Chai algorithm identified more evolutionarily

constrained bases compared to binCons (Fig. 2A and table S1). For example, at a 5% FDR, binCons identified 6.7% of bases as constrained, corresponding with previous findings of sequence-based mammalian constraint (8–11, 15, 22, 23). Chai identified nearly twice as many bases as constrained (12%) (Fig. 2B) while identifying 89% of the regions identified with the binCons method.

We next determined the extent to which evolutionarily constrained regions identified by Chai and binCons harbor functional elements. We examined DNase I hypersensitive sites (15) and predicted transcriptional enhancers identified using chromatin modification patterns (24) from published studies (21) as examples of non-coding functional elements that frequently exhibit little sequence similarity. Compared to binCons regions, we found that Chai-detected regions overlap a higher proportion of DNase I hypersensitive sites (78% for Chai versus 50% for binCons) and predicted enhancers (84% for Chai versus 59% for binCons) (Fig. 2C). To test whether the increased correlation with functional sequences was due to the additional territory identified by Chai relative to binCons, we tuned each algorithm for equal base coverage as opposed to equal statistical confidence and found that Chai-detected regions still overlap a significantly greater number of functional non-coding elements compared to binCons-detected regions ($P < 10^{-12}$; Fisher's exact test) (figs. S1, 2A, table S1).

Focusing our analysis on regions identified only by Chai and not by binCons (Chai-only regions) resulted in a statistically significant overrepresentation of noncoding functional sequences ($P < 0.01$; GSC statistic) [see (21) for a description of the GSC statistic] in the Chai-only regions and a statistically significant underrepresentation ($P < 10^{-9}$; GSC statistic) of coding regions (Fig. 2C). This suggests that some of the functional information in the non-coding portion of the genome is conferred by DNA structure as well as by the nucleotide sequence.

We thus examined if non-coding nucleotide substitutions inducing changes in DNA structure impact the biological function of a sequence. We focused on the DNA-binding properties of the Zif268 protein, a mammalian transcription factor that consists of three zinc fingers that wrap around DNA in the major groove (25). Binding affinity data (26) were compared with structural profiles for the 15 sites identified to bind wild type Zif268 and the Zif268 REDV mutant (21). For both proteins the structural profiles of the high-affinity sequence motifs were similar, while low-affinity motifs had a different structural profile (figs. S2A and S3A).

We ranked each Zif268 REDV motif by the extent of DNA structural difference relative to the best binding site, and found that binding affinity correlated with this ranking ($r = 0.75$; $P = 6.9 \times 10^{-4}$; t test; fig. S2B). In other words, low-affinity binding sites differed dramatically in structure from

high-affinity sites, while sites with intermediate affinity showed less structural difference. We observed a similarly high correlation for the wild type protein, which has different binding preferences ($r = 0.74$; $P = 7.5 \times 10^{-4}$; t test; fig. S3B). Analysis of DNA-binding site structural profiles and binding affinity for the archaeal transcriptional regulator Ss-LrpB revealed similar trends to those found for Zif268 (fig. S4).

We next tested if mutations that change the molecular topography of non-coding genomic regions have phenotypic consequences with data from the PhenCode Project (27); which organizes information from several locus-specific mutation databases. We gathered 734 non-coding single-nucleotide variants in the human genome associated with a phenotype. For each variant, we calculated the difference in the structural profile between the mutant and non-mutant sequence over an 11-bp window centered on the variant nucleotide (similar to the analysis above). For comparison, a distribution of baseline variation in DNA topography was computed for 16,832 neutrally-evolving single-nucleotide polymorphisms (SNPs) (21). The phenotype-associated distribution was significantly correlated with larger changes in structure ($P < 3 \times 10^{-4}$; Wilcoxon rank sum test) relative to the baseline distribution, and contained an additional high structure-change peak (Fig. 3). These results indicate that phenotype-associated mutations tend to induce larger changes in the structural profile of non-coding DNA compared to baseline neutral variation.

We also identified phenotype-associated SNPs causing the top 10% structural changes and the lowest 10%. The phenotype-associated mutations with the highest changes in DNA structure occurred significantly more often ($P < 10^{-2}$, Fisher's exact test) in evolutionarily constrained regions of the genome (56% for high structure-change regions versus 29% for low structure-change regions) (21). This suggests that non-coding DNA may be under selective constraint, which may prevent changes in DNA structure. Because the severity of structural change might help identify functional SNPs, we constructed a database of changes in the structural profile for all known SNPs in the human genome (21).

Finally, we demonstrated that structural changes affect biological function in non-coding evolutionarily constrained regions identified by the Chai algorithm. We chose 12 predicted enhancer-containing regions (24): 5 regions overlap elements only detected by the Chai algorithm; 7 regions overlap a combination of Chai- and binCons-detected elements (table S2). We cloned the 300 bp surrounding each of these genomic regions in a luciferase reporter construct and transfected them into 293T cells. Eight of the twelve constructs displayed luciferase activity that was significantly greater ($P \leq 0.05$; Wilcoxon rank sum test) than that of random control sequences (Fig. 4). Three of the 5 constructs

only overlapping Chai-detected elements were positive (table S2).

Given the plethora of regulatory functions that a genome encodes and the three-dimensional genomic architecture required to orchestrate these events (28), it may not be surprising that there is widespread conservation of local DNA topography. Perhaps only a subset of local structural configurations can accommodate the functional requirements of a genomic locus [for example see (29)]. Once the molecular topography of a locus is permissive to a regulatory function, this structure may be maintained within the genome. Our high-resolution topography-based constraint-detection method reveals that structure-informed constraint is widespread in the human genome, and that these regions overlap known non-coding functional sites. Because different DNA sequences can have similar local structures (20), these regions might escape detection with sequence-based conservation-identification methods.

References and Notes

1. International Human Genome Sequencing Consortium, *Nature* **431**, 931 (2004).
2. M. D. Wilson *et al.*, *Science* **322**, 434 (2008).
3. L. Elnitski *et al.*, *Genome Res.* **13**, 64 (2003).
4. M. Kellis, N. Patterson, M. Endrizzi, B. Birren, E. S. Lander, *Nature* **423**, 241 (2003)
5. G. G. Loots *et al.*, *Science* **288**, 136 (2000).
6. L. A. Pennacchio, E. M. Rubin, *Nat. Rev. Genet.* **2**, 100 (2001).
7. W. W. Wasserman, M. Palumbo, W. Thompson, J. W. Fickett, C. E. Lawrence, *Nat. Genet.* **26**, 225 (2000).
8. E. H. Margulies, M. Blanchette, D. Haussler, E. D. Green, *Genome Res.* **13**, 2507 (2003).
9. A. Siepel *et al.*, *Genome Res.* **15**, 1034 (2005).
10. G. M. Cooper *et al.*, *Genome Res.* **15**, 901 (2005).
11. S. Asthana, M. Roytberg, J. Stamatoyannopoulos, S. Sunyaev, *PLoS Comp. Biol.* **3**, e254 (2007).
12. S. Fisher, E. A. Grice, R. M. Vinton, S. L. Bessling, A. S. McCallion, *Science* **312**, 276 (2006).
13. D. M. McGaughey *et al.*, *Genome Res.* **18**, 252 (2008).
14. H. M. Petrykowska, C. M. Vockley, L. Elnitski, *Genome Res.* **18**, 1238 (2008).
15. ENCODE Project Consortium, *Nature* **447**, 799 (2007).
16. T. Ohyama, Ed., *DNA Conformation and Transcription* (Landes Bioscience/Eurekah.com, Georgetown, TX, 2005).
17. W. K. Olson, A. A. Gorin, X. J. Lu, L. M. Hock, V. B. Zhurkin, *Proc. Natl. Acad. Sci. U.S.A.* **95**, 11163 (1998).
18. R. E. Dickerson, *Methods Enzymol.* **211**, 67 (1992).
19. M. A. Price, T. D. Tullius, *Methods Enzymol.* **212**, 194 (1992).
20. J. A. Greenbaum, B. Pang, T. D. Tullius, *Genome Res.* **17**, 947 (2007).
21. Materials and methods are available as supporting material on *Science Online*.
22. E. H. Margulies *et al.*, *Genome Res.* **17**, 760 (2007).
23. Mouse Genome Sequencing Consortium, *Nature* **420**, 520 (2002).
24. N. D. Heintzman *et al.*, *Nature Genet.* **39**, 311 (2007).
25. N. P. Pavletich, C. O. Pabo, *Science* **252**, 809 (1991).
26. M. L. Bulyk, P. L. F. Johnson, G. M. Church, *Nucleic Acids Res.* **30**, 1255 (2002).
27. B. Giardine *et al.*, *Human Mutation* **28**, 554 (2007).
28. T. Misteli, *Bioessays* **27**, 477 (2005).
29. R. Joshi *et al.*, *Cell* **131**, 530 (2007).
30. We thank E. D. Green and L. C. Brody for feedback on the manuscript, E. Bishop and D. Landsman for discussion, and G. K. McEwen for assistance with experimental steps. Funded by the National Human Genome Research Institute (NHGRI) of the National Institutes of Health (NIH) (R01 HG003541) to T.D.T. E.H.M. was supported by the Intramural Research Program of the NHGRI, NIH. S.C.J.P. was supported by a National Academies Ford Foundation Dissertation Fellowship.

