# The Homeodomain Resource: sequences, structures, DNA binding sites and genomic information

**Sharmila Banerjee-Basu, Daniel W. Sink and Andreas D. Baxevanis\***

Genome Technology Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD 20892, USA

## ABSTRACT

**The Homeodomain Resource is an annotated collection of non-redundant protein sequences, three-dimensional structures and genomic information for the homeodomain protein family. Release 3.0 contains 795 full-length homeodomain-containing sequences, 32 experimentally-derived structures and 143 homeobox loci implicated in human genetic disorders. Entries are fully hyperlinked to facilitate easy retrieval of the original records from source databases. A simple search engine with a graphical user interface is provided to query the component databases and assemble customized data sets. A new feature for this release is the addition of DNA recognition sites for all human homeodomain proteins described in the literature. The Homeodomain Resource is freely available through the World Wide Web at http:// genome.nhgri.nih.gov/homeodomain.**

## INTRODUCTION

The homeodomain is a common DNA-binding structural motif found in many eukaryotic regulatory proteins (1,2). Homeodomain proteins are involved in the transcriptional control of many developmentally important genes, and 143 human loci have been linked to various genetic and genomic disorders. X-ray crystallographic and NMR spectroscopic studies (3–11) on several members of this family have revealed that the homeodomain motif is comprised of three α-helices that are folded into a compact globular structure. Helices I and II lie parallel to each other and across from the third helix. This third helix is also referred to as the 'recognition helix', as it confers the DNA-binding specificity of individual homeodomain proteins. The homeodomain has been evolutionarily conserved at the structural level (12); this is most evident upon examination of divergent members of the homeodomain family.

The Homeodomain Resource represents a comprehensive collection of information about the homeodomain family. The database contains all available full-length and homeodomain-only sequence data and structures as of October, 2000. The genetic data contained in this database includes information on human diseases in which homeodomain-containing proteins are implicated, cytogenetic map locations and specific mutation data underlying the disease condition. Since its last release

(13), 27 new loci for genetic disorders have been identified. The sequence information within the database is automatically updated on a monthly basis, with each entry in the database rigorously selected to assure non-redundancy.

## DATABASE DESCRIPTION

The current version of the database contains 795 full-length homeodomain protein sequences isolated from 83 different species (Table 1). The complete full-length sequence data as well as the homeodomain portion of the sequence is available in FASTA format. The database can be searched on the basis of SWISS-PROT ID, GenBank accession number, gene names (both common and alternative), protein description, sequence and organism name. The search engine also supports Boolean queries, allowing users to search on individual fields or on all fields at once. Search results are returned in a tabular format, with hyperlinks to the original records in GenBank and SWISS-PROT, respectively. Individual sequences can also be retrieved in FASTA format from a pop-up window.

**Table 1.** Homeodomain Resource statistics

| | |
|---|---|
| Total sequences available | 3796 |
| Non-redundant full-length sequences | 795 |
| Genes/gene symbols | 375 |
| Distinct organisms | 83 |
| Three-dimensional structures | 32 |
| Homeobox loci implicated in human genetic disorders | 143 |

The genetic information available for the homeodomain protein family has increased ~23% since the last release, with the inclusion of 27 new loci. The genetic data are compiled from both the literature and from the Online Mendelian Inheritance in Man (OMIM) database at NCBI (http://www.ncbi.nlm.nih.gov/ Omim/). As before, a search engine is available for querying this genomic information. Search results are presented in a tabular format with hyperlinks to the original records and can be sorted by disease name, map location, gene symbol, protein name or OMIM identifier. This value-added format allows users easy access to related information.

\*To whom correspondence should be addressed. Tel: +1 301 496 8570; Fax: +1 301 402 6858; Email: andy@nhgri.nih.gov

