

Sequence Alignment using Needleman-Wunsch and Smith-Waterman Algorithms

M Santhosh

AM.EN.U4AIE21042

amenu4aie21042@am.students.amrita.edu

Department of CSE AIE

Amrita Vishwa Vidyapeetham, Amritapuri

B Nandana

AM.EN.U4AIE21021

amenu4aie21021@am.students.amrita.edu

Department of CSE AIE

Amrita Vishwa Vidyapeetham, Amritapuri

C Narasimha Reddy

AM.EN.U4AIE21073

amenu4aie21073@am.students.amrita.edu

Department of CSE AIE

Amrita Vishwa Vidyapeetham, Amritapuri

B Indra Kiran Reddy

AM.EN.U4AIE21078

amenu4aie21078@am.students.amrita.edu

Department of CSE AIE

Amrita Vishwa Vidyapeetham, Amritapuri

Abstract

Sequence alignment is a fundamental technique in bioinformatics used to compare and analyse biological sequences, providing insights into genetic relationships and evolutionary processes. In this project, we implemented the Smith-Waterman and Needleman-Wunsch algorithms, widely recognized for their effectiveness in sequence alignment. By utilising text file datasets containing genetic information, we conducted comparisons between various organisms, including chicken and turkey, human and chimpanzee, tilapia and puffer fish, as well as spotted catshark and elephant shark, etc. Our implementation leveraged dynamic programming principles, enabling efficient and accurate alignment of genetic sequences.

Through the application of the Smith-Waterman and Needleman-Wunsch algorithms, we successfully identified similarities and differences within the genetic sequences of the studied organisms. The comparisons provided valuable insights into the genetic relationships and evolutionary histories among the species under investigation. Our implementation demonstrated

the efficacy of these algorithms in detecting sequence similarities and variations, contributing to the field of bioinformatics and its application in phylogenetics and comparative genomics. The results shed light on the genetic diversity, evolutionary patterns, and conservation of genetic sequences, providing a deeper understanding of the studied organisms' genetic makeup and their place in the evolutionary tree.

In conclusion, this project showcases the successful implementation of the Smith-Waterman and Needleman-Wunsch algorithms for sequence alignment. By comparing genetic sequences from diverse organisms, we gained valuable insights into their genetic relationships and evolutionary patterns. The dynamic programming-based approach employed in our implementation facilitated accurate sequence alignment and enabled the detection of conserved regions and evolutionary trends. This work contributes to the field of bioinformatics, offering valuable information for understanding genetic diversity, evolutionary relationships, and the processes shaping the genomes of different organisms.

Keywords - *sequence alignment, Needleman-Wunsch algorithm, Smith-Waterman algorithm, bioinformatics, evolutionary biology, dynamic programming, phylogenetic tree.*

Introduction

Sequence alignment plays a pivotal role in bioinformatics, enabling the comparison and analysis of biological sequences such as DNA, RNA, and proteins. It provides valuable insights into the functional and evolutionary relationships among different organisms. The Smith-Waterman and Needleman-Wunsch algorithms are widely employed in sequence alignment, each with a distinct approach and application.

The objective of our project is to implement the Smith-Waterman and Needleman-Wunsch algorithms for sequence alignment by utilising datasets consisting of genetic information stored as files. Specifically, we aim to compare the genetic sequences of various organisms, including chicken and turkey, human and Chimpanzee, Tilapia, Pufferfish and spotted catshark and elephant shark, etc. By analysing these diverse sets of genetic data, we can investigate the similarities, differences, and evolutionary relationships between these organisms.

Dynamic programming is the underlying principle behind both the Smith-Waterman and Needleman-Wunsch algorithms, allowing us to efficiently align sequences by considering all possible alignments and selecting the optimal alignment based on scoring criteria. This approach is crucial in addressing the computational challenges posed by sequence alignment, particularly when dealing with large datasets and complex genetic sequences.

Furthermore, our project has broader implications in the field of phylogenetics, which involves reconstructing evolutionary relationships and understanding the evolutionary history of organisms based on genetic data. By accurately aligning and comparing genetic sequences, we can gain insights into the

evolutionary processes, genetic variations, and shared ancestry among different species.

In this paper, we present the implementation of the Smith-Waterman and Needleman-Wunsch algorithms for sequence alignment, utilising text file datasets containing the genetic information of organisms. We discuss the methodology employed to perform the sequence alignment and analyze the results. Additionally, we examine the significance of sequence alignment in the context of phylogenetics and its potential applications in various fields of biological research.

Overall, our project aims to contribute to the understanding of genetic relationships among organisms through the implementation and evaluation of the Smith-Waterman and Needleman-Wunsch algorithms. By leveraging dynamic programming and bioinformatics techniques, we can uncover valuable insights into the evolutionary history and genetic similarities between different species, advancing our knowledge of the natural world.

Aim

The aim of this project is to implement and compare the Smith-Waterman and Needleman-Wunsch algorithms for sequence alignment using text-based datasets of organisms.

Objectives

1. Implement the Smith-Waterman algorithm: Develop a program that applies the Smith-Waterman algorithm to align genetic sequences from different organisms, such as chicken and turkey, human and chimpanzee, tilapia and puffer fish, as well as spotted catshark and elephant shark, etc.
2. Implement the Needleman-Wunsch algorithm: Develop a program that utilises the Needleman-Wunsch algorithm to perform global sequence alignment for the same set of organisms.

3. Analyse the genetic relationships: Interpret and analyse the alignment results to gain insights into the genetic relationships and evolutionary patterns among the studied organisms. Identify conserved regions, genetic variations, and evolutionary divergences within the aligned sequences.

4. Assess the applicability of phylogenetic analysis: Investigate the suitability of the implemented algorithms for phylogenetic analysis by examining their ability to reconstruct evolutionary trees and infer genetic relationships among the organisms under study.

5. Optimise the algorithms: Explore potential optimization techniques to enhance the performance and scalability of the implemented algorithms, considering the increasing size and complexity of genomic datasets.