Supporting Online Material

www.sciencemag.org/cgi/content/full/1169050/DC1

Materials and Methods

Figs. S1 to S5

Tables S1 and S2

References

26 November 2008; accepted 23 February 2009

Published online 12 March 2009; 10.1126/science.1169050

Include this information when citing this paper.

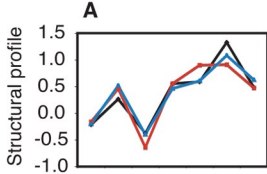
Fig. 1. DNA topography vs. nucleotide sequence. (A to C) Hydroxyl radical cleavage patterns (referred to here as the Structural profile) and corresponding color-matched sequence alignments (below). Asterisks indicate identical columns in the sequence alignment. (D to G) Single-base substitutions have a range of effects on local DNA structure. For all possible 11 mer nucleotide sequences, we computationally changed the middle position to each possible base and measured its effect on the structural profile. These changes were quantitatively measured by calculating the mean Euclidean distance (referred to here as the Average structural change) (21), where low values indicate similar structure profiles and large values indicate different structure profiles. (G) The distribution of average structure changes observed for all 11 mers is shown. Arrows indicate what we classify as low, moderate, or high average structure changes, with representative 11 mers for each shown in (D), (E), and (F), respectively. For (D) to (F), the y axis indicates the structure profile at each nucleotide position in the 11-mer (x axis). The

structure profile for each 11-mer containing a different base in the middle position (noted by an 'N') is plotted in a different color. Note how the structure profile patterns are increasingly different from low to high structure change [(E) to (G)].

Fig. 2. Structure-informed evolutionarily constrained regions. (A) Plot showing the fraction of bases identified as evolutionarily constrained at various False Discovery Rates (FDRs) (21) by the DNA sequence-based algorithm binCons (blue line) and the structure-informed algorithm Chai (red line). (B) Venn diagram showing nucleotides identified by each algorithm [Chai (red) and binCons (blue)] at 5% FDR [gray vertical line in (A)]. Chai detects nearly the same set of bases that binCons identifies (intersection), plus almost twice as many more (Chai-only). (C) Bar graph showing the fraction of different types of genomic regions that overlap binCons- (blue), Chai- (red), and Chai-only-detected (white) constrained elements. Black points are the mean of a null distribution constructed with the GSC method (21); error bars represent 95% confidence intervals.

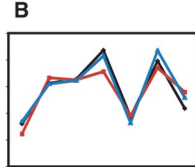
Fig. 3. The distribution of DNA structural changes for phenotype-associated variants (red line) and neutral variants (black line) (21) with the real structure prediction algorithm (A) and a randomized version (B). (A) The phenotype-associated distribution was shifted significantly to the right ($P < 3 \times 10^{-4}$; Wilcoxon rank sum test) and contains an additional high structure change peak relative to the neutral distribution. (B) No significant shift ($P > 0.05$; Wilcoxon rank sum test) was observed with a randomized version of the structure prediction algorithm.

Fig. 4. Luciferase-based reporter activity of 12 regions containing Chai-detected elements (21). These regions overlapped predicted enhancer regions (see text). Plotted for each element is the luciferase activity relative to the median activity from 100 random control constructs (y axis; see fig. S5). Error bars represent one standard deviation from the mean of four experimental replicates and asterisks denote $P \leq 0.05$; Wilcoxon rank sum test.