| Common Name | DNA Binding Sequence | Common Name | DNA Binding Sequence |
|---|---|---|---|
| ARIX | G T C **A A T T A** G | OCT 1 | A A **A T G C A A A T** A C C |
| BARX2 | Y Y **T A A T G R** T T T T Y | OCT 2 | **A T G C A A A T** |
| CDX1 | C **T T T A T** G | OCT3 | **A T G C A A A T** |
| CDX2 | A · · · · A / C **T T T A T** G | OTX2 | C C **T A A T** C C T G G G **T T A T** C |
| CDX4 | A / **T T T A T** T | PAX2 | C · · G / **T C A T G C A T G A** C |
| CRT1 | **T A A T** N N N **A T T A** | PAX3 | **T A A T** N N **A T T A** |
| CRX | C / **T A A T C A** | PAX4 | C / A N N N **T C A C C** C |
| CSX | **T N A A G T** G | PAX5 | A · · C · · · · · T / G **T C A C G C T T G A T G** C |
| DLX2 | C A C **T A A T T G A** G | PAX6 | A · · C · · · · TA · T / A A N N **T T C A C G C T T G A T G C A** C |
| DLX3 | C A N **T T A T C T C A** G | PAX7 | **T A A T T** |
| DLX4 | T C **A A T T A A T T G** A | PBX1 | **T T G A T T G A** T |
| DLX5 | **A T A A T T A** G | PBX1-HOXA10 | **A T G A T T T A T G** A |
| DLX6 | **A T A A T T A** G | PBX1-HOXA9 | T / **T G A T T T A** C |
| EMX2 | C A C **T A A T T G A** G | PBX1-HOXB1 | **A T G A T T G A T** C G |
| EVX1 | **T N A T T A** N N N N N **T A A T** N G | PBX1-HOXB4 | **A T G A T T G A T** G A |
| GBX2 | A G T G A G **T A A T T G** G | PBX1-HOXB6 | **A T G A T T T A T T** A |
| HEX | T · · · T / **C A A G** | PBX1-HOXB7 | **A T G A T T T A T G** C T C T A |
| HK31 | **T A A G T** A | PBX1-HOXB8 | **A T G A T T T A T G** A C T C T A |
| HLX1 | T C **T A T T A A T T G** A | PBX1-HOXB9 | **A T G A T T T A** C G A C |
| HME1 | T C T G **T A A T T A** C A | PBX1-HOXC6 | G · · · · · · A / **A T G A T T T A T T** G C T T T |
| HNF1 | A G T **T T A T A T T T G A** C A | PBX1-HOXC6 | A / G **T G A T T T A T T A** C T T T |
| HNF4 | C G C T T G G C A A A G G T C A C C T | PBX1-MEIS1 | G / **T G A T T G A** C A G C T |
| HNF6 | G A **T A T T G A T** T T T T | PBX2 | **T T G A T T G A** T |
| HOX 11 | C G **T T A A T T G** G | PBX2-HOXC6 | G **T G A T T T A T T** G C T T T |
| HOX11L1 | C G **T T A A T T G** G | PBX3 | **T T G A T T G A** T |
| HOX11L2 | C G **T T A A T T G** G | PBX3-HOXC6 | G / G **T G A T T T A T T A** C T T T |
| HOXA10 | T / **T A T T G** A | PITX1 | **T A A T C** C |
| HOXA9 | T / **T T T C** | PITX2 | **T A A T C** C |
| HOXB7 | **T T T A T G A** C | PMX1 | **T A A T** A |
| HOXB8 | **T T T A T G A** C | PNX1 | **T G A C A** G |
| HOXB9 | T / **T T A A C G A** C | POU1F1 | **T A A A** T |
| HOXC11 | T / **T T T A** C | POU2F1 | **A T G C A A A T** |
| HOXD12 | T / **T T A A C G A** C | POU3F1 | T G G A G A N **A A T A G A** G |
| HOXD13 | **T T T A C G A** G | POU3F2 | **A T G C A A A T** |
| IPF1 | T / C T C **T A A T G** G G | POU3F4 | **A T G C A A A T** |
| ISL-1 | G C **T T A A T A T C T** G | POU4F1 | A · · · · · · T / G C **T C A T T A A** C |
| ISL-2 | T G G / **T T A A G T** A | POU4F2 | A · · · · · · T / G C **T C A T T A A** C |
| LH2 | **T T A A T T** A | POU4F3 | A · · · · · · T / G C **T C A T T A A** C |
| MEI1 | G / **T G A C A G** C T | POU5F1 | **A T G C A A A T** |
| MOX2 | T C **A A T T A A T T G** A | POU6F1 | **T A A T G A G C T G C X T A A T** |
| MSX-1 | C **T A A T T** G | SATB1 | C / A **T A A T** A |
| MSX-2 | T C **A A T T A A T T G** A | TCF-1 | A A **C A A A** G |
| N-OCT5 | **A T G C A A A T** | TGIF | **T G A C A** A |
| NK61 | **T T A A T T A** C | TTF-1 | T C **A A G T G** T T |
| NKX2.2 | T · · · · · · · A C / T C **A A G T G G** T T | | |

**Figure 1.** DNA binding sites for human homeodomain loci implicated in genetic or genomic disorders. Data is derived from the published literature, citations in OMIM, and entries for DNA-bound homeodomain structures from the Protein Data Bank. The core regions of each of the DNA binding sites are shown in bold type. An expanded version of this table is available online (at http://genome.nhgri.nih.gov/homeodomain); the expanded version includes a listing of alternative gene names, as well as citation information.