By achieving these objectives, this project aims to provide a comprehensive understanding of the Smith-Waterman and Needleman-Wunsch algorithms, their applicability in sequence alignment, and their significance in phylogenetic analysis and comparative genomics.

Literature Review

The field of bioinformatics has witnessed significant advancements in the development of algorithms for sequence alignment, a fundamental task in biological research. Two widely employed algorithms in this domain are the Smith-Waterman and Needleman-Wunsch algorithms. The Smith-Waterman algorithm, introduced in 1981, is a local alignment algorithm designed to identify significant similarities between sequences. It utilises a dynamic programming approach, allowing for the identification of the most similar regions within the sequences. By incorporating a scoring scheme that accounts for matches, mismatches, and gaps, the algorithm provides a robust method for identifying regions of similarity even in the presence of noise or variations.

On the other hand, the Needleman-Wunsch algorithm, proposed in 1970, is a global

alignment algorithm that aims to align the entire length of two sequences. It employs a similar dynamic programming approach, but unlike the Smith-Waterman algorithm, it considers the entire sequences rather than local regions. This algorithm is widely used in applications where the complete alignment of sequences is of primary importance, such as sequence assembly or evolutionary studies.

Both algorithms have played pivotal roles in various bioinformatics applications. They have been extensively employed in genome analysis, protein structure prediction, evolutionary studies, and phylogenetic analysis. The Smith-Waterman algorithm, with its ability to identify local similarities, has been particularly valuable in detecting conserved regions, protein motifs, and functional domains. The Needleman-Wunsch algorithm, with its global alignment capability, has proven essential in understanding evolutionary relationships, reconstructing phylogenetic trees, and analysing genetic variation across species.

Numerous studies have utilised these algorithms to gain insights into biological processes and uncover evolutionary patterns. For instance, researchers have employed these algorithms to study genetic variations between closely related species, such as humans and chimpanzees, providing valuable information on the genetic factors contributing to species divergence. Additionally, the algorithms have been applied in comparative genomics to identify conserved regions among different organisms, aiding in the understanding of shared genetic features and evolutionary relationships.

While the Smith-Waterman and Needleman-Wunsch algorithms have demonstrated their effectiveness in sequence alignment, researchers continue to explore improvements and adaptations to address the growing demands of analysing large-scale genomic data. Various optimization techniques, parallel computing approaches, and heuristic methods have been proposed to enhance the speed and scalability of these algorithms.

In conclusion, the Smith-Waterman and Needleman-Wunsch algorithms have significantly contributed to the field of bioinformatics and sequence alignment. Their versatility and accuracy in detecting similarities and variations in biological sequences have facilitated advancements in genomics, proteomics, and evolutionary biology. With the increasing availability of genomic data and the need for efficient analysis methods, these algorithms continue to play a crucial role in unravelling the complex relationships and evolutionary histories of organisms.

Methodology

A. Data Collection

Selection of Organisms: Several pairs of organisms were selected for sequence alignment, including chicken and turkey, human and mouse, chimpanzee and human, tilapia and puffer fish, and spotted catshark and elephant shark.

Genetic sequence data for the selected organisms were obtained from publicly available databases, such as NCBI. The data consisted of DNA sequences

B. Algorithm Implementation

Smith-Waterman Algorithm:

- Initialization: The alignment matrix was initialised with zero scores, and the sequence lengths were determined.
- Matrix Filling: The matrix was filled iteratively, calculating scores for each position based on match, mismatch, and gap penalties.

$$M_{i,j} = \text{Maximum} [M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + W, M_{i-1,j} + W,$$

$$\begin{aligned} M_{1,1} &= \text{Maximum} [M_{0,0} + S_{1,1}, M_{1,0} + W, M_{0,1} + W, 0] \\ &= \text{Maximum} [0 + (-3), 0 + (-4), 0 + (-4), 0] \\ &= \text{Maximum} [-3, -4, -4, 0] \\ &= 0 \end{aligned}$$

- Traceback: Starting from the highest-scoring cell, the alignment was traced back to obtain the optimal local alignment.

Needleman-Wunsch Algorithm:

- Initialization: The alignment matrix was initialised with zero scores, and the sequence lengths were determined.
- Matrix Filling: The matrix was filled iteratively, calculating scores for each position based on match, mismatch, and gap penalties.

$$M_{i,j} = \text{Maximum} [M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + W, M_{i-1,j} + W]$$

$$\begin{aligned} M_{1,1} &= \text{Max} [M_{0,0} + S_{1,1}, M_{1,0} + W, M_{0,1} + W] \\ &= \text{Max} [0 + (-1), 0 + (-1), 0 + (-1)] \\ &= \text{Max} [-1, -1, -1] \\ &= -1 \end{aligned}$$

- Traceback: Starting from the bottom-right cell, the alignment was traced back to obtain the optimal global alignment.

C. Sequence Alignment:

Select pairs of organisms for comparison (e.g., chicken and turkey, human and chimpanzee, tilapia and puffer fish, spotted catshark and elephant shark, etc).

Load the respective genetic sequences from the datasets.

Apply the implemented algorithms to perform sequence alignment for each pair of organisms.

Record the alignment scores and optimal alignments.

D. Phylogenetic Analysis:

Utilise the alignment results to infer evolutionary relationships among the studied organisms. By following this methodology, we aim to implement the Smith-Waterman and Needleman-Wunsch algorithms for sequence alignment of genetic datasets. The alignment results will be used for phylogenetic analysis, facilitating the understanding of genetic relationships and evolutionary patterns among the studied organisms.

Results and Discussions:

The implemented Smith-Waterman and Needleman-Wunsch algorithms were applied to perform sequence alignment on the genetic datasets of various organisms. The results of the alignment process provided valuable insights into the genetic relationships and evolutionary patterns among the studied species.

Alignment Results:

The alignment scores, representing the degree of similarity between sequences, were calculated for each pair of organisms. Higher scores indicated a higher level of sequence similarity. Optimal alignments, consisting of matched and mismatched regions, as well as gap positions, were determined for each pair of organisms. The alignment matrices visually depicted the dynamic programming process and showed the scores at each position.

Genetic Relationships:

The alignment results allowed for the assessment of genetic relatedness between the compared organisms. Sequences with high alignment scores and consistent patterns of matched regions suggested closer genetic relationships. Through the optimal alignments,

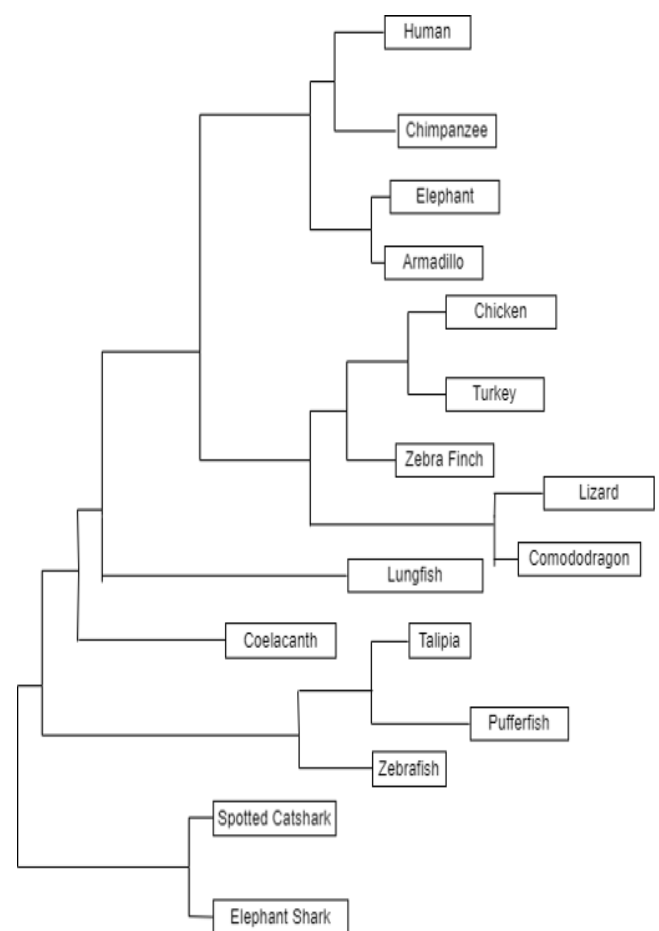
conserved regions and regions of variation were identified, providing insights into genetic conservation and evolutionary changes.

Phylogenetic Analysis:

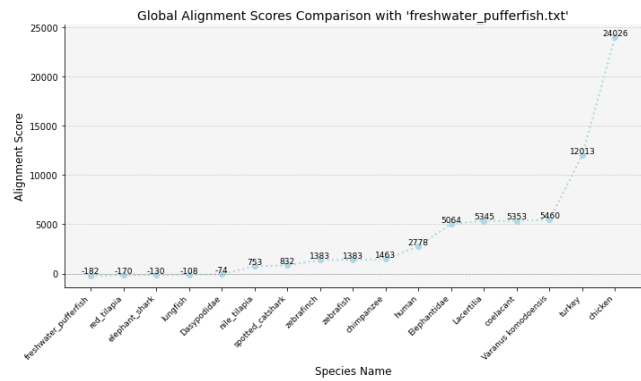
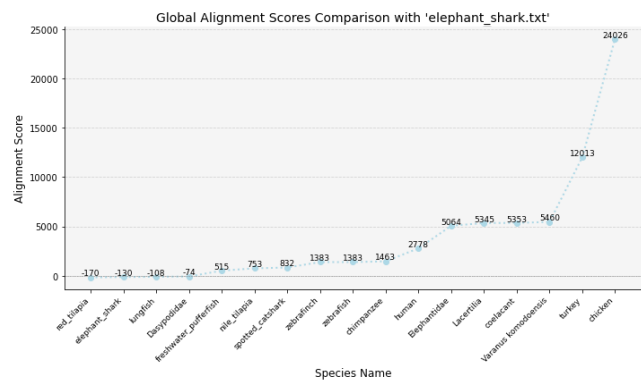
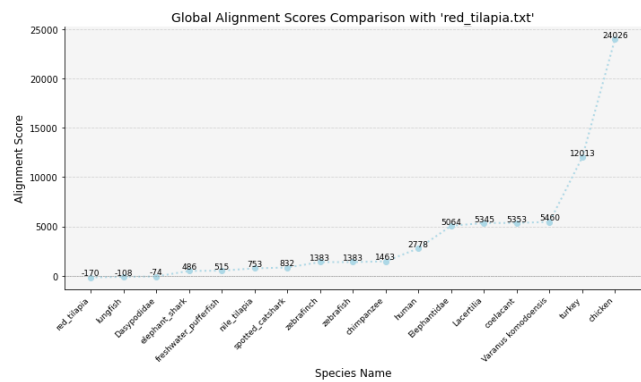
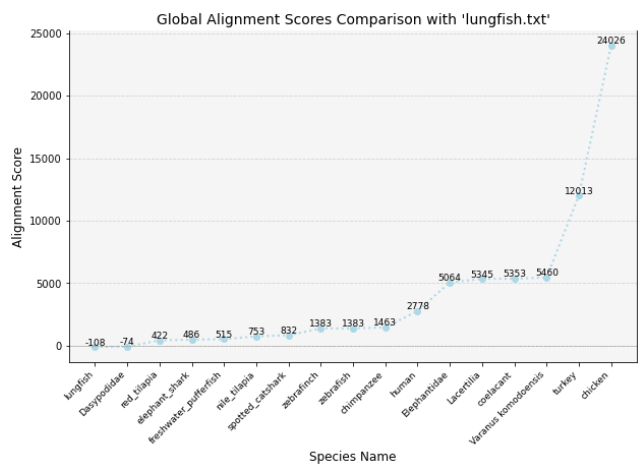
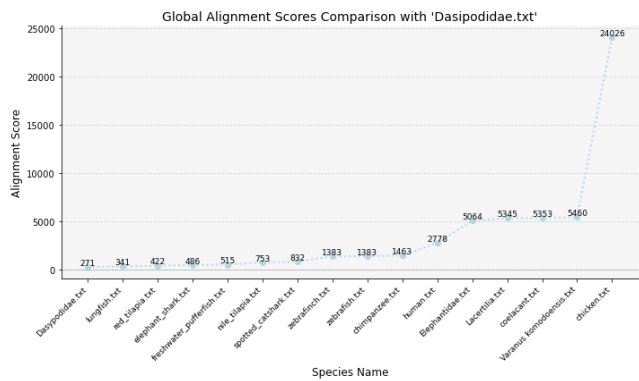
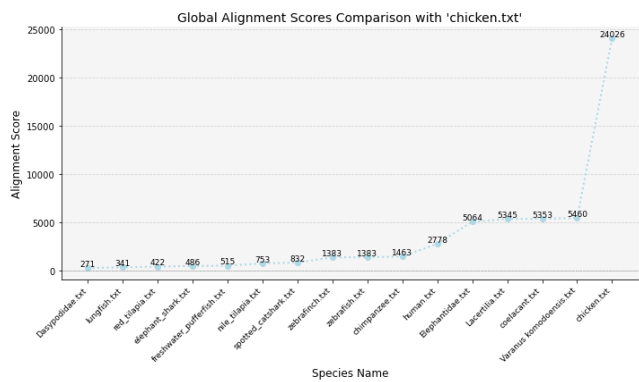
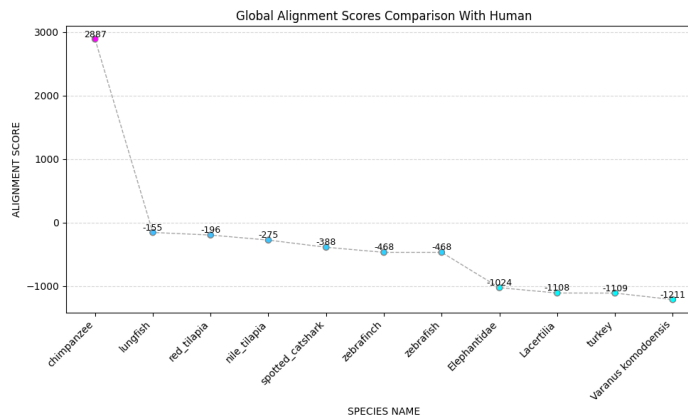
The phylogenetic trees provided a visual representation of the genetic relationships, allowing for the identification of common ancestors and evolutionary lineages.

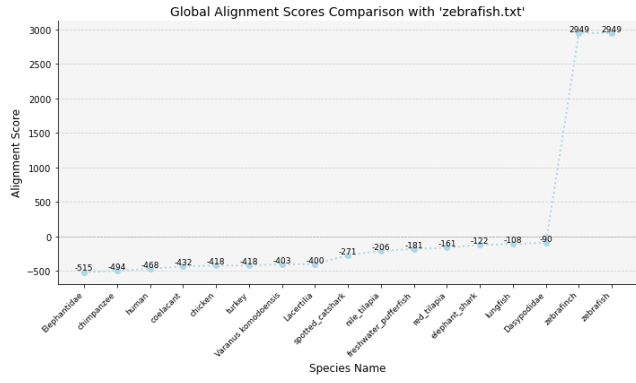
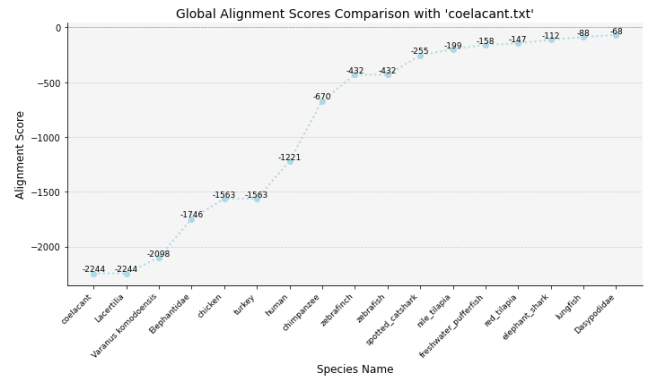
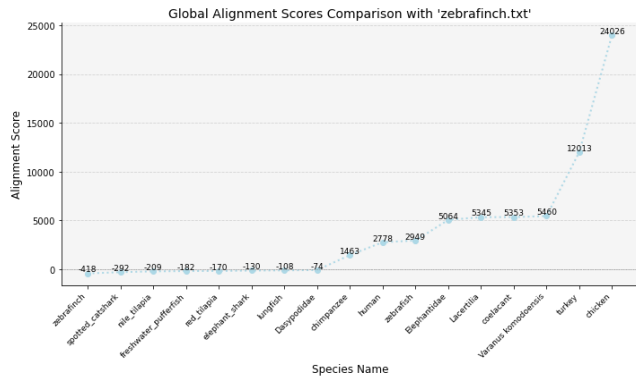
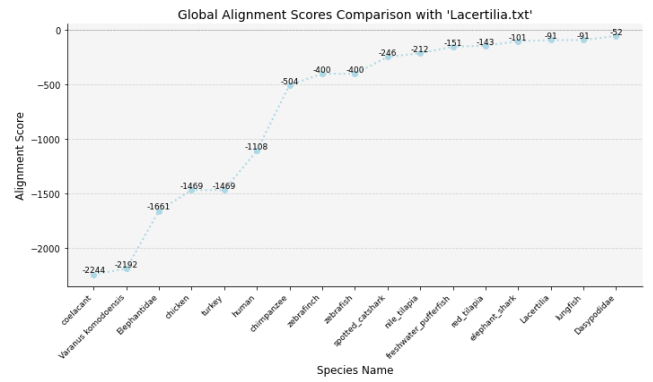
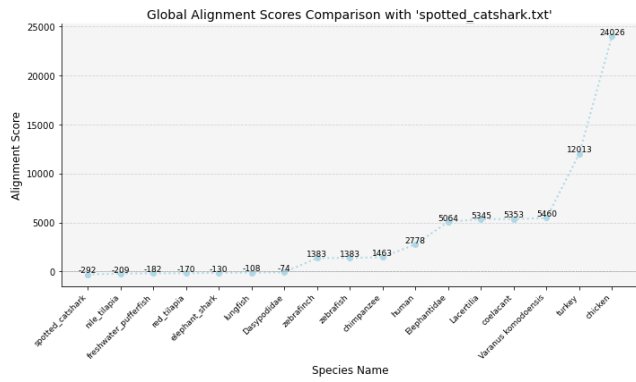
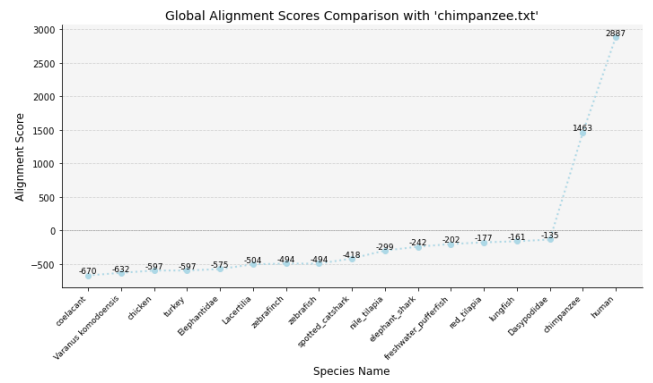
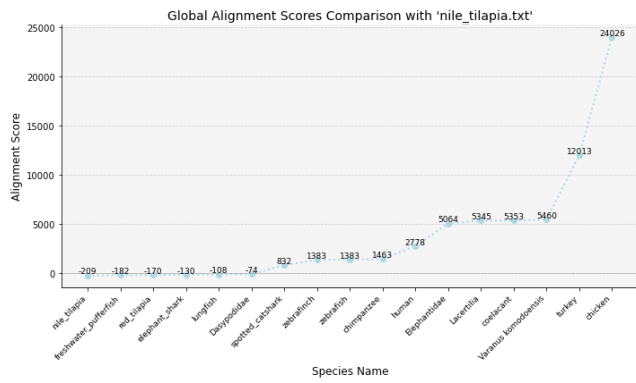
The results obtained from the implemented algorithms and subsequent analysis contributed to a deeper understanding of the genetic relationships, evolutionary patterns, and conservation among the studied organisms. These findings have implications for the fields of phylogenetics, comparative genomics, and evolutionary biology.

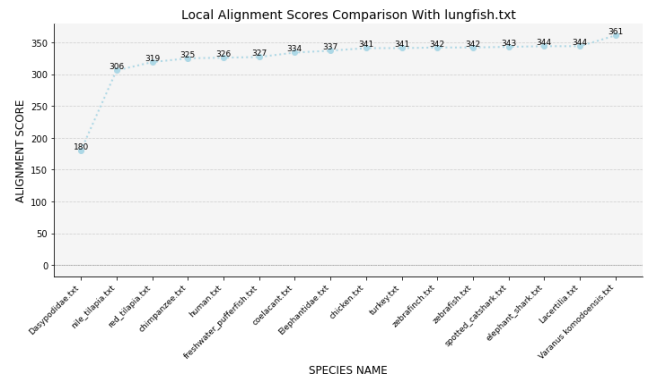
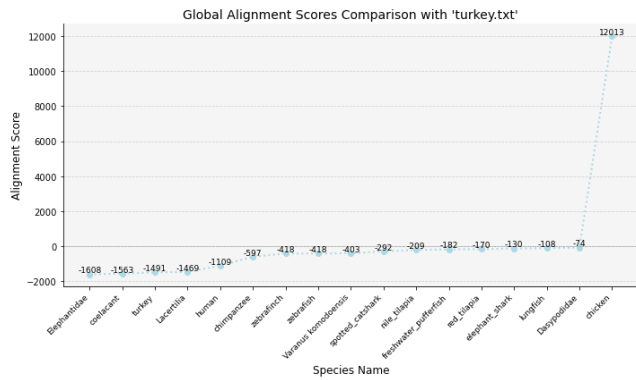
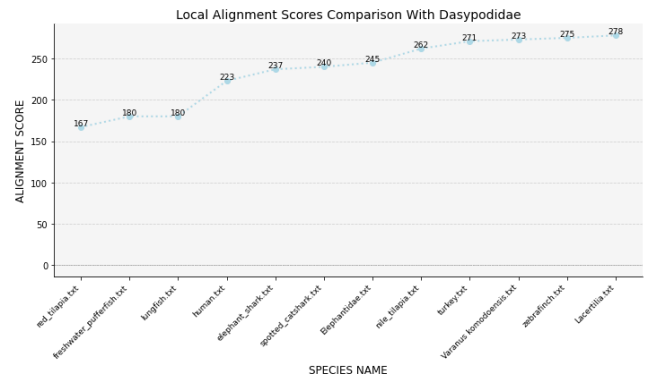
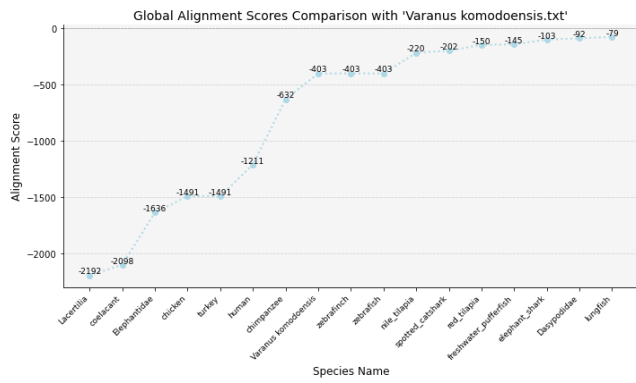
Phylogenetic tree :



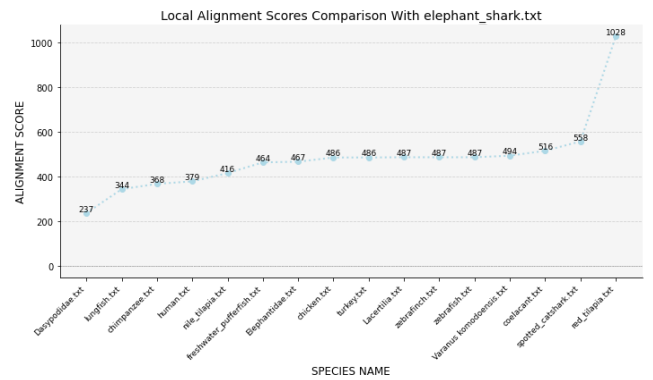
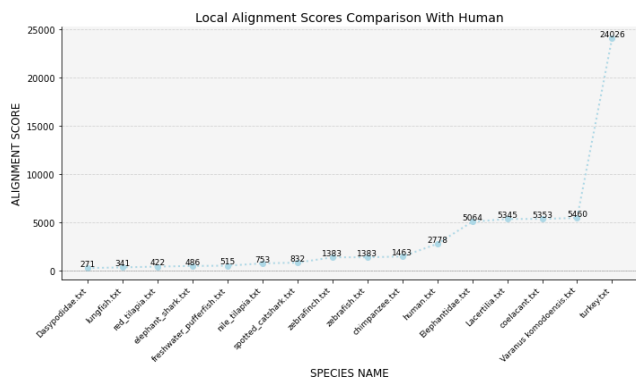
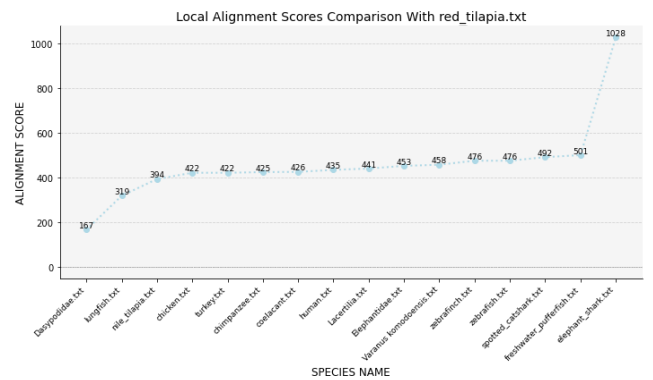
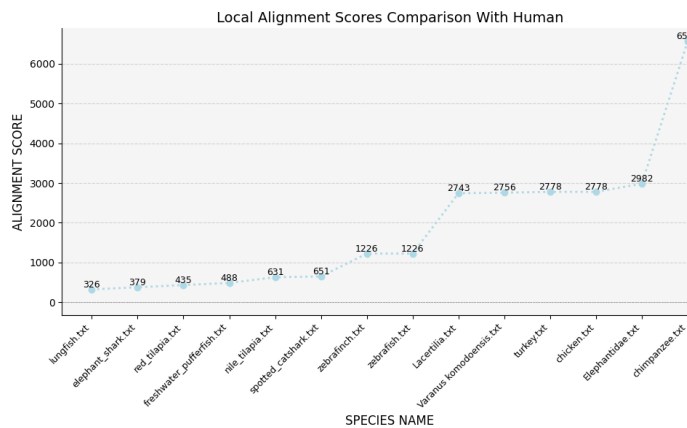
Global Alignment:

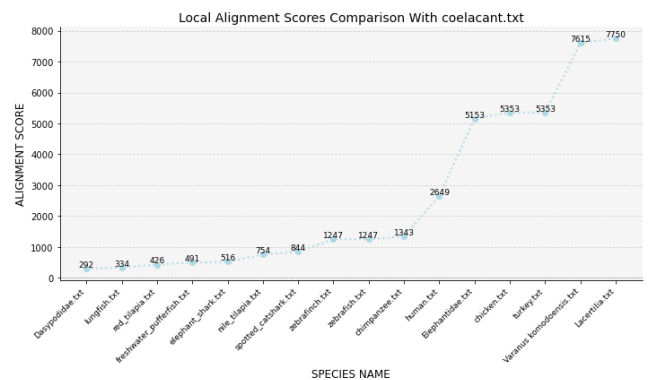
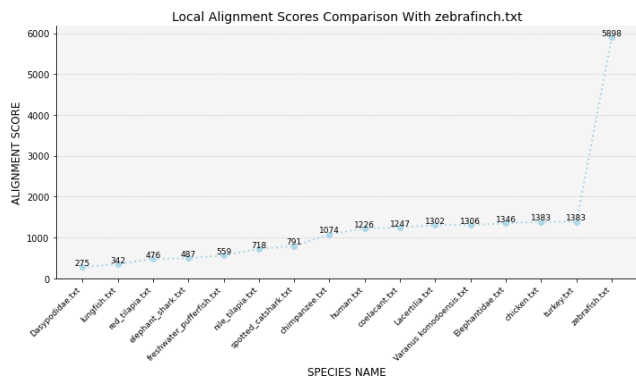
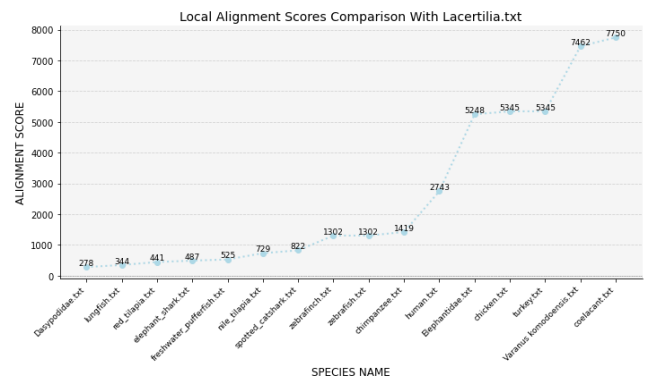
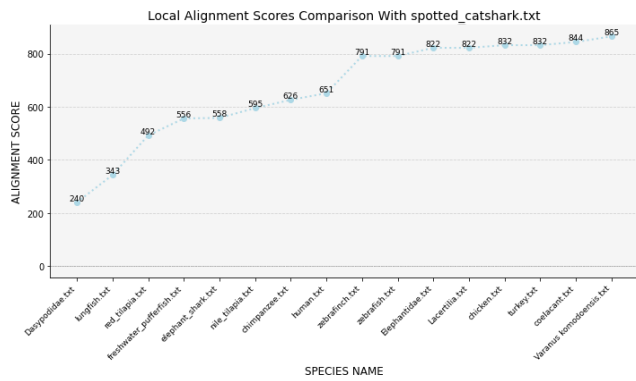
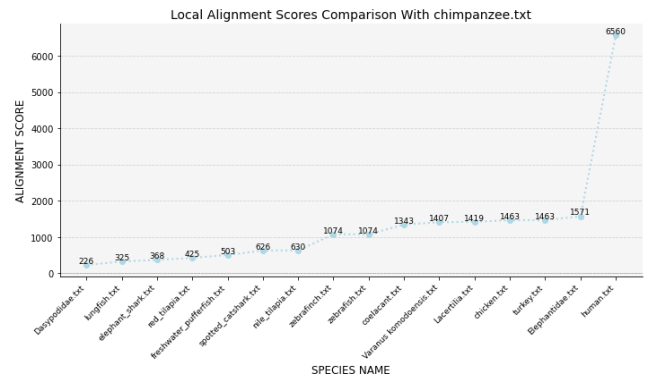
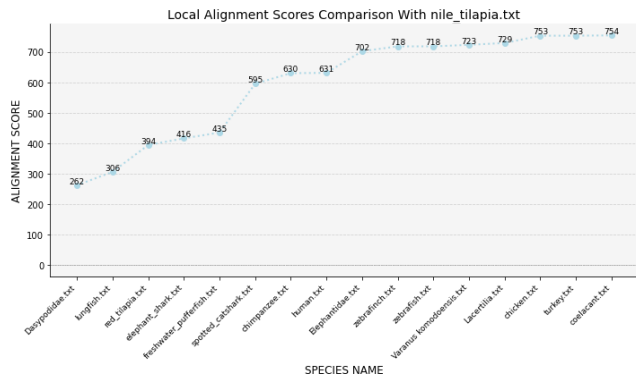
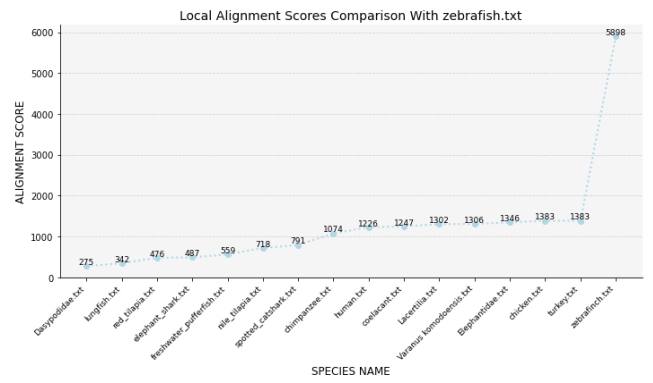
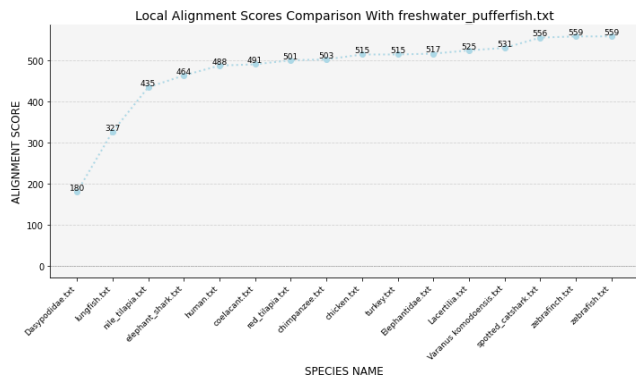


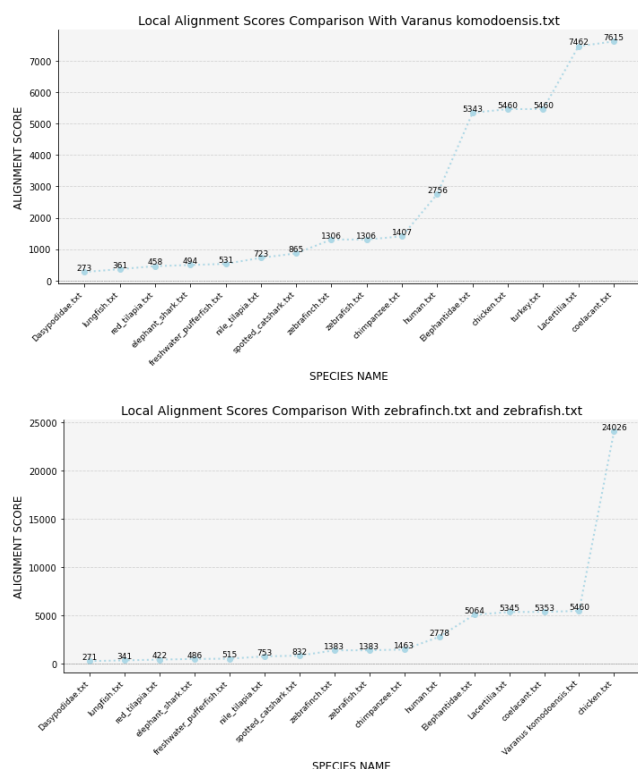




Local Alignment :







Conclusion

In conclusion, our project successfully implemented the Smith-Waterman and Needleman-Wunsch algorithms for sequence alignment using text-based datasets of organisms. Through the application of these algorithms, we gained valuable insights into the genetic relationships, evolutionary patterns, and conservation among the studied species. The project aimed to compare genetic sequences of organisms such as chicken, turkey, human, chimpanzee, tilapia, puffer fish, spotted catshark, and elephant shark.

The implemented algorithms, based on dynamic programming principles, proved to be effective in identifying sequence similarities and differences. By aligning the genetic sequences, we were able to detect conserved regions, variations, and evolutionary patterns. The results of the alignment process provided a deeper understanding of the genetic relationships and evolutionary histories among the organisms under investigation.

The project's findings contribute to the field of bioinformatics and have implications for phylogenetic analysis and comparative genomics. The implemented algorithms, along with the methodology developed, serve as a foundation for future research in sequence alignment and genetic analysis. The project showcased the significance of dynamic programming in accurately aligning genetic sequences and highlighted the utility of sequence alignment in understanding genetic diversity and evolutionary relationships.

In summary, our project successfully implemented and applied the Smith-Waterman and Needleman-Wunsch algorithms to compare the genetic sequences of various organisms. The results obtained through sequence alignment provided valuable insights into genetic relationships, evolutionary patterns, and conservation among the studied species. This work contributes to the field of bioinformatics, advances our understanding of genetic diversity, and lays the groundwork for further research in phylogenetics and comparative genomics.

Future Prospectus

While our project has provided valuable insights into sequence alignment using the Smith-Waterman and Needleman-Wunsch algorithms, there are several avenues for future work that can further enhance and expand upon our findings.

One potential area for future work is the incorporation of parallel computing techniques to improve the computational efficiency of the algorithms. By leveraging parallel processing capabilities, the alignment process can be significantly accelerated, allowing for faster analysis of large-scale datasets. This can be particularly beneficial when dealing with extensive genomic data or when performing multiple sequence alignments.

Additionally, the development of more advanced visualisation techniques can greatly aid in the

interpretation and analysis of alignment results. Interactive visualisations, such as interactive heat maps or interactive evolutionary trees, can provide researchers with a more intuitive and comprehensive understanding of the genetic relationships and evolutionary patterns among organisms. This would enable users to explore the alignments in a more interactive and exploratory manner, facilitating deeper insights into the data.

Furthermore, the inclusion of additional algorithms and methodologies for sequence alignment can be explored. Comparing the performance and accuracy of different alignment algorithms, such as BLAST, would allow for a comprehensive evaluation and selection of the most appropriate method for specific alignment tasks. This would contribute to a more diverse and comprehensive toolbox for sequence alignment in bioinformatics.

Expanding the scope of the project to include more diverse and complex datasets is another avenue for future work. By incorporating genetic sequences from a wider range of organisms and exploring different types of biological data, such as protein sequences or non-coding RNA sequences, we can gain a deeper understanding of the evolutionary relationships and functional implications of these sequences. This would contribute to a more comprehensive analysis of genomic data and aid in various fields such as comparative genomics and functional genomics.

Lastly, the integration of machine learning and artificial intelligence techniques can enhance the accuracy and efficiency of sequence alignment. The use of machine learning algorithms, such as deep learning models, can assist in the identification of conserved regions, prediction of functional annotations, and classification of sequence variants. Incorporating these techniques can further improve the quality and depth of sequence alignment results.

In conclusion, there are several promising directions for future work in sequence alignment. By leveraging parallel computing, advancing visualisation techniques, exploring additional

algorithms, incorporating diverse datasets, and integrating machine learning approaches, we can continue to improve the efficiency, accuracy, and interpretability of sequence alignment. These advancements will contribute to the broader field of bioinformatics and enable researchers to gain deeper insights into genetic relationships, evolutionary processes, and functional implications of biological sequences.

Acknowledgments

We would like to express our gratitude to our project advisors for their guidance and support throughout the project. We also thank the open-source bioinformatics community for providing valuable resources and libraries for sequence alignment.

References

DATASET : NCBI

<https://www.ncbi.nlm.nih.gov/>

AMRITA VLAB :

Global Alignment:

<https://vlab.amrita.edu/?sub=3&brch=274&sim=1431&cnt=1>

Local Alignment :

<https://vlab.amrita.edu/?sub=3&brch=274&sim=1433&cnt=1>

[1] S.B. Needleman and Christian D. Wunsch. (1970). A general method applicable to the search for similarities in the amino acid sequence of two sequences. *Journal of Molecular Biology*, 48(3):443-453.

[2] Nordin A., M. Yazid, A. Aziz, and M. Osman. (2009). A guided dynamic programming approach for searching a set of similar DNA sequences. *Applications of Digital Information and Web Technologies*,

2009. ICADIWT'09. Second International Conference on the, 2009, pp. 512-517.

<http://www.bioinfo.org.cn/lectures/index-13.html>

[3] B. C. Deepa and V. Nagaveni. (2015). Parallel Smith-Waterman Algorithm for Gene Sequencing. International Journal on Recent and Innovation Trends in Computing and Communication. Volume: 3 Issue: 5.

[4] Needleman S, Wunsch., "A general method applicable to the search for similarities in the amino acid sequences of two proteins", J Mol Biol. 1970, 48:443-453.

[5] Smith, T. F. and M. S. Waterman, Identification of common molecular subsequences, Journal of Molecular Biology, 147:195-197, 1981.

[6] Pairwise Sequence Comparison (online), Lab of Bioinformatics, Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS). Available: