

# Computational Biophysics: Algorithms to Applications (CS61060)

Spring 2022-23



## Term Project Report

**Topic: Fast multiple similar genetic sequence alignment based on the center star strategy**

### Group Members:

Rahul Mandal (20CS30039)

Rohit Ranjan (20CS30066)

Burra Nithish (20CS10018)

# Table of contents

- Introduction
- Problem statement
- Existing work/ Literature review
- Methodology
- Datasets
- Result
- Further Work & Conclusion

# Introduction

Sequence alignment is considered the 'Holy Grail' problem in computational biology and is of vital importance for molecular function prediction. The widely used databases PFAM (Robert et al., 2010) and RFAM (Gardner et al., 2011) are constructed based on multiple sequence alignment (MSA). Molecular function prediction sometimes depends on evolutionary information (Liu et al., 2014). MSA is also required for evolutionary tree reconstruction. Most of the available phylogenetic tree construction software tools require previously aligned sequences as input. When addressing the evolutionary analysis of bacterial and viral genomes, large-scale similar DNA and RNA sequences often prevent these MSA tools from functioning (Wang et al., 2013). The evolution of viruses is rapid, and massive viral DNA sequences often appear in phylogenetic reconstructions. Therefore, it is necessary to improve the scalable capacity of MSA tools.

# Problem Statement

Q). Multiple sequence alignment (MSA) is important work, but bottlenecks arise in the massive MSA of homologous DNA (deoxyribonucleic acid), RNA or genome sequences. Try to implement trie trees to accelerate the center star MSA strategy. The expected time complexity will be decreased to linear time from square time. The algorithm for center star strategy is also discussed below.

## Algorithm 1. Improved Centre Star Algorithm Based on Trie Trees

**Input:**  $n$  DNA Sequences,  $S_1, S_2, \dots, S_n$

**Output:**  $n$  aligned DNA Sequences  $S'_1, S'_2, \dots, S'_n$

1. For each DNA Sequence,  $S_i$ ,
2. Partition  $S_i$  into  $k$  segments  $\{S_{i1}, S_{i2}, \dots, S_{ik}\}$  with equal lengths;
3. Construct trie tree  $T_i$  for the segments  
set  $S_i = \{S_{i1}, S_{i2}, \dots, S_{ik}\}$
4. for  $j$  from 1 to  $n$ ,  $j \neq i$
5. search  $T_i$  in  $S_j$ , and set  $m_{ij}$  as the segment appearance times; record all of the appearances in  $\mathcal{A}_{ij}$
6. end for
7. calculate  $m_i = \sum_{j=1, j \neq i}^n m_{ij}$
8. end For
9.  $m^* = \operatorname{argmax}_{i=1,2,\dots,n} m_i$ , set  $S_{m^*}$  as the centre star sequence
10. For each  $i$  from 1 to  $n$ ,  $i \neq m^*$
11. Partition  $S_i$  and  $S_{m^*}$  according  $\mathcal{A}_{im^*}$ , align the mismatched regions and obtain the pairwise alignment; record all of the positions of inserted spaces in  $\mathcal{P}_{im^*}$  and  $\mathcal{P}_{m^*i}$ .
12. end For
13. For  $i$  from 1 to  $n$ ,  $i \neq m^*$
14. sum  $\mathcal{P}_{m^*i}$  to  $\mathcal{P}_{m^*}$
15. end For
16. obtain the final result,  $S'_{m^*}$ , according to  $\mathcal{P}_{m^*}$
17. For  $i$  from 1 to  $n$ ,  $i \neq m^*$
18. compare  $\mathcal{P}_{m^*i}$  with  $\mathcal{P}_{m^*}$ , and update  $\mathcal{P}_{im^*}$ , then obtain the final result,  $S'_i$
19. end For

# Existing Work and Literature Review

## **Zou et.al. - HAlign: Fast multiple similar DNA/RNA sequence alignment based on the centre star strategy**

The authors develop two software tools to address the MSA problem. The first employed trie trees to accelerate the centre star MSA strategy. Our problem statement uses this proposed algorithm. The expected time complexity was decreased to linear time from square time. Another major contribution of the paper was incorporating parallelism using Hadoop platform to address large-scale data.

## **Aho Corasick algorithm for Pattern Searching using Trie Trees**

Proposed by Alfred Aho and Margaret Corasick in 1975, the Aho–Corasick algorithm is a string-searching algorithm that utilizes Trie Trees data structure. Let there be a set of strings with the total length  $m$  (sum of all lengths). The Aho-Corasick algorithm constructs a trie with some additional links, and then constructs a finite state machine (automaton) in  $O(m \cdot k)$  time, where  $k$  is the size of the used alphabet. This automaton contains failure links where the longest suffix of current state which happens to be a proper prefix of another word in Trie is utilized. Final time complexity is  $O(n + m + z)$ , where ' $n$ ' is the length of the text, ' $m$ ' is the length of keywords, and ' $z$ ' is the number of matches in text.

# Methodology and Code Flow

1. Sequences in the input are broken into segments of size `SEGMENT_LENGTH`; and respective Trie trees are constructed with these segments.
2. Each sequence is now searched using the Aho-Corasick algorithm in the trie trees of other sequences and occurrences are used to score the trie tree sequence.
3. We now choose the centre sequence as the sequence with the highest calculated score.
4. Aho-Corasick is again used to match all the sequences with the chosen centre sequence and finding the indices of matched segments. At this stage we have all sequences partially aligned with the centre sequence because we have found matching segments. Between each of these matching segments, the sequence is still unaligned.
5. Call the normal dynamic programming based global alignment solution for these unmatched sequences.
6. Finally insert blanks into centre sequences from each of the instances to globally align them. Output results.

## File structure:

- `msa_star.py` -> Implementation of DP based MSA using centre star strategy for benchmarking.
- `msa_star_trie.py` -> Implementation of the current problem statement as described above.
- `aho_corasick.py` -> Implementation of Trie tree and utility algorithms for pattern match
- `extseq.py` and `extdbn.py` -> Code for preparing dataset in seq and dbn format for input.

# Datasets

The datasets used were the RNA sequences of distinct species like

- Archaea
- Bacteria
- Mitochondria
- Prokaryotes
- Plastids

The sequences were stored in the .seq files for each species.

Results for some species are attached below.

# Results

Following are the results of sequence alignment

## 1. 30 Sequences at a time for Species Archaea with Normal MSA algorithm

Time taken to align sequences: 9.773279428482056

The aligned sequences are:

```
-U-GAUAC-GG-CGGCCAUAGC-GGA-G-GUGUC-C-CAU--CC-G-AUCCCAUUCGGAUCUCGGAA-AUUAAGCC-CU-CCA-GCGAU-UUC---U-U--A-AGUAC-----UGC-C-A--UAUGGUG--GG-A-A-C-AAG-AU-GA-C---GCUGCCGAUC---AC
-U-C-A-AUAG-C-GGCCACAGCAGU--GUGUC-A-C-A-CCC-GUUC-CCAUUCCGAACACGGAAGUUA-AG-AC-AC-CUCACG-U---G-G-AGUAC-GGUACU-G-A-GGUACGCGAGUCCU-C-GGGA-A-AU--C--AUCCU-C-GC-U-G-CUA-UUGUU-
-----AA-GG--CGGCCAUAGC-GGC-G-GGGUC-C-C-U-CCC-GUA-CCCAUCCCGAACACGGAAGA-UAAGCC-CG-CCU-GCG-UAU-U-G--G-U--G-AGUAC--UGGAG-UGGGAGAC-CCU-CUG--GGAG-AGC-U-G-AU--U-C---GCCGCCU-----U-
----GAA-GG--CGGCCAGAGC-GGU-A-GGGAA-A-C-A-CCC-GUA-CCCAUUCGGAACACGGAAGU-UAAGCC-UA-CCA-GCG-UAU-C---G-U--GAAGUAC--UGGAG-UGAGCGAU-CCU-CUG--GG-A-ACCAC-GAG--U-C---GCCGCCU-C---CC
--U-A-A-GG-C-GGCCAUAGCG-GC-G-GGGUU-C-C-U-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCCC-GC-CU-GCGUUA--U-GG--U--G--AGUACUGG-A-GUG-GGAGACCCU-CUG-GGA-G--AGC-U-G-AUUC--G--C-U-G-CU--U--UA-
-N-U----GG-C-GGUCAUAGC-GC-G-GGGAA-A-C-A-CCC-GUAC-UCAUUCCGAACACGGAAGUUA-AGCCC-GC-CU-GCGU-U-UC-GG----C-GAGUACUGG-A-GUG-GCGAGCCU-CUG-GGA-G--AGC-U-G-AUUC--G--C-C-GUC-A-C--U--
-UAGGUUU-GG-C-GGUCAUAGC-GAU-G-GGGUU-A-C-A-CCUGGU-C-UCGUUUCGAUCCGAGAAGUUA-AG-UC-U-UUUCGCGU-U-UU-G-U--U--UGUGUACUUAUGG-G-UUC-CG-G-UCU-AUGGGA-A--UUU-C-A-UU--U-A-GC-U-GCC-AGCUUUUU
-G-CCCA-CC--CGGCCAUAGU-GG--GAGGGA-A-C-A-CCC-GGA-CUCAUUCGGAACCCGGAAGU-UAAGCCUUC--CC--ACG-U-U--G--G---A-AGGGCAGUGGGG-UCCGAGAGGAC--CU---G-C-AGC-C--CU--UCC-AAGCCG-GGA-U--GGG
--U----GA-C-GGUCAUAGC-GC-G-GGGAA-A-C-A-CCCGU-C-UCAUUCCGAACCCGGAAGUUA-AGCCC-GC-CC-GCGU-A-CC-GUG-U--GC--GUACUGG--GAUGCCCGAG-CU-CCCGGG--A--AGCAC-G-U-UC-G--C-U-GUC-ACCA-UC-
----GCC-GA--CGGCCAUAGG-GGUCG-GGGAA-A-C-A-CCC-GGA-CUCAUUCGGAACCCGGAAGU-UAAGCC-CGACC--CCG-U-UCC-G-CG-U---GGUAC--U-GUGUUCGAGAGGGGA-C-G--GG-A-AGCGC-G-G--A--A---GCCGUCGG-C---A-
----U---GGCCCGGCCAUAGC-UGC-C-GGGUA-A-C-A-CCC-GGA-CUCGUUUCGAACCCGGAAGU-UAAGCCG-G-CC--GCG-U-U---GAAG-UUGCCAG----U-GAGUUCGGAAGGGCU-C-GCAGGCA-CUU-C-AAGC--U---G--G-GCC---C--G-
GG--UCA-GA--CGGCCAUAGC-AGC-G-GGGUU--C-A-CCCGU--CCCAUUCGGAACCCGGAAGU-UAAGCC-CG--CUCGCG-UAU-C-U--G-U--C-GCGUAC--UGUUA-UGCGCAAG-UGUAC-G--GGAA-AGC-A-G-AC--A-C---GCUGUUAACC-ACU-
-----GA-CC--CGGCCAUAGU-GGC-C-GGGCA-A-C-A-CCCGU--CUCAUUUCGAACCCGGAAGU-UAAGCC-GG-CC--ACG---UCA-G-AG---C-GGCAG--U-GAGGUCCGAGAGGCCU-C-----G-C-AGC-C-G-CU-CUGA---GCUG-GGAUC---GG
----CGA-CC--CGGCCAUAGU-GGC-C-GGGCA-A-C-A-CCCGU--CUCGUUUCGAACCCGGAAGU-UAAGCC-GG-CC--ACG---UCA-G--A-A-C--GGCC-GU-GAGGUCCGAGAGGCCU-C-----G-C-AGC-C--GU--U-CUGAGCUG-GGAUC---GG
-U-U--A--GG-C-GGCCACAGCG-GU-G-GGGUUC-C-U-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCCC-AC-CA-GCGU-U-CC-GG--G--G--AGUACUGG-A-GUG-GCGAGCCU-CUG-GGA-A--A-C-CG-G-GUUC--G--C-C-GCC-ACCA---
AG-U----GG-U-GGCCAUUUCG-GC-G-GGGUU-C-C-U-CCCC-GUAC-CCAUUCCGAACACGGAAGUUA-AGCCC-GC-CA-GCG--U-CC-GG----C-AGUACUGG-A-GUG-GCGAGCCU-CUG-GGA-A--AUC-C-G-GUUC--G--C-C-GCC-A-CN-N-
-U--U-A-GG--CGGCCACAGC-GGU-G-GGGUUC-C-U-CCC-GUA-CCCAUCCCGAACACGGAAGA-UAAGCC-CA-CCA-GCG-U-UCCAG--G---G-AGUAC--UGGAG-UGCGCGAG-CCU-CUG--GGAA-AUC-C-G-GU--U-C---GCCGCC-A-C---C-
-U--U-A-GG--CGGCCACAGC-GGU-G-GGGUUC-C-U-CCC-GUA-CCCAUCCCGAACACGGAAGA-UAAGCC-CA-CCA-GCG-U-UCC-G--G--G-AGUAC--UGGAG-UGCGCGAG-CCU-CUG--GGAA-A-CG-C-GU--U-C---GCCGCC-ACC--A-
-G-U--A--G-C-GGCCACAGCG-GU-G-GGGUU-C-C-U-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCCC-AC-CA-GCGU-U-CC-GG--G--G--AGUACUGG-A-GUG-GCGAGCCU-CUG-GGA-A--A-C-CG-G-GUUC--G--C-C-G-CUA-CN-N-
-U--UAA-GG--CGGCCAUAGC-GGU-G-GGGUU-A-C-U-CCC-GUA-CCCAUCCCGAACACGGAAGA-UAAGCC-CG-CCU-GCG-U-UCC-G--G--U--C-AGUAC--UGGAG-UGCGCGAG-CCU-CUG--GGAA-AUC-C-G-GU--U-C---GCCGCCU-C---U-
-U-U----GG-C-GACC AUAGC-G--C-GAGUG-ACC-U-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCUC-GC-CU-GCGU-U-UC-GG--U--C--AGUACUGG-A-UUG-GCGAGCCU-CUG-GGA-A--AUC-U-G-AUUC--G--C-C-GCC-A-C--C-
-----C-GG--CGGCCAGAGC-GGU-G-AGGUU-C-C-A-CCC-GUA-CCCAUCCCGAACACGGAAGU-UAAGCU-CG-CCU-GCG-U-UUCU-G--G-U--C-AGUAC--UGGAG-UGAGCGAU-CCU-CUG--GGAA-AUC-C-A-GU--U-C---GCCGCCU-----U-
AU--UAC-GG--CGGCCAGAGC-GGU-G-AGGUU-C-C-A-CCC-GUA-CCCAUCCCGAACACGGAAGU-UAAGCU-CG-CCU-GCG-U-UUCU-G--G-U--C-AGUAC--UGGAG-UGAGCGAU-CCU-CUG--GGAA-AUC-C-A-GU--U-C---GCCGCC-C-C---U-
-U-U-A-A-GG-C-GGCCAGAGCG-GU-G-AGGUU-C-C-A-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCUC-AC-CU-GCGU-U-UC-GG--U--C--AGUACUGG-A-GUG-AGCGAUCCU-CUG-GGA-A--AUC-C-A-GUUC--G--C-C-GCC-C-CU---
-U--UAA-GG--CGGCCACAGC-GGU-G-AGGCA-A-C-A-CCC-GUA-CCCAUCCCGAACACGGAAGU-UAAGCU-CG-CCA-GCG--U--G--A--UAACAAGUAC--UGGAG-UGUGCGAA-CCU-CU---GGAA-A-C-A-UUUA-CU-C---GCCGCCU-C---U-
-----A-GG--CGGCCAUAGC-GGC-A-GGGAA-A-C-A-CCC-GUA-CCCAUCCCGAACACGGAAGU-UAAGCC-UG-CCA-GCG-U-UGC-G--G-U--G-AGUAC--UGGGG-UGUGCGAA-CCC-CUG--GGAA-AGC-C-G-GU--U-C---GCCGCCU-C---CA
-U--U-A-GG--CGGCCACAGC-GGC-G-AGGUU-C-C-U-CCC-GUA-CCCAUCCCGAACACGGAAGUUA-AGCUC-CG-CCU-GCG-UAU-C-G--G---CAUUAUC--UGGAG-UGGGCGAC-CCU-CUG--GGAACGU--C-G-AU--U-C---GCCGCCCC-C---A-
----UAA-GG--CGGCCAUAGC-GGC-G-GGGCA-A-C-A-CCC-GUA-CCCAUCCCGAACACGGAAGU-UAAGCC-CG-CCU-GCG-U-UCC-G--G--G---CGAGUAC--UGGAG-UACGAGAG-UCU-CUG--GGAA-AAC-C-G-GU--U-C---GCCGCCU-C---A-
--G-A-A-GG-C-GGCCAGAGCG-GU-G-GGGAA-A-C-A-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCCC-AC-CA-GCGU-A-CC-G--U--G--AAGUACUGG-A-GUG-AGCGAUCCU-CUG-GG--A--ACCAC-GAG-UC--G--C-C-GCCU-C--C-
-G-U--A-GG-C-GGCCAGAGCG-GU-A-GGGAA-A-C-A-CCC-GUAC-CCAUUCCCGAACACGGAAGUUA-AGCCU-AC-CA-GCGUUA--C--G--U--GAAGUACUGG-A-GUG-AGCGAUCCU-CAG-GG--A--ACCAC-GAG--UC-G--C-C-GCCUA-C--C--
```



### 30 Sequences at a time for Species Archaea with MSA using trie tree algorithm

```
Time taken to align sequences: 0.2355635166168213

The aligned sequences are:

---UGA-U-ACGG--CGGCUAUGC-GGA-GGUG-UC-C--CAU--CCG-AUCCCAUCCGAUCUCGGAA-AUUUAG-C-C-C-UC-CAGCGAU-UUC-U-U--UAA--G-U-AC-----UGC-C---A-U--UAUGGUGG--AA-C-A-A-G-AUGA-C---GCUG-CCGA-U---CAC
---UCAU-A-G--CGGCCACAGCAGGU--GUG-U--CA-C-A-CCCGUU--CCCAUCCGGAACACGGAAGU-UAAGACAC-C-U-CA-CG--U--G-GA-UGACGG-U-AC--U--GAGG-UACGCG--A-GUCU-C-GGG-AAAUCAU-C---C-U-C---GCUG-CU-AUUG---UU
-----A-A-GG--CGGCCAUAGC-GGC-GGGG-UC-C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-GC-CUGCG-UAU-CG-U--UGA--G-U-AC--UG-GA-G-UGGAG--A-C-CU-CUGGG-AGA--G-CUG-AU-U-C---GCCG-UC-A---C---UU
---G---A-A-GG--CGGCCAGAGC-GGU-AGGG-AA-A--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-U-AC-CAGCG-UAU-CG-U--GAA--G-U-AC--UG-GA-G-UGGAG--A-C-CU-CUGGG-AAC-C-A-C-GAG--U-C---GCCG-CCUU--C---CC
---U---A-A-GG--CGGCCAUAGC-GGC-GGGGUU--C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-GC-CUGCG-UAU-CG-U--UGA--G-U-AC--UG-GA-G-UGGAG--A-C-CU-CUGGG-AGA--G-CUG-AU-U-C---GCCG-CC-U--U---UA
---N---U---GG--CGGCUAUGC-GGC-GGGG-AA-A--C-A-CCCGUA--CUCAUCCGGAACACGGAAGU-UAAG-C-C-C-GC-CUGCG-U-UUCG-G--CGA--G-U-AC--UG-GA-G-UGCGCG--A-G-CU-CUGGG-AGA--G-CUG-AU-U-C---GCCG-UC-A-----CU
UAGGU---U-U-GG--CGGCUAUGC-GAU-GGGGUU-CA-C---CUGGU--CUCGUUCCGAUCCCGGAAGU-UAAGU-U-U-U-C-GCG-U-UUGU--UGU--G-U-ACUAG-G--G-UUC-CG---G-UUC-AUGG-AU-U-U-C-A-U-U-U-A---GCUG-CC-AG-CUUUUU
GC-C---C-A-CC--CGGCCAUAGU-GG--GAGG--G-CAAC-A-CCCGA--CUCAUCCGGAACCCGGAAGU-UAAG-C-CUC--C-CA-CG-U-U--G-GAAG--G-C-AG--UG-GG-G-UCCGAG--AGG-AC-CU--G--CA--G-C-C-U-CCAA-GCCG-GG-AU-G---GG
-----U---GA--CGGCUAUGC-GGC-GGGGAA-A-C---CCGGU--CUCAUCCGGAACCCGGAAGU-UAAG-C-C-GC-CCCG-U-ACCGUG--UGC--G-U-AC--UG-G--GAUGCGG--A-G--CUCCCGG-AGG-C-A-C-G-GU-U-C---GCCG-UC-A--C---CAUC
---G---C-C-GA--CGGCCAUAGG-GGUGGGG-AA-A--C-A-CCCGA--CUCAUCCGGAACCCGGAAGU-UAAG-C-C-GC-C-CG-U-UCCGCG--UG--G-U-AC--UGU--U-UCCGAG--AGG-GCA-C-GGG-AGG-C-G-C-G-G-A-A---GCCG-UC-G--G---CA
-----U---GCCCGCCAUAGC-UGC-CGGG-UA-A--C-A-CCCGA--CUCGUUCCGAACCCGGAAGU-UAAG-C-C-G-GC-C-GCG-U-U-----GAA--GUUGCC--AGUGA-G-UUC-CGAAAGG-GCU-C--G--CA-G-G-C-A-CU-U-CAA-GCUG-GG-G--C---CG
GG-U---G-A-GA--CGGCCAUAGC-AGC-GGGGUU--CA-C---CCGGU--CCCAUCCGAACCCGGAAGU-UAAG-C-C-C-GC-C-GC-U-ACUG-U--CGC--G-U-AC--UG-UA-U-UGCGA--A-G-UGUAG-GGG-AAA---G-CAG-AC-A-C---GCUUCU-A--C---CACU
-----G-A-CC--CGGCCAUAGU-GGC-CGGGCA-A-CA-C---CCGGU--CUCAUCCGGAACCCGGAAGU-UAAG-C-C-G-GC-CA-CG--U-CAAG--CG--G-C-AG--U--GAGG-UCCGAG--AGG-CCU-C--G-C-AGC-C-G-C-U-CU-G-A---GCUG-GG-A--U---CG
---C---G-A-CC--CGGCCAUAGU-GGC-CGGGCA-A-CA-C---CCGGU--CUCGUUCCGAACCCGGAAGU-UAAG-C-C-G-GC-CA-CG--U-C-A--GAAC-G--GC--CGUGAGG-UCCGAG--AGG-CCU-C--G--CA--G-C-C-GU-U-CUGAGCUG-GG-A--U---CG
---U---U-A-GG--CGGCCACAGC-GGU-GGGGUU-C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-AC-CAGCG-U-UCCG-G--GGA--G-U-AC--UG-GA-G-UGCGCG--A-G-CU-CUGGG-AAA-C-G-C-G-GU-U-C---GCCG-CC-A--C---CA
A--G--U---GG--UGGCCAUUUC-GGC-GGGGUU--C--C-UCCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-AC-CAGCG--U-UCCG-G--CAA--G-U-AC--UG-GA-G-UGCGCG--A-G-CU-CUGGG-AAUUC---C-G-GU-U-C---GCCG-CC-A--C---NN
---U---U-A-GG--CGGCCACAGC-GGU-GGGGUU-C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-AC-CAGCG-U-UCCA-G--GGA--G-U-AC--UG-GA-G-UGCGCG--A-G-CU-CUGGG-AAUUC---C-G-GU-U-C---GCCG-CC-A--C---C
---U---U-A-GG--CGGCCACAGC-GGU-GGGGUU-C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-AC-CAGCG-U-UCCG-G--GGA--G-U-AC--UG-GA-G-UGCGCG--A-G-CU-CUGGG-AAA-C-G-C-G-GU-U-C---GCCG-CC-A--C---CA
---G---U-A-G--CGGCCACAGC-GGU-GGGGUU--C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-AC-CAGCG-U-UCCG-G--GGA--G-U-AC--UG-GA-G-UGCGCG--A-C-CU-CUGGG-AAA-C---CGG-GU-U-C---GCCG-CU-A--C---NN
---U---U---GG--CGGCCAUAGC-GGU-GGGGUU--A--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-GC-CUGCG-U-UCCG-G--UCA--G-U-AC--UG-GA-G-UGCGCG--A-G-CU-CUGGG-AAUUC---C-G-GU-U-C---GCCG-CCUA-----CU
---U---U---GG--CGGCCAUAGC-GG--CGAG-UGAC--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-U-C-GC-CUGCG-U-UUCG-G--UCA--G-U-AC--UG-GA-U-UGGCG--A-C-CU-CUGGG-AAUUC---U-G-AU-U-C---GCCG-CC-A--C---C
-----C-GG--CGGCCAGAGC-GGU-GAGGUU--C--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-U-C-GC-CUGCG-U-UCUG-G--UCA--G-U-AC--UG-GA-G-UGAGCG--A-U-CU-CUGGG-AAUUC---C-A-GU-U-C---GCCG--C---C---UU
A---U--U-ACGG--CGGCCAGAGC-GGU-GAGGUU--C--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-U-C-GC-CUGCG-U-UCUG-G--UCA--G-U-AC--UG-GA-G-UGAGCG--A-U-CU-CUGGG-AAUUC---C-A-GU-U-C---GCCG-CC---C---CU
---U---U---AA--GG--CGGCCAGAGC-GGU-GAGGUU--C--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-U-C-AC-CUGCG-U-UCUG-G--UCA--G-U-AC--UG-GA-G-UGAGCG--A-U-CU-CUGGG-AAUUC---C-A-GU-U-C---GCCG-CC---C---CU
---U---U---AA--GG--CGGCCACAGC-GGU-GAGG--CA-A--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-U-C-GC-CAGCG-U-GAUA-A--CAA--G-U-AC--UG-GA-G-UGUGCG--A-A-CU-CU-GG-AAA-C-A-U-U-AUCU-C---GCCG-CC-U--C---CU
-----A-GG--CGGCCAUAGC-GGC-AGGG-AA-A--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-U-GC-CAGCG-U-UGCG-G--UGA--G-U-AC--UG-GG-G-UGUGCG--A-A-CCC-CUGGG-AAA---GCC-G-U-U-C---GCCG-CCUG--C---CA
---U---U-A-GG--CGGCCACAGC-GGC-GAGGUU--C--C-U-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-U-C-GC-CUGCG-UAU-CG-G--CAA--U-U-AC--UG-GA-G-UGGCG--A-C-CU-CUGGG--AA-C-GUC-G-AU-U-C---GCCG-CC-C--C---CA
---U---A-A-GG--CGGCCAUAGC-GGC-GGGG-CA-A--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-GC-CUGCG-U-UCCG-G--CGA--G-U-AC--UG-GA-G-UACGAG--A-G-UUCU-CUGGGAAAA-C---C-G-GU-U-C---GCCG-CC-U--C---CA
---G---A-A-GG--CGGCCAGAGC-GGU-GGGG-AA-A--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-C-AC-CAGCG-U-ACCG-U--GAA--G-U-AC--UG-GA-G-UGAGCG--A-U-CU-CUGGG-AAC-C-A-C-GAG--U-C---GCCG-CCUG--C---C-
---G---U-A-GG--CGGCCAGAGC-GGU-AGGG-AA-A--C-A-CCCGUA--CCCAUCCGGAACACGGAAGU-UAAG-C-C-U-AC-CAGCG-UAU-CG-U--GAA--G-U-AC--UG-GA-G-UGAGCG--A-U-CU-CAGGG-AAC-C-A-C-GAG--U-C---GCCG-CCUA--C---C-
```

Here for this result set we can see the time difference to align the set of 30 sequences. The time taken by normal MSA is 9.773 s and for aligning the same set of 30 sequences the time taken by the MSA using trie tree is 0.2355.

We can see the improvement in time when MSA with trie trees is used.

## 2. 100 sequences at a time for species Bacteria using Normal MSA

Time taken to align sequences: 105.10149693489075

The aligned sequences are:

```
--G-AC-CU-GG-U--GGCUA--U-G--U-C--G--G--G-U-U-G-G-UUCC--CC-AC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-UCG-UUAA-G-CC-C-U-C--C-A--GA-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-G-C--GC--GG-GAGAG-UA-GG-UC--G-CC-GCCGGGU-C-U
--G-A-C--CC-GG-C-G-GCUA--U-G--U-C-G-G--A-G--G-U--UCC--CC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-UCG-UUAA-G-CC-C-U-C--C-A--GA-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-G-C--GC--GG-GAGAG-UA-GG-UC--G-CC-GCCGGGU-C-U
--G-AC-CU-GG-U--GGCUA--U-G--U-C-G-G--A-G--G-U--UCC--CC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-UCG-UUAA-G-CC-C-U-C--C-A--GA-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-G-C--GC--GG-GAGAG-UA-GG-UC--G-CC-GCCGGGU-C-U
--G-A-C--CC-GG-C-G-GCUA--U-G--U-C-G-G--A-G--G-U--UCC--CC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-UCG-UUAA-G-CC-C-U-C--C-A--GA-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-G-C--GC--GG-GAGAG-UA-GG-UC--G-CC-GCCGGGU-C-U
--G-AC-CU-GG-C--GGCUA--U-G--U-C-G-G--G--G--G-U--UCC--CC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-UCG-UUAA-G-CC-C-U-C--C-A--GA-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-G-C--GC--GG-GAGAG-UA-GG-UC--G-CC-GCCGGGU-C-U
--G-A-C--CC-GG-C-G-GCUA--U-G--U-C-G-G--A-G--G-U--UCC--CC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-UCG-UUAA-G-CC-C-U-C--C-A--GA-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-G-C--GC--GG-GAGAG-UA-GG-UC--G-CC-GCCGGGU-C-U
--G-AC-CU-GG-U--GAUUA--U-G--G-C-G-G--G-U-G-G--CU--GC-AC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-CCG-UGAA-A-CG-C-C-C--C-U--GC-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-A-C--GC--GG-GAGAG-UA-GG-UC--G-UU--GCCAGU-C-U
--G-AC-CU-GG-U--GAUUA--U-G--G-C-G-G--G-U-G-G--CU--GC-AC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-CCG-UGAA-A-CG-C-C-C--C-U--GC-GC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-A-C--GC--GG-GAGAG-UA-GG-UC--G-UU--GCCAGU-C-U
--G-A-C--CU-GG-U--GAUUA--U-G--G-C-G-G--G-U-G-G--CU--GC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-GG-C-GU-GAA-A-CG-CCC--C-U--G--GCC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-A-C--GC--GG-GAGAG-UA-GG-UC--G-UU--GCCAGU-C-U
--G-A-C--CU-GG-U--GAUUA--U-G--G-C-G-G--G-U-G-G--CU--GC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-GG-C-GU-GAA-A-CG-CCC--C-U--G--GCC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-A-C--GC--GG-GAGAG-UA-GG-UC--G-UU--GCCAGU-C-U
--G-A-C--CU-GG-U--GAUUA--U-G--G-C-G-G--G-U-G-G--CU--GC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-GG-C-GU-GAA-A-CG-CCC--C-U--G--GCC--C-GAUGGUA-CU-U-C--GU-C-U-UA-A-G-A-C--GC--GG-GAGAG-UA-GG-UC--G-UU--GCCAGU-C-U
--G-AC-CU-GG-U--GAUUA--U-G--G-C-G-G--G-U-G-G--CU--GC-ACC-C-C-G-AUC-CCAU-UCCGAACUC-G-G-CCC-UGAA-A-CG-C-C-C--C-A--GC-GC--C-GAUGGUA-CU--C-UGU-C-U-UA-A-G-A-C--GC--GG-GAGAG-UA-GG-CC--G-CC-GCCAGG--U
--G-UA--U--GGAGCCGA-UA--G--G--A-C-U-G--G--G--U--UC-AC-U-C--UUCCCCAU-CUCGAAC-A-GAG-UAG-U-UA-A--G-C--C--C-A--GU-ACCGUUGAUGUA-CU--G-C-AAU-C--ACAUG--C--GG-GAAG-C-A--CCAAG-C-UCCAUC-C--
--U--NN-GCU--G-U--GACUA--U--A--A-C-A-G--AUG-U-G-GAAAC--GC--C-U-U-G-CU--CCAU-CCCCGAAC-C-AAG-AAG-CCAA-G-CA-CAU-C--G-U--G--GG--U-GAUGUA-CUAUG-C--C-U-UA-UUG-G-C--AGG--GG-GAAG-UA-GC-UC--A-UU-G-C-GGA-C--
--U--NN-GCU--G-U--GACUA--U--A--A-C-A-G--AUG-U-G-GAAAC--GC--C-U-U-G-CU--CCAU-CCCCGAAC-C-AAG-AAG-CCAA-G-CA-CAU-C--G-U--G--GG--U-GAUGUA-CUAUG-C--C-U-UA-UUG-G-C--AGG--GG-GAAG-UA-GC-UC--A-UU-G-C-GGA-C--
--G--U-CU-GG-U--GGCCA--AAG-CA-C--G-A--G--G--G--CAAAAC-AC-C-C-G-AUC-CCAU-CCCCGAACUC-G-G-CCG-UUAA-G-UGCCGU--C--GC-GC--C-GAUGGUA-CU--G-C-GU-C-A-AA-A-G-A-C--GU--GG-GAGAG-UA--GGAUC--A-CC-GCCA-GA-C-C
--G--C-CU-GA-C-G-ACCA--U--A--A--GCG-A-GU-CG-G--U--CCCCUCC--UUC-CCAU-CCCCGAAC-A-GGA-CCGU-GAA-A-CGAC-UC--U--A--C-GCC--GAUGUA-GU--G-C-GGA-UU--C--CCGU--G--U--GAAAGUA-G-G-UA-A--UC-GUCAGGC-C-U
--G--U-UU-GG-G--GACAA--UAG--C-G-G--U-U-U-G-G-AAC--C-AC-C-C-C-UUC-CCAU-CUCGAACAG-G-G-CCG-UGAA-A-CG-A-A-C--U--U--GC-GC--C-GAUGUA-GU--G-U-ACU-C-U-U--C-GUA--U--GC-GAAG-UA-GG-UC--A-UC-CCCAGG--C-U
--U-C-CU-GG-C-G-ACUA--U--A--A--GCG-A-UU-UG-G-AAC--C-AC-CU--G-AUA-CCAU-CUCGAACUC-A-G-AAG-UGAA-A-CA-UUUC--C--GC-GCC--GAUGUA-GU--G-UGAGG-C-UU--CC--U-CA-U-G--C--GAAAGUA-G-UC-A--UC-GCCAGG--G-U
--UC-CU-GG-C--GACUA--U--A--A--GCG-A-UU-UG-G-AAC--C-AC-C-U-G-AUA-CCAU-CUCGAACUC-A-G-AAG-UGAA-A-CA-UUUC--C--GC-GC--C-GAUGUA-GU--G-UGAGG-C-UU--C-U-U-C--AU--GC-GAAG-UA-GG-UC--A-UC-GCCAGG--G-U
--U-C-UU-GG-C-G-ACUA--U--A--A--GCG-A-UU-UG-G-AAC--C-AC-CU--G-AUA-CCAU-CUCGAACUC-A-G-AAG-UGAA-A-CA-UUUC--C--GC-GCC--GAUGUA-GU--G-UGAGG-C-UU--CC--U-CA-U-G--C--GAAAGUA-G-G-UC-A--UC-GCCAGG--A-U
--UC-CU-GG-C--GACCA--U--A--A--AGC-G-G--U-U-U-G-G-AAC--C-AC-C-U-G-A-CUCCAU-CUCGAACUC-A-G-AAG-UGAA-A-CG-A-A-C--C--C--GC-GC--C-GAUGUA-GU--GUGAGU-C--U--C-C-U-C--AU--GU-GAAG-UA-GG-UC-AG--C-GCCAGG--G-U
--U-C-CU-GG-U-G-ACUA--UAG--C-G-G--U-U-U--G-G-AAC--C-AC-CU--GAU--CCAU-CUCGAACUC-A-G-AAG-UGAA-A-CG--AA--C--C--GC-GCC--GAUGUA-GU--G-UGAGG-C-UU--CC--U-CAU--G--U--GAAAGUA-G-G-UC-A--UC-GCCAGG--G-U
--U-C-CU-GG-C-G-ACCA--UAG--C-G-G--U-U-U--G-G-AAC--C-AC-CU--G-AUJ--CCAU-CUCGAACUC-A-G-AAG-UGAA-A-CG-A-A-C--C--C--GC-GCC--GAUGUA-GU--GUGAGU-C--U--CC--U-CA-U-G--U--GAAAGUA-G-G-UC-A--UC-GCCAGG--G-U
--AG--C-CU-GG-U-G-CCCA--UAG--C-AUG--A-GU--G--A--A--AC-ACA-C--G-AUC-CCAU-CCCCGAACUC-GA-C-GU-GAA-A-C--C-U--A--UAG--C-CU--GAUGUA-GU--A-U-GUC--AU-AAG-U-C-AU--G--G--GAGAGUA-G-G-UC--ACU-GCCAGG--C-U
--A--CA-GG-U-G-ACUA--UAG-CAUC--A-G--G-G--U--CC-AC-CUC--UUC-CCAU-UCCGAAC-A-GAG-AAG-UUAA-G-C-C-C-UGA-UC--C--GCC--GAUGUA-CU--G-C-G--U--A-AC-A--G-U-G--GAGAGUA-G-G-UA-G-UC-CCC--GU--U
--U--CA-GG-U-G-ACUA--UAG-CAUC--A-G--G-G--U--CC-AC-CUC--UUC-CCAU-UCCGAAC-A-GAG-AAG-UUAA-G-C-C-C-UGA-UC--C--GCC--GAUGUA-CU--G-C-G--U--A-AC-A--G-U-G--GAGAGUA-G-G-UA-G-UC-CCC--GU--U
--C-UCAG-GG-U-G-UGUA--U-U--A-C-G-U-U--G-G--G--U--CC-AC-CUC--UUC-CCAU-UCCGAAC-A-GAG-AAG-UUAA-G-C-C-C-CAACGU-GCC--GAUGUA-CU--G-C-G--U--A--A--G--G--GAGAGUA-G-G-AC--G-CC-GCC--G-C-C
--C-UU-GG-U--GAGAA--GAGCUA-C-G-G--G--G--U--AC-AC-C-CAG-A-A-ACAU-UCCGAAC-CUG-G-AAG-U-UA-A--G-C--C--CGUA-AAC-GC--UGAA-AGUA-CU-UG--GA-G-G-GA-A-G-C-C-UCCU--GG-GAG-GAUA-GGAAC--U--GCCAAG--U-U
--C-UU-GG-U--GAGUA--UAGCUA-U-G-G--G--G--U--AC-AC-C-UAG--UJ-ACAU-UCCGAAC-C-UAG-AAG-U-UA-A--G-C--C--CAUA-UAC-GCU--GAUGUA-CU-UG--G--UGGA-A-GCCG--C--U--GG-GAGAG-UAGGAU--U--GCCAAG--C--
--C-UU-GG-U--GAAGAUAAGCU-G--U-G-G--G--G--U--AC-AC-C-U-G-GUC-CCAU-UCCGAACCC-A-GAG-U-UA-A--G-C--C--CGUA-AAC-GC--UGAA-AGUA-CU-UG--GA-G-G-GA-A-G-C-C-UCCU--GG-GAG-GAUA-GGAAC--U--GCCAAG--C--
--UG--U-CU-GG-C--GGCCA--U--A--A--AGC-G-C--A-G-U-G-G-AAC--C-AC-C-C-C-UUC-CCAU-CUCGAAC-A-G-GAGCG-UGAA-A-CG-C-U-G--C--A--GC-GC--C-UGAUGUA-GU-UG-AGGU--C--U--C--C-CU-CCG-GAAG-UC-GG-UC--A-CC-GCCA-GA-C-AC
--U--UU-GG-U--GGUCA--U--A--G--G-C-U-U-G-G--C--UAAAC-AC-C-C-G-AUC-CCAU-CCCCGAACUC-G-G-CAG-UUAA-G-GG-CCA-A-C--A--C--GC-C--GAUGGUA-CU--G-C-GU-C-U-CA-A-G-A-C--GU--GG-GAGAG-UA-GG-UC--A-CC-GCCA-AA-C-C
--C-AA--C--GG-C--GACAA--U--U--U--C--C--C-U-U-G-G--U--GA-ACAC-C-UUCU-CCAU-UCCGAAC-A-GAG-UCG-U-UA-A--G-C--C--CAAG-GA--GA-GC--C-GAUGGUA-CU--G-C-UU-CAU-U--G--C--GG-GAGAG-UA-GG-UC--G-UC-GCC--GU-G-U
--GUUA--C--GG-C-G-GUCA--U--A--A--AGC-GUG--G--G--G--GAAACGCCC-G-GUU-CCAU-UCCGAACCC-G-G-AGCU--A-A-GG-C-C-C-AC--A--GC-GCC--GAUGUA-CU--G-C-AAC--CG-GGAG--G-U--UGU--GG-GAGAGUA-G-G-UC--G-CC-GCC-GGA-C-AA
--A-AC-CCC-CG-C--GCCCA--UAG-CA-C--U--U-G-U-G-G-AAC--C-AC-C-C-C-ACC-CCAU-GCCGAACUC-G-G-UCG-UGAA-A-CA-C-A-G--C--A--GC-GC--C-GAUGUA-CU--GGGCCG--C--A--G-G-G-C--CC--CG-AAAAG-UC-GG-UC-AG--C-GCGGGG-C-U
--U--CU-GG-U-G-AUGA--U-G--G-C-G-A--A-G-A-G--G--U-CAC-ACC-C-C-G-UUC-CCAU-GCCGAACUC-GGA-A-GU--U-UA-A--G-C-C-UUC-AG--C-GCC--GAUGUAU--G-G-GGG--U--C--CCCU--G--C--GAGAGUA-G-G-AC-A--UC-GCCAGG--C--
--U--CU-GG-U--GAUGA--U-G--G-C-G-A--A-G-A-G--G--U-CAC-ACC-C-C-G-UUC-CCAU-GCCGAACUC-G-G-AAG-U-UA-A--G-C-U-UUC--A--GC-GC--C-GAUGGUA-GU-CG-GGG-C-U-U--C-C-C-C--U--GU-GAGAG-UA-GG-AC--A-UC-GCCAGG--C--
--U--CU-GG-U--GAUGA--U-G--G-C-G-A--A-G-A-G--G--U-CAC-ACC-C-C-G-UUC-CCAU-GCCGAACUC-G-G-AAG-U-UA-A--G-C-U-UUC--A--GC-GC--C-GAUGGUA-GU-CG-GGGU--U-U--C-C-C-C--U--GU-GAGAG-UA-GG-AC--A-UC-GCCAGG--C--
```

## 100 Sequence at a time for species Bacteria with MSA with trie tree algorithm

```
Time taken to align sequences: 1.2540323734283447

The aligned sequences are:

--G--A-C--CU-GG-UGGC-UA-U-G-U--C--G-GAG-GU--U-C--C--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-UCG-UU-AA-G-C--C-CU-C--C-A--G--A-GCCG-AUGGUA-C-U-U--C--G-UC-U-U-A-A-G-G-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CC-GG-CGGC-UA-U-G-U--C--G-GAG-GU--U-C--C--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-UCG-UU-AA-G-C--C-CU-C--C-A--G--A-GCCG-AUGGUA-C-U-U--C--G-UC-U-U-A-A-G-G-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGGC-UA-U-G-U--C--G-GAG-GU--U-C--C--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-UCG-UU-AA-G-C--C-CU-C--C-A--G--A-GCCG-AUGGUA-C-U-U--C--G-UC-U-U-A-A-G-G-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CC-GG-CGGC-UA-U-G-U--C--G-GAG-GU--U-C--C--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-UCG-UU-AA-G-C--C-CU-C--C-A--G--A-GCCG-AUGGUA-C-U-U--C--G-UC-U-U-A-A-G-G-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGGC-UA-U-G-U--C--G-GAG-GU--U-C--C--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-UCG-UU-AA-G-C--C-CU-C--C-A--G--A-GCCG-AUGGUA-C-U-U--C--G-UC-U-U-A-A-G-G-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CC-GG-CGGC-UA-U-G-U--C--G-GAG-GU--U-C--C--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-UCG-UU-AA-G-C--C-CU-C--C-A--G--A-GCCG-AUGGUA-C-U-U--C--G-UC-U-U-A-A-G-G-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-A--C--G-G-A-GU-G-G-C--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCG-UG-AA-A-C-G-C-U--C--A--G--C-GCCA-AUGGUA-C-U--G--C--G-UC-U-U-A-A-G-A-C--G--U--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-G--C--G-G-G-GU-G-G-C--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCG-UG-AA-A-C-G-C-C--C--U--G--C-GCCA-AUGGUA-C-U--U--C--G-UC-U-U-A-A-G-A-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-G--C--G-G-G-U-G-GCC--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCG-UG-AA-A--CGC--C--C--U--G--C-GCCA-AUGGUA-C-U--U--C--G-UC-U-U-A-A-G-A-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-G--C--G-G-G-GU-G-G-C--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCG-UG-AA-A-C-G-C-C--C--U--G--C-GCCA-AUGGUA-C-U--U--C--G-UC-U-U-A-A-G-A-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-G--C--G-G-G-GC-G-G-C--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCG-UG-AA-A-C-G-C-C--C--U--G--C-GCCA-AUGGUA-C-U--U--C--G-UC-U-U-A-A-G-A-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-G--C--G-G-G-U-G-GUC--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCG-UG-AA-A--CGC--C--C-A--G--C-GCCA-AUGGUA-C-U--U--C--G-UC-U-U-A-A-G-A-C--G--C--G--G--AGA-G-UA-GG-L
--G--A-C--CU-GG-UGAU-UA-U-G-G--C--G-G-G-GU-G-G-C--U--G--C-A--C-C-G-A-U-C-CCAU-UC-CGAACU-C-G-G-CCC-UG-AA-A-C-G-C-C--C--A--G--C-GCCA-AUGGUA-C--U--U--CU--G-UC-U-U-A-A-G-A-C--G--C--G--G--AGA-G-UA-GG-C
--GU-A----U-GG-AGCC-GA-UAG-GA-C--U-G-GU--G--U--U--C--A--C-U-C--U--UCC-CCA--UCUCGAAC--A-GAG-UAG-U-UA-A--G-C-C--CA-G-UA-C-C-GUUG-AUGUA-C--U--U--CAA-U-C-A-C-A-U--G--C--G--G--AAA-G-CA-A-C
--UN--N-G--CU--G-UGAC-UA-U--A--CA-G-A-U-GU-G-G-A--A--A--C-G--C-CUU-G-C-U--CCAUC-CC-GAAC--CAA-G-AAGC-AA-G-C-A-CAU--C--G-U--G--GUG-AUGUA-C-UA-U--G--CC-U-U-AUU-G-G-CA--G--G--G--AAA-G-UA-GC-L
--UN--N-G--CU--G-UGAC-UA-U-A--A--G-A-U-GU-G-G-A--A--A--C-G--C-CUU-G-C-U--CCAUC-CC-GAAC--CAA-G-AAGC-AA-G-C-A-CAU--C--G-U--G--GUG-AUGUA-C-UA-U--G--CC-U-U-AUU-G-G-CA--G--G--G--AAA-G-UA-GC-L
--G--U--CU-GG-UGGC-CA-A-A-G--A--C-G-A--G-CAAA--A--C-A--C-C-G-A-U-C-CCAU-CC-GAACU-C-G-G-CCG-UU-AA-G-U-G--C--C-G-UC-G--C-GCCA-AUGGUA-C-U--G--C--G-UC-A-A-A-G-A-C--G--U--G--G--AGA-G-U--GGAL
--G--C--CU-GA-CGAC-CA-U-A-G--C--G-A-GUCG-G--U--C--C-A--CUC--C--U--U-C-CCAU-CC-GAACU-C-G-G-CCG-UG-AA-A-C-G-A-C--U--C--U--A--C-GCCA-AUGUA-G-U--G--C--G-GA-U-U--C--C-C--G--U--G--U--AAA-G-UA-GG-L
--G--U--UU-GG-GGAC-AA-U-A-G--C--G-G-U-UU-G-G-A--A--C--C-A--C-C-C-U-U-C-CCA--UCUCGAACU-C-G-G-CCG-UG-AA-A-C-G-A-C--U--U--G--C-GCCA-AUGUA-G-UG-UA--C--U--U--C--GUA-U--G--C--G--G--AAA-G-UA-GG-L
--U--C--CU-GG-CGAC-UA-U-A-G--C--G-A-U-UU-G-G-A--A--C--C-A--C-C-U-G-A-U-A-CCA--UCUCGAACU-C-A-G-AAG-UG-AA-A-CAU-U-U--C--C--G--C-GCCA-AUGUA-G-UG-UG--A--G-GC-U-UCC-U-C-A-U--G--C--G--G--AAA-G-UA-GG-L
--U--C--CU-GG-CGAC-UA-U-A-G--C--G-A-U-UU-G-G-A--A--C--C-A--C-C-U-G-A-U-A-CCA--UCUCGAACU-C-A-G-AAG-UG-AA-A-CAU-U-U--C--C--G--C-GCCA-AUGUA-G-UG-UG--A--G-GC-U-UCC-U-C-A-U--G--C--G--G--AAA-G-UA-GG-L
--U--C--CU-GG-CGAC-CA-U-A-G--C--G-A-U-UU-G-G-A--A--C--C-A--C-C-U-G-A--U--CUCCA--UCUCGAACU-C-A-G-AAG-UG-AA-A-C-GAA--C--C--G--C-GCCA-AUGUA-G-UG-UGA--G--G-UC--U--C--C-U-C--A--U--G--U--AAA-G-UA-GG-L
--U--C--CU-GG-CGAC-UA-U-A-G--C--G-G-U-UU-G-G-A--A--C--C-A--C-C-U-GAA-U--CCA--UCUCGAACU-C-A-G-AAG-UG-AA-A-C-GAA--C--C--G--C-GCCA-AUGUA-G-UG-UG--A--G-GC-U--C--C-U-C--A--U--G--U--AAA-G-UA-GG-L
--U--C--CU-GG-CGAC-CA-U-A-G--C--G-G-U-UU-G-G-A--A--C--C-A--C-C-U-G-A-U-U-CCA--UCUCGAACU-C-A-G-AAG-UG-AA-A-C-GAA--C--C--G--C-GCCA-AUGUA-G-UG-UGA--G--G-UC--U--C--C-U-C--A--U--G--U--AAA-G-UA-GG-L
--AG--U--UU-GG-UGGC-UA-U-A-G--CA-U-G-A-GU--G-A--A--A--C-A--C-A-C-G-A-U-C-CCAU-CC-GAACU-C-G-A-AGC-UG-AA-A--CGC--U--C--A--UA-G--C-GCUA-AUGGUA-C-UA-U--G--G-UC-A-U-A-A-G-U-C--A--U--G--G--AGA-G-UA-AG-L
--A--CA--CG-CGAC-UA-U-A-G--C--GUU-G-G-G-G-A--U--C--C-A--C-CUC--U--U-C-CCAU-UC-CGAAC--A-GAG-AAG-U-UA-A--G-C-C-CU-G--A-UC-A--C-GCCG-AUGGUA-C--U--U--G--C-G-U-A-A-C-A--G--U--G--G--AGA-G-UA-GG-L
UC--A----GG-UGAC-UA-U-A-G--CA-U-C-A-G-G-G-U--U--C--C-A--C-CUC--U--U-C-CCAU-UC-CGAAC--A-GAG-AAG-U-UA-A--G-C-C-CU-G--A-UC-A--C-GCCG-AUGGUA-C--U--U--G--C-G-U-A-A-C-A--G--U--G--G--AGA-G-UA-GG-L
C-----U--CA-GG-UGGU-UA-U-A-G--GUU-G-G-G-G-A--U--C--C-A--C-CUC--U--U-C-CCAU-UC-CGAAC--A-GAG-AAG-U-UA-A--G-C-C-CU-G--A-UC-A--C-GCCG-AUGGUA-C--U--U--G--C-G-U-A-A-C-A--G--U--G--G--AGA-G-UA-GG-L
C-----C--UU-GG-UGAG-AA-U-A-G--CU-GCG-G-G-G-G--U--A--C-A--C-C-U-G-G-U-C-CCAU-UC-CGAAC--A-GAG-AAG-U-UA-A--G-C-C-CU-G--A-UC-A--C-GCCG-AUGGUA-C--U--U--G--C-G-U-A-A-C-A--G--U--G--G--AGA-G-UA-GG-L
C-----C--UU-GG-UGA--GA-A-GAG--CUAG-G-G-G-G-G--U--A--C-A--C-C-CAG-A--A-ACAU-UC-CGAAC--CUG-G-AAG-U-UA-A--G-C-C-CU-G--A-UC-A--C-GCCG-AUGGUA-C--U--U--G--C-G-U-A-A-C-A--G--U--G--G--AGA-G-UA-GG-L
C-----C--UU-GG-UGAG-UA-U-A-G--CUAU-G-G-G-G-G--U--A--C-A--C-C-U-AGU-U-ACAU-UC-CGAAC--CUA-G-AAG-U-UA-A--G-C-C-CU-A--U-A--C-GCUG-AUGGUA-C-U--U--G--C-G-UG-A-A-G--G--GC--C--G--G--AGA-G-UAUGAL
C-----C--UU-GG-UGAG-AA-U-A-G--CU-GCG-G-G-G-G--U--A--C-A--C-C-U-G-G-U-C-CCAU-UC-CGAAC--A-GAG-AAG-U-UA-A--G-C-C-CU-G--A-UC-A--C-GCCG-AUGGUA-C--U--U--G--C-G-U-A-A-C-A--G--U--G--G--AGA-G-UAUGAL
UG--U--U--CU-GG-CGGC-CA-U-A-G--C--G-C-A-GU-G-G-A--A--C--C-A--C-C-C-U-U-C-CCA--UCUCGAAC--A-G-GACCG-UG-AA-A-C-G-C-U--G--C--A--G--C-GCU-AUGUA-G-U--UGAGG--G-UC--U--C--C-C-U--C--G--C--G--AGA-G-UC-GG-L
-----U--UU-GG-UGGU-CA-UAG-G-C-----UU-G-G-C--UAA--C-A--C-C-G-A-U-C-CCAU-CC-GAACU-C-G-G-CAG-UU-AA-G-G-G-C-CA-A-C--A--C-GCCG-AUGGUA-C-U--G--C--G-UC-U-C-A-A-G-A-C--G--U--G--G--AGA-G-UA-GG-L
C--A--A--C--GG-CGAC-AA-U--U--C--C--C-CUU-G-G--U--GAMC-A--C-CUC--U--U-C-CCAU-UC-CGAAC--A-GAG-UC-U--UA-A--G-C-CAGC--G-A--G--G--A--G--G--G--AGA-G-UA-GG-L
--GUUA--C--GG-CGG-CA-U-A-G--C--RUG-G-G-G-G-A--A--C--G--C-C-G-G-U-U-CCAU-UC-CGAAC--C-G-G-AGCU--A--A--G--C--C--CA-A--G--C-GCCG-AUGGUA-C-U--G--CA--A--CCGG-G-A-G-G-U-U--G--U--G--G--AGA-G-UA-GG-L
--A--ACC--CC-CG-CGCC-CA-U-A-G--C--A-C-U-GU-G-G-A--A--C--C-A--C-C-C-A-C-C-CCAU-CC-GAACU-C-G-G-UCG-UG-AA-A-C-A-C-A--G--C--A--G--C-GCCA-AUGUA-C-U--G--G--GCCG-G-C-AGG-G-C--C--C--G--GA-AAA-G-UC-GG-L
-----U--CU-GG-UGAU-GA-U-G-G--C--G-A-A-GA-G--U--C-A--C-A--C-C-G-U-U-C-CCAU-CC-GAACU-C-G-G-AAG-UU-AA-G-C-U-C-U--U--C-A--G--C-GCCG-AUGGUA-G-U--UGG--G--G-GC-U-U--C--C--C--C--CU--G--CG-AGA-G-UA-GG-L
-----U--CU-GG-UGAU-GA-U-G-G--C--G-A-A-GA-G--U--C-A--C-A--C-C-G-U-U-C-CCAU-CC-GAACU-C-G-G-AAG-UU-AA-G-C-U-C-U--U--C-A--G--C-GCCG-AUGGUA--G--U--C--CGGG-GC-U--U--C--C--C--C--CU--G--UG-AGA-G-UA-GG-L
-----U--CU-GG-UGAU-GA-U-G-G--C--G-A-A-GA-G--U--C-A--C-A--C-C-G-U-U-C-CCAU-CC-GAACU-C-G-G-AAG-UU-AA-G-C-U-C-U--U--C-A--G--C-GCCG-AUGGUA--G--U--C--CGGG-GC-U--U--C--C--C--C--CU--G--UG-AGA-G-UA-GG-L
```

In this result set we have taken a set of 100 sequences of species Bacteria and the time taken by normal MSA for aligning these 100 sequences is 105.101 s (about little less than 2 minutes) and the time taken by the MSA with trie tree is only 1.2540 s which is significantly smaller than the former.

**Further results on all five species datasets with different segment lengths are available in the submitted repository.**

## Further work and Conclusion

Scalable and fast MSA is the need of the hour. While the HAlign algorithm decreases the time complexity to linear from squared time, the quality of results may vary because the algorithm uses some heuristics. Once we get the first level of matches on running Aho-Corasick with the centre sequence and any one sequence at a time, we get the indices of these matches of segments. However, the optimal arrangement of these segment matches is again a DP problem. Here we achieve fast times by using a greedy FIFO allotment which might not be optimal always. Practically however, we find the results to be suitable.

Further work should further focus on increasing the stability of the MSA algorithms in the linear time complexity space. Our experiments showed decreasing match quality as the number of sequences input was increased. In conclusion, we found the algorithm to be a great steppingstone towards linear time MSA algorithms. Working directly with code and datasets has given us a very clear image of the computations involved in sequence alignment and its several bottlenecks.