# 天津大学

# 2020-2021 学年第一学期
# 《高级算法》结课作业

题　　目：**Pair Wised Sequence Alignment using Suffix Tree**

学　　院：　　　　智能与计算学部　　　　

专　　业：　　　计算机科学与技术专业　　　

教　　师：　　　　宫秀军老师　　　　

学　　号：　　　　2020244002　　　　

姓　　名：　　　　刘冬冬　　　　

邮　　箱：　　　3463264078@qq.com　　　

**2021 年　1 月 12　日**

# Abstract

[Abstract]：With the continuous deepening of bioinformatics research, a large amount of bioinformatics has been produced and is still in constant flow production. If the processing of this information can be speeded up, the mystery of life can be better revealed to get a promising future.

For this reason, researchers use indexing techniques to improve the processing of sequences. For the smaller sequence, the suffix tree is undoubtedly a good choice, but there is a "memory bottleneck" problem that makes it unsuitable for long genes sequence; suffix array is another good choice, but it occupies more storage space than suffix tree occupies storage space. Although time is less, but the performance of the suffix array is much lower than that of the suffix tree; the search based on and although the citation structure can be used for fast matching, it cannot do anything about gene sequences with low similarity.

In this paper, based on the deep study of the classical suffix tree[1], we present Binary-Tree based on a suffix tree index structure. Binary-Tree restores all short biology sequences, which divided by Longest common sequence[2] which is calculated using suffix-tree. The Binary-Tree based algorithm solve the "memory bottleneck" that the classical suffix tree exists, and then in order to construct a suffix tree, we sort the short biology sequences by alphabetical order and according to the biggest public prefix of the two adjacent biology sequences confirm the branch point. During constructing the suffix tree, and align all the sequence in one time, through experiment, we can see that the length of sequence data is two long, that we can't align sequence in one time, considering computer requirements, such as the memory, and computing time. So we use suffix tree to find the longest common sequence and use the longest common sequence find aforementioned to split original sequence and construct the binary tree. After

constructing the binary tree, we using the middle order algorithm to go over the whole tree to get the alignment sequence and the alignment cost. And for each node, we use Needleman-Wunsch algorithm[3] to align partial sequence and calculate the score using dynamic algorithm.

# Introduction

[Introduction]:

Bioinformatics is not only one of the most important frontiers of life sciences today, but also one of the core areas of science. With the deepening of bioinformatics research and the development of the human genome blueprint Completed, the amount of biological information has increased sharply, especially the number of gene sequences has increased exponentially, with the amount of gene, it is particularly important to retrieve the required information quickly and efficiently in the sequence. The most typical feature of a gene sequence database is that it has an extremely large data size. For example, GenBank[4], a database established by the National Biotechnology Information Center is increasing exponentially. The number of nucleotide bases doubles every month.

Sequence alignment aim to discover maximally similar (or identical) amino acid or identical nitrogenous base positions across the aligned query set of sequences. Sequence wised alignment is an essential preprocess to many bioinformatics analyses, including homology modeling, secondary structure prediction, phylogenetic reconstruction, and protein structure, and function prediction. Sequence alignment is visualized as two-dimensional matrix in which the rows are the individual sequences and the columns are maximally similar or identical amino acid or nitrogenous base positions arranged to correspond by inserting gap characters in appropriate positions. Sequence alignment can provide a wealth of information about the structure or function relationships within a set of sequences, such as the evolutionary

conservation of functionally or structurally important sites, and conserved hydrophobicity patterns in precise regions.

# Methods

[Methods]:

Question: Given genome sequences of mycobacterium tuberculosis (MTB) 35 spices, we are required to design an algorithm to make alignment between any two of spices using suffix tree methods, and the requirements are as follows:

➢ Presenting the principle, workflow and codes of the algorithm;

➢ Implementing the algorithm using c/c++ or python(can modify existing algorithms to meet your needs);

➢ Reporting your results for 35*34 experiments including spend time, used memory of suffix tree, aligned scores and the final alighment;

➢ Comparing your results with BLAST and FASTA algorithm.

Suffix tree pairwise alignment: the basis of our method is a data structure known as a suffix tree, which is a compressed (Aho and Corasick, 1975) [5] containing all the suffixes of a given text as keys, and positions in the text as values, suffix trees allow particularly fast implementations of many important string operations(Baeza-Yates and Gonnet, 1996) [6]. And the definition is that, the tree has exactly n+1 leaves numbered from 0 to n. Except for the root, every internal node has at least two children. Each edge is labeled with a nonempty substring of S, No two edge starting out of a node can have string labels beginning with the same character. The string obtained by concatenating all the string labels found on the path from the root to leaf I spells out a suffix, a simple substring that begins at any position in S and extends to the end of S. The efficient algorithms given by Ukkonen can construct a suffix tree in linear time with O(n) computational space from a string, where n is the length of the string. This methods also reduces suffix tree construction to O(n) time, for

constant-size alphabets, and O(nlogn) in general.

Original methods proposed in the work[7], given two sequences which are needed to be aligned, we call the sequence S1 and sequence S2 for convenience, it is assumed that the sequence to be compared are highly similar. Therefore, there are many common substrings in both S1 and S2. The alignment process consists of four steps, which is show in figure 1.
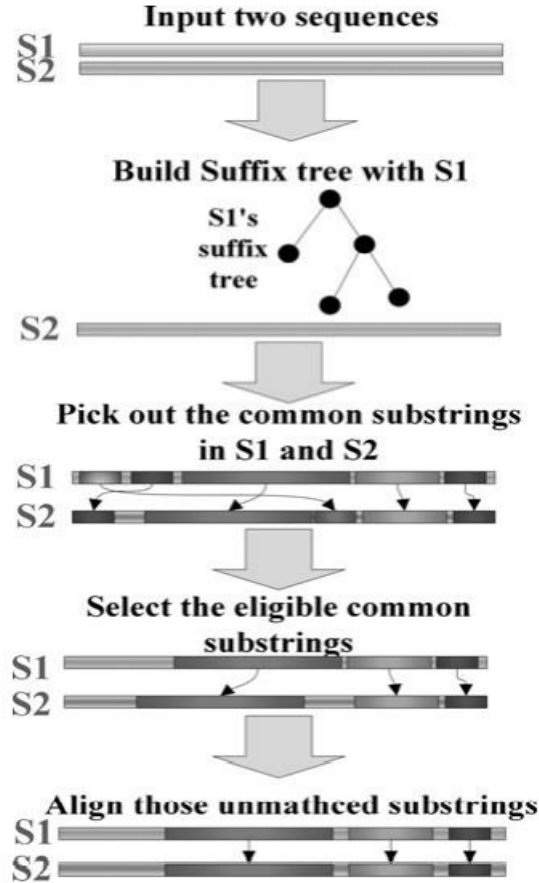


Figure 1: Alignment

The alignment mainly contain four steps, which are as follows:

1. Build a suffix tree from one sequence, here S11 is assumed as the chosen one, and the tree's name is tree-S1;

```python
def _build_generalized(self, xs):
    """Builds a Generalized Suffix Tree (GST) from the array of sequences provided.
    """
    terminal_gen = self._terminalSymbolsGenerator()

    # _xs = [x + next(terminal_gen) for x in xs]
    _xs = []
    n = 0
    beg = 0
    for x in xs:
        end = beg + len(x)
        self.begs.append(beg)
        self.ends.append(end)
        _xs += x + next(terminal_gen)
        n += 1
        beg = end + 1
    self.word = _xs
    self._generalized_word_starts(xs)
    self._build(_xs)
    self.root._traverse(self._label_generalized)
```

Figure 2: Suffix-Trees Construction Code

2. Pick out common sequences.

```python
def find_commom_substrings(seq1,seq2):
    seqindex=0
    index=0
    resultAlignList=[]
    while(index<len(seq2) and seqindex<len(seq1)):
        st=STree.STