# Sequence and pharmacophore analysis of DNA recognition helices in HTH family of proteins

**A PROJECT REPORT**

*Submitted by:*

**Suvan Kumar Sahu**

**(235HSBB034)**

*to*

**Institute of Bioinformatics and Applied Biotechnology**

*in partial fulfilment of the requirements for*

Master of Science in Biotechnology and Bioinformatics (Degree to be awarded by Bangalore University, Bengaluru)

*under the guidance of:*

**Prof. S Thiyagarajan**



IBAB

Institute of Bioinformatics and Applied Biotechnology

**INSTITUTE OF BIOINFORMATICS AND APPLIED BIOTECHNOLOGY BENGALURU**

February 2025

# DECLARATION

I certify that

a) the work contained in this report is original and has been done by me under the guidance of my supervisor(s).

b) I have followed the guidelines provided by the Department in preparing the report.

c) I have conformed to the norms and guidelines given in the Honor Code of Conduct of the Institute.

d) whenever I have used materials (data, theoretical analysis, figures, and text) from other sources, I have given due credit to them by citing them in the text of the report and giving their details in the references. Further, I have taken permission from the copyright owners of the sources, whenever necessary.

*Signed by:*

**Suvan Kumar Sahu**

# CERTIFICATE

It is certified that the work contained in this report titled "Sequence and pharmacophore analysis of DNA recognition helices in HTH family of proteins" is the original work done by Suvan Kumar Sahu and has been carried out under my supervision.

*Signed by:*

Prof. S Thiyagarajan

**Date:** 23.02.25

# **<u>Abstract</u>**

The helix-turn-helix (HTH) motif is a key DNA-binding domain found in transcription factors, playing a crucial role in regulating genes. Despite its importance, the diversity in the structure and sequence of these motifs makes it difficult to predict how they recognize and bind to specific DNA sequences. This study aims to analyse the sequences of HTH proteins and model their chemical properties to understand how they interact with DNA.

Using a large collection of HTH proteins from NAKB database, we grouped these proteins based on their mode of binding to DNA. We created a new pharmacophore-based scoring matrix to identify the key chemical features that drive these interactions.
Our findings describe the unique roles of certain amino acids, such as cysteine and tryptophan, in determining binding specificity. While the study provides valuable insights, it also highlights some limitations in the current approach, paving the way for future improvements. By bridging the gap between sequence analysis and pharmacophore modelling, this research enhances our understanding of how HTH proteins interact with DNA and opens new possibilities for developing therapies targeting gene regulation.

# TABLE OF CONTENTS

# <u>List of figures</u>

# List of Abbreviations

| Abbreviation | Description |
| --- | --- |
| DNA | Deoxy-ribonucleic acid |
| RNA | Ribonucleic acid |
| HTH | Helix-turn-helix |
| PDB | Protein Data Bank |
| CAP | Catabolite gene activator protein |
| NAKB | Nucleic Acid Knowledge Base |
| CSV | Comma-separated values |
| CRP | Cyclic-AMP Receptor |
| Fis | Factor for inversion stimulation |
| BLOSUM | Block Substitution Matrix |

# Introduction

## 1.1 Helix-turn-helix motifs

DNA-binding motifs are structural domains present within proteins that allow them to interact with DNA, thus playing an important role in the regulation of gene expression, replication, and DNA repair. These motifs have a unique ability to recognize and bind to specific DNA sequences through highly specific molecular interactions (Luscombe et al., 2000). One of the most common and well-characterized DNA-binding motifs is the helix-turn-helix (HTH), which is widely found in transcription factors and plays a significant role in gene regulation in both prokaryotic and eukaryotic organisms (Brennan et al., 1989).

The HTH motif typically comprises of two alpha helices that are connected by a short, flexible turn (Pabo & Sauer, 1984). The first helix, known as the stabilizing helix, functions as a structural scaffold that helps stabilize the protein-DNA complex, while the second helix, known as the recognition helix, makes sequence-specific contacts with the base pairs located in the major groove of the DNA (Anderson et al., 1981). The stabilizing helix interacts with the minor groove of the DNA, thereby assisting in proper positioning the recognition helix so that it fits into the major groove, while the turn provides the necessary flexibility to maintain the overall structure of the HTH motif (Steitz, 1990).

The recognition helix carries out sequence-specific interactions with DNA bases, through hydrogen bonding, van der Waals forces, and electrostatic interactions (Pabo & Sauer, 1984). The sequence and structural variations within the recognition helix determine the binding specificity to DNA motifs, which ultimately influences gene regulatory mechanisms (Brennan, 1992). As a result, performing sequence analysis on the recognition helix can reveal conserved residues that are critical for both DNA binding and sequence specificity, offering important insights into protein-DNA interactions and gene regulation (Rosinski et al., 1999).

## 1.2 Mode of DNA binding

Different HTH proteins have different modes of mechanism to bind to the DNA. Some bind as monomers while some bind as dimers. Out of the dimers, some are homodimers and some are heterodimers (Harrison & Aggarwal, 1990). Apart from this, the HTH proteins can also be classified on the basis of the site on DNA to which they bind – groove binding or phosphate backbone binding.

## 1.3 Role of Pharmacophore modelling

Pharmacophore properties refer to the key structural and chemical features of a molecule (such as a ligand or protein) that are responsible for its biological activity by enabling it to interact with a specific target (protein, enzyme, or receptor). These include properties like aromaticity, hydrogen bond donors, hydrogen bond acceptors, positive and negative charges, etc. (Guner, 2000).

Investigating the sequence and pharmacophore characteristics of these proteins is critical to understanding their DNA-binding mechanisms and sequence specificity. Advances in computational biology have made it possible to systematically analyse HTH protein sequences, revealing conserved motifs, crucial residues, and structural patterns linked to their functional roles.

Furthermore, pharmacophore modelling, which identifies the spatial arrangement of chemical features required for molecular interactions, offers a powerful means to study how HTH proteins bind to DNA and how small molecule inhibitors could disrupt this interaction. Such details are fundamental not only to study protein-DNA interactions but also for developing therapeutic strategies for diseases associated with transcriptional dysregulation. By combining sequence analysis with pharmacophore modelling, this research seeks to provide a comprehensive understanding of the structural and functional determinants of DNA-binding HTH proteins, paving the way for future studies into their biological significance and drug development potential.

## 1.4 Objectives

- Perform Sequence analysis of helix-turn-helix proteins
- Identify the pharmacophore pattern in the HTH recognition helix
- Derive a new scoring matrix based on pharmacophore properties
- Identify the nucleotide sequence or pattern to which HTH binds

# Chapter 2
# **Literature Survey**

## 2.1 Previous knowledge

The helix-turn-helix (HTH) motif is a structural domain present in most DNA-binding proteins and is essential for the regulation of transcription. Since the discovery, many studies have been dedicated to understanding the sequence and structural characteristics of this motif, providing details about its mechanisms of recognizing DNA and its evolutionary role.

Initial investigations described the structural mechanism for DNA recognition using the HTH motif, focusing on the importance of the recognition helix, which typically inserts into the major groove of the DNA to form sequence-specific hydrogen bonds and van der Waals interactions with the nucleotide bases (Brennan, 1992). It set the ground for what can affect protein-DNA binding specificity based on DNA-binding helices' sequence diversity. But then, sequence-based HTH motif detection was in its early stages, and there was an obvious need for computational methods that could detect these motifs with greater accuracy (Brennan, 1992).

The helix-turn-helix motif was first discovered in lambda repressor and Cro proteins from lambda bacteriophage. It consists of two alpha helices that are separated by a short turn. The two alpha helices are nearly perpendicular to each other. The HTH motifs have an average length of 22 amino acids, with a minimum length of 20 and maximum length of 32. The DNA-binding domains containing the HTH motifs are much larger (25-150 residues) (Pellegrini-Callace et al., 2005).

Experiments have shown that the two α-helices project out from the surface of the protein, creating a convex unit that is inserted into the DNA groove, and the phosphate backbone and recognition helix interactions help in stabilizing the DNA-binding process (Steitz et al., 1982). Early structural characterization of bacteriophage λ lambda repressor and Cro proteins initially recognized the HTH motif as a universally conserved DNA-binding component (Anderson et al., 1981). Further studies have determined that the HTH structure is quite well-conserved but its sequence is very variable among protein families so it is challenging to identify by sequence-based approaches alone (McLaughlin et al., 2003).

Several HTH-containing proteins have been found to create structural distortions in DNA, such as smooth bends and sharp kinks, that can have functional importance for transcriptional control. DNA bending helps in strengthening of protein-protein interactions and to attract transcription machinery in proteins like CAP (catabolite gene activator protein) and Fis (factor for inversion stimulation) (Brennan, 1992).

Even though the HTH motif is mainly related to DNA-binding proteins, research has also indicated that some RNA-binding proteins have HTH-like folds, suggesting a wider functional application outside of DNA binding (Anantharaman et al., 2010).

Due to the variability in the sequence of HTH motifs, most research indicates that structure-based methods are more accurate than sequence-based methods to identify HTH-containing proteins (Pellegrini-Callace et al., 2005). Through the combination of structure and sequence data, computational models have greatly enhanced the identification of HTH motifs and have enabled the identification of new DNA-binding proteins (Pellegrini-Callace et al., 2005).

## 2.2 Research gaps

In spite of these developments, there are a number of research gaps. One of the gaps is that the sequence of the recognition helix in HTH motifs is highly variable among different HTH-containing proteins. Therefore, no exact pattern or consensus sequence has been identified that characterizes the recognition helices of all HTH proteins
Such variability complicates the prediction of DNA-binding specificity through sequence analysis alone, restricting the accuracy of computational models for motif detection.

Addition to this, the precise DNA motifs to which HTH proteins bind are not fully characterized, as binding specificity appears to be influenced not only by the recognition helix sequence but also by the context of the flanking residues, protein-protein interactions, and chromatin environment.

Another significant gap in the field is the lack of a complete understanding of the various modes of DNA binding adopted by different types of HTH proteins. While the classical model involves insertion of the recognition helix into the major groove of the DNA, structural studies have revealed alternative binding orientations and multi-protein assemblies that contribute to gene regulation, yet these mechanisms remain incompletely understood.

# Materials and Methods

## 3.1 Sequence Collection

The list of proteins containing DNA-binding helix-turn-helix motifs were obtained from NAKB database. NAKB contains a total of 14320 annotated protein structures out of which 6193 are regulatory proteins. 2693 are transcription factors and out of that 1579 are DNA-binding. Out of the 1579 DNA-binding transcription factors, 765 are helix-turn-helix DNA-binding transcription factors, which are again subdivided into 5 different structurally different types – homeodomain(213), lambda repressor-like(103), tetR family(38), trp repressor(4), and winged helix/forkhead(402).
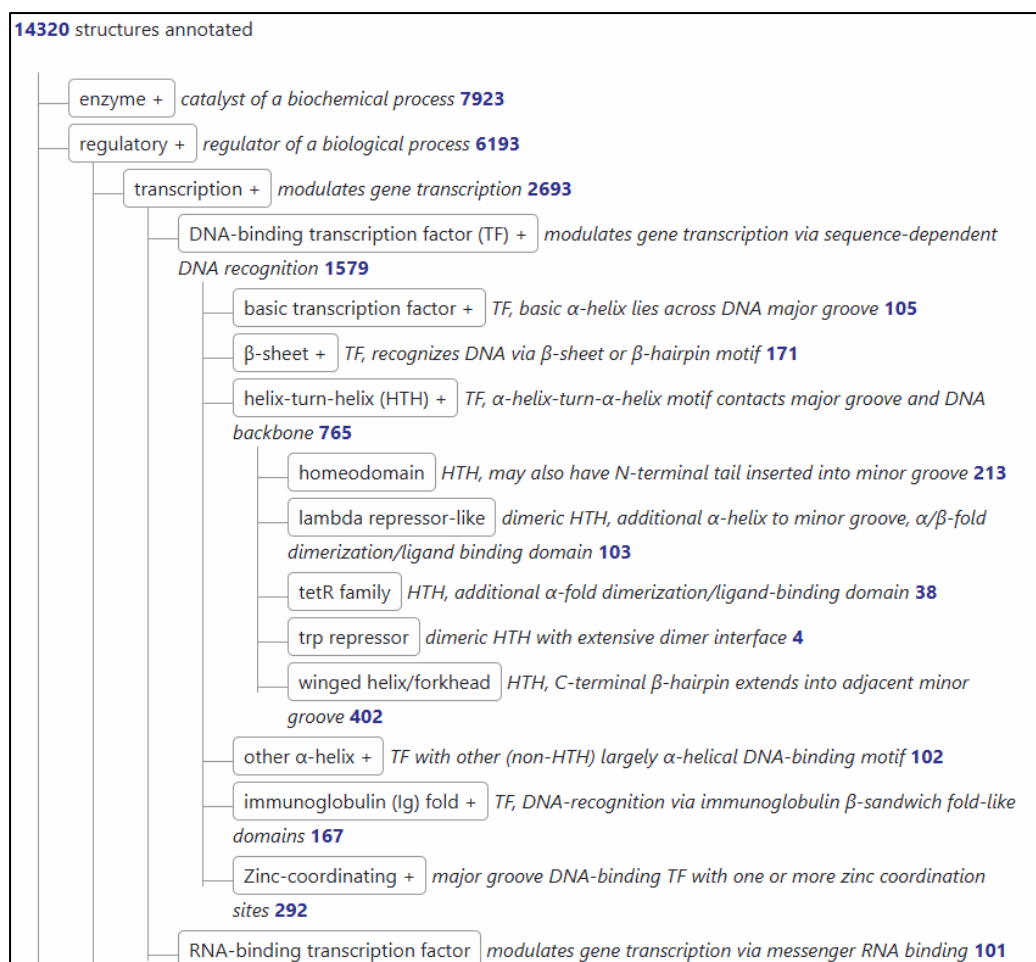


**Fig.1** Protein Annotation tree in NAKB

All of these HTH proteins are well-annotated and present as protein-DNA complexes. The full list of PDB-IDs for all the available DNA-binding HTH proteins was downloaded in the form of a CSV file. Then using an in-house developed python script, the protein sequences for all the proteins mentioned in the csv file were downloaded and stored in a FASTA file.

## 3.2 Classifying the proteins based on mode of DNA binding

The HTH protein sequences obtained from the NAKB database were then separated into groups based on their mode of DNA-binding. They were separated as monomer or dimer, homodimers or heterodimers, and binding to groove or to the phosphate backbone. This classification was done using an in-house developed python script which uses API request to determine their type from the PDB webpage.

## 3.3 Assigning Pharmacophore descriptors

A total of 9 pharmacophore descriptors were chosen on the basis of unique chemical features that play distinct role in binding to the DNA bases. The pharmacophores chosen were: aromatic, hydrophobic, hydrophilic, positively charged, negatively charged, and proline, cysteine, tryptophan and glycine were taken separately as they show unique characteristics.

Phenylalanine and tyrosine were grouped as aromatic; alanine, valine, leucine, isoleucine, and methionine were grouped as hydrophobic; asparagine, glutamine, serine, and threonine were grouped together as hydrophilic as these four amino acids have both hydrogen bond donor and acceptor groups in their side chain; histidine, arginine, and lysine were grouped as positively charged; aspartic acid and glutamic acid were grouped as negatively charged.

| AROMATIC {F, Y} | R | PROLINE {P} | P |
|---|---|---|---|
| HYDROPHOBIC {A, V, L, I, M} | H | CYSTEINE {C} | C |
| HYDROPHILIC {N, Q, S, T} | F | TRYPTOPHAN {W} | W |
| POSITIVELY CHARGED {H, R, K} | O | GLYCINE {G} | G |
| NEGATIVELY CHARGED {D, E} | X | | |

**Fig.2** Notations for the pharmacophores assigned

### 3.4 Building the pharmacophore matrix

To build the pharmacophore scoring matrix, we obtained the substitution frequencies for all the pairs of amino acids and also, the individual amino acid frequencies, which were used by Henikoff & Henikoff to create the BLOSUM62 scoring matrix, from the BLOCKS database.

The amino acids were grouped into their respective pharmacophore and then the substitution frequencies were recalculated for each pharmacophore using an in-house developed python script. The individual frequencies of the amino acids were also recalculated for each pharmacophore.

Then, using the observed frequencies and the expected frequencies for each pharmacophore obtained, the log odds score for each pair of pharmacophores was calculated using an in-house developed python script.
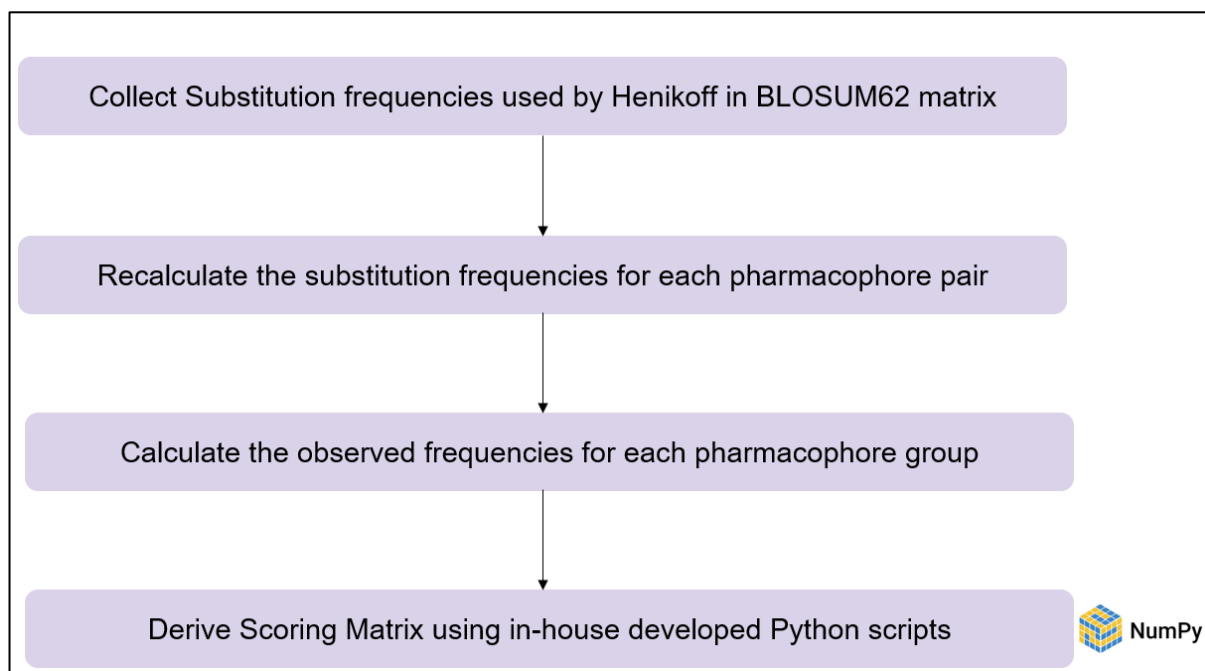


**Fig.3** Workflow for deriving the pharmacophore scoring matrix

Chapter 4

# Results and Discussion

## 4.1 Pharmacophore Matrix

A new pharmacophore scoring matrix was developed using in-house developed python script. The matrix contains the log-odds score for each pharmacophore pair mentioned in chapter 3.

```
AA     C      F      G      H      O      P      R      W      X
 C   9.00  -2.00  -3.00  -1.00  -3.00  -3.00  -2.00  -2.00  -4.00
 F  -2.00   2.00  -1.00  -1.00   0.00  -2.00  -2.00  -3.00   0.00
 G  -3.00  -1.00   5.00  -2.00  -2.00  -3.00  -3.00  -3.00  -2.00
 H  -1.00  -1.00  -2.00   2.00  -2.00  -2.00  -1.00  -2.00  -2.00
 O  -3.00   0.00  -2.00  -2.00   4.00  -2.00  -1.00  -3.00   0.00
 P  -3.00  -2.00  -3.00  -2.00  -2.00   7.00  -4.00  -4.00  -2.00
 R  -2.00  -2.00  -3.00  -1.00  -1.00  -4.00   5.00   1.00  -3.00
 W  -2.00  -3.00  -3.00  -2.00  -3.00  -4.00   1.00  10.00  -4.00
 X  -4.00   0.00  -2.00  -2.00   0.00  -2.00  -3.00  -4.00   4.00
```

**Fig.4** The pharmacophore scoring matrix derived using python script

Since Cysteine is unique, its high score suggests that it has distinct properties, likely due to its ability to form disulfide bonds. Negative scores with most other residues (e.g., C-X = -4.00, C-W = -2.00) indicate that Cysteine's interactions are highly selective.

Tryptophan (W) has the highest self-score (10.00). This suggests strong self-affinity, likely due to its bulky aromatic nature and involvement in π-π interactions. Negative interactions with negatively charged residues (X = -4.00) suggest that Tryptophan is not favourable in negatively charged environments.

Glycine is the smallest amino acid, providing structural flexibility. The matrix shows it has a high self-score (5.00) but does not interact well with other groups, likely because it lacks a side chain and contributes little to binding specificity. Thus, it has negative scores with all other groups.

Proline is a structural disruptor due to its rigid cyclic structure, which often introduces kinks in proteins. It has a high self-score of 7.00 but strong negative scores with positively charged (O = -4.00) and aromatic (R = -3.00) residues, which suggests that Proline disrupts stable charged/aromatic interactions.

Aromatic residues (R) have mixed interactions. The self-score is moderate, implying that π-π stacking is possible but not as strong as for Tryptophan (W). R-X score (-3.00) suggests that aromatic residues do not interact well with negatively charged residues.

Most interactions involving hydrophobic residues (H) are negative, meaning that hydrophobic residues do not significantly contribute to specificity. This aligns with their role in protein folding rather than direct molecular recognition.

There are several problems in this matrix that need to be fixed. There are some zero values such as for X-O, X-F, and O-F. This happened due to grouping of the amino acids which affected their observed frequency values used to calculate the log odds score.

## 4.2 Classification of the sequences

After classifying the HTH protein sequences, a total of 19 heterodimers and 263 homodimer sequences were found. The classification of the sequences into groove-binding and phosphate backbone-binding is still in progress.

## 4.2 Aligning HTH protein sequences

Out of all the helix-turn helix sequences obtained from NAKB, 34 closely related HTH motif sequences were taken and aligned using ClustalW offline tool. First the sequences were aligned using the BLOSUM62 matrix and then, the sequences were converted to pharmacophore sequences and realigned using the new pharmacophore-based matrix. The differences in the alignments were observed. It was observed that due to the inappropriate scores between O-X, X-F, and O-F, the alignment had unnatural gaps in case of the pharmacophore-based matrix, whereas those parts were well-aligned in case of the BLOSUM62 matrix.
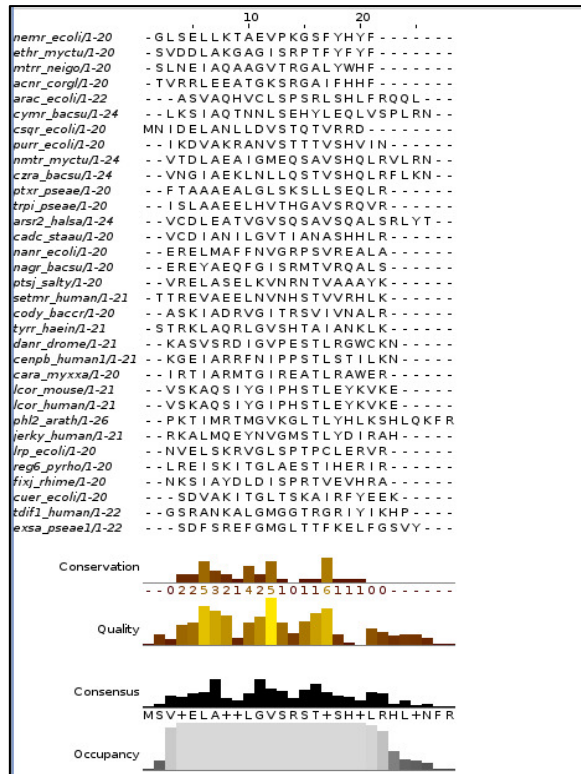
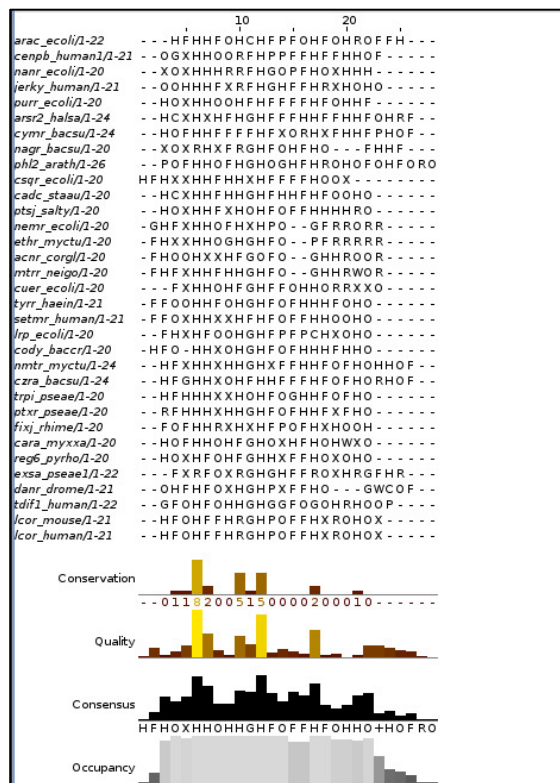**Fig.5** Multiple sequence alignment using BLOSUM62 matrix



**Fig.6** Multiple sequence alignment using pharmacophore-based matrix

<div align="right">

Chapter 5

# Summary and Conclusion

</div>

---

## 5.1 Summary

The pharmacophore scoring matrix developed in this study gives a detailed perspective on the unique chemical characteristics of amino acids within the context of helix-turn-helix (HTH) protein-DNA interactions. By classifying amino acids into pharmacophore categories—such as aromatic, hydrophobic, hydrophilic, positively charged, negatively charged, and distinct individual residues like proline, cysteine, tryptophan, and glycine, the matrix illustrates the biochemical details that dictate binding to DNA. Notably, cysteine's high self-score reflects its unique chemical properties, particularly its ability to form disulfide bonds, which contributes to structural stability, while its negative scores with most other residues suggest a selective and highly context-dependent role in protein-DNA interactions. Similarly, tryptophan exhibits the highest self-score in the matrix, indicating a strong self-affinity, likely driven by its bulky aromatic side chain and propensity for π-π stacking interactions. This high score suggests that tryptophan plays a significant role in stabilizing protein-DNA complexes, while its negative interaction with negatively charged residues implies an aversion to electrostatically unfavourable environments. Glycine, the smallest and most flexible amino acid, displays a high self-score, underscoring its importance in maintaining the structural flexibility of the HTH motif's turn region. However, its lack of a side chain and minimal involvement in direct binding specificity are reflected in its generally weak interactions with other pharmacophore groups. The development of this matrix, using recalculated substitution frequencies and log-odds scores, offers a novel lens through which to evaluate the role of amino acid properties in DNA-binding specificity, providing a powerful tool for understanding and predicting protein-DNA interactions. By integrating sequence analysis and pharmacophore-based scoring, this research not only highlights the conserved yet variable nature of HTH recognition helices but also paves the way for identifying subtle sequence-structure-function relationships, which are crucial for deciphering transcriptional regulation and designing potential therapeutic interventions.