A new feature of the Homeodomain Resource is the inclusion of DNA binding sites for human homeodomain loci implicated in genetic or genomic disorders (Fig. 1). These data were obtained by extensive review of the published literature, citations in OMIM, and entries for DNA-bound homeodomain structures from the Protein Data Bank. The online version of Figure 1 (available via http://genome.nhgri.nih.gov/homeodomain) includes alternate gene names and references to the primary citation from which the information was retrieved.

## REFERENCES

1. Gehring,W., Qian,Y., Billeter,M., Furukubo-Tokunaga,K., Schier,A., Resendez-Perez,D., Affolter,M., Otting,G. and Wuthrich,K. (1994) Homeodomain-DNA recognition. *Cell*, **78**, 211–223.

2. Laughon,A. (1991) DNA binding specificity of homeodomains. *Biochemistry*, **30**, 11357–11367.
3. Dekker,N., Cox,M., Boelens,R., Verrijzer,C., van der Vliet,P. and Kaptein,R. (1993) Solution structure of the POU-specific DNA-binding domain of Oct-1. *Nature*, **362**, 852–855.
4. Viglino,P., Fogolari,F., Formisano,S., Bortolotti,N., Damante,G., Di Lauro,R. and Esposito,G. (1993) Structural study of rat thyroid transcription factor 1 homeodomain (TTF-1HD) by nuclear magnetic resonance. *FEBS Lett.*, **336**, 397–402.
5. Gruschus,J., Tsao,D., Wang,L., Nirenberg,M. and Ferretti,J. (1997) Interactions of the vnd/NK-2 homeodomain with DNA by nuclear magnetic resonance spectroscopy: basis of binding specificity. *Biochemistry*, **36**, 5372–5380.
6. Kissinger,C., Liu,B., Martin-Blanco,E., Kornberg,T. and Pabo,C. (1990) Crystal structure of an engrailed homeodomain-DNA complex at 2.8 Å resolution: a framework for understanding homeodomain-DNA interactions. *Cell*, **63**, 579-590.
7. Ceska,T., Lamers,M., Monaci,P., Nicosia,A., Cortese,R. and Suck,D. (1993) The X-ray structure of an atypical homeodomain present in the rat liver transcription factor LFB1/HNF1 and implications for DNA binding. *EMBO J.*, **12**, 1805–1810.
8. Liu,B., Kissinger,C. and Pabo,C. (1990) Crystallization and preliminary X-ray diffraction studies of the engrailed homeodomain and of an engrailed homeodomain/DNA complex. *Biochem. Biophys. Res. Commun.*, **171**, 257–259.
9. Qian,Y., Billeter,M., Otting,G., Muller,M., Gehring,W. and Wuthrich,K. (1989) The structure of the Antennapedia homeodomain determined by NMR spectroscopy in solution: comparison with prokaryotic repressors. *Cell*, **59**, 573–580.
10. Qian,Y., Furukubo-Tokunaga,K., Resendez-Perez,D., Muller,M., Gehring,W. and Wuthrich,K. (1994) Nuclear magnetic resonance solution structure of the fushi tarazu homeodomain from Drosophila and comparison with the Antennapedia homeodomain. *J. Mol. Biol.*, **238**, 333–345.
11. Wolberger,C., Vershon,A., Liu,B., Johnson,A. and Pabo,C. (1991) Crystal structure of a MAT alpha 2 homeodomain-operator complex suggests a general model for homeodomain-DNA interactions. *Cell*, **67**, 517–528.
12. Buerglin,T. (1994) In Duboule,D. (ed.), *Guidebook to the Homeobox Genes*. Oxford University Press, Oxford, UK, pp. 25–71.
13. Banerjee-Basu,S., Ryan,J.F. and Baxevanis,A. (2000) The Homeodomain Resource: sequences, structures and genomic information. *Nucleic Acids Res.*, **28**, 329–330.