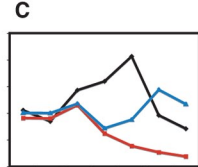
Sequence 1 A T A C G C G
 Sequence 2 A T A G G C G
 Sequence 3 A T A T G C G
 * * * * * * *

% identity 86%



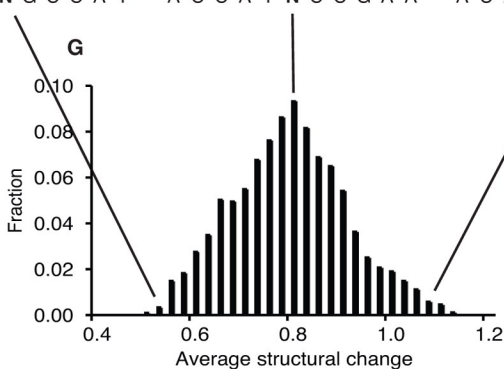
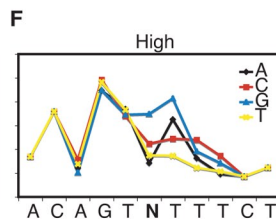
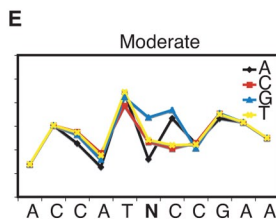
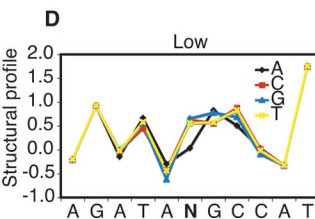
A C G T A C C
 A G G G A G C
 A T G C A T C
 * * * *

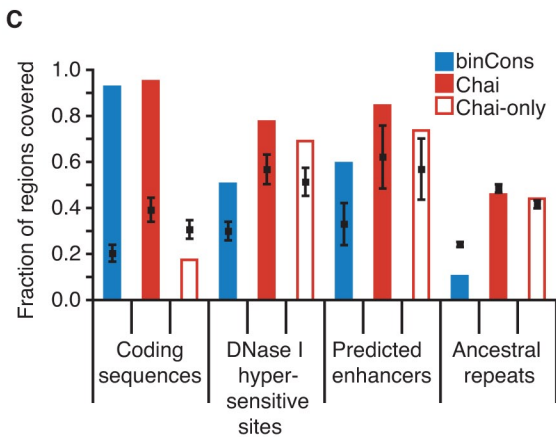
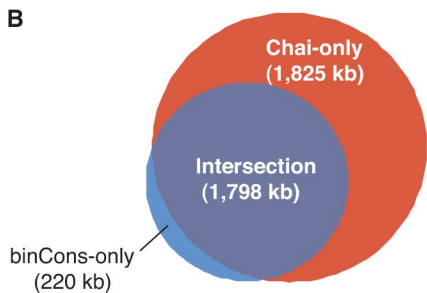
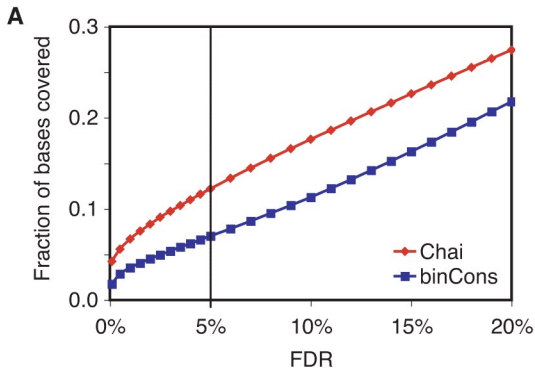
57%



A A T G T T T
 A A T T T T T
 A A T C T T T
 * * * * * *

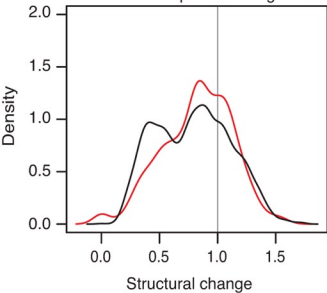
86%





A

Real structure prediction algorithm

**B**

Random structure prediction algorithm