**5.2 Conclusion**

At this point in the study, we can conclude that the pharmacophore scoring matrix that has been derived still has several limitations and needs a lot of modifications. The scores in the matrix do not seem to be appropriate, so the process of calculating the log-odds scores needs to be inspected.

The probability values and the frequencies which are obtained after the grouping of the amino acids need to be normalized in a certain manner. This is important as the pharmacophore groups contain varying number of residues which results in inappropriate probability values as well as frequency values. The normalization method is to be worked upon in the remaining period and a proper pharmacophore scoring matrix is to be developed which can be used to determine the pattern of DNA binding in the HTH proteins.

# References

1) Brennan, Richard G. "DNA recognition by the helix-turn-helix motif: Current Opinion in Structural Biology 1992, 2: 100…-108." *Current Opinion in Structural Biology* 2.1 (1992): 100-108.

2) Pellegrini-Calace, Marialuisa, and Janet M. Thornton. "Detecting DNA-binding helix–turn–helix structural motifs using sequence and structure information." *Nucleic acids research* 33.7 (2005): 2129-2140.

3) Steitz, T. A., et al. "Structural similarity in the DNA-binding domains of catabolite gene activator and cro repressor proteins." *Proceedings of the National Academy of Sciences* 79.10 (1982): 3097-3100.

4) Anderson, W. F., et al. "Structure of the cro repressor from bacteriophage λ and its interaction with DNA." *Nature* 290.5809 (1981): 754-758.

5) McLaughlin, William A., and Helen M. Berman. "Statistical models for discerning protein structures containing the DNA-binding helix-turn-helix motif." *Journal of molecular biology* 330.1 (2003): 43-55.

6) Rosinski, James A., and William R. Atchley. "Molecular evolution of helix–turn–helix proteins." *Journal of molecular evolution* 49 (1999): 301-309.

7) Harrison, Stephen C. "A structural taxonomy of DNA-binding domains." *Nature* 353.6346 (1991).

8) Steitz, Thomas A. "Structural studies of protein–nucleic acid interaction: the sources of sequence-specific binding." *Quarterly reviews of biophysics* 23.3 (1990): 205-280.

9) Luscombe, Nicholas M., et al. "An overview of the structures of protein-DNA complexes." Genome biology 1 (2000): 1-37.

10) Steitz, Thomas A. "Structural studies of protein–nucleic acid interaction: the sources of sequence-specific binding." Quarterly reviews of biophysics 23.3 (1990): 205-280.

11) Pabo, Carl O., and Robert T. Sauer. "Transcription factors: structural families and principles of DNA recognition." *Annual review of biochemistry* 61.1 (1992): 1053-1095.

12) Harrison, Stephen C., and Aneel K. Aggarwal. "DNA recognition by proteins with the helix-turn-helix motif." Annual review of biochemistry 59.1 (1990): 933-969.

13) Güner, Osman F., ed. Pharmacophore perception, development, and use in drug design. Vol. 2. Internat'l University Line, 2000.

14) Anantharaman, Vivek, Dapeng Zhang, and L. Aravind. "OST-HTH: a novel predicted RNA-binding domain." Biology direct 5 (2010): 1-8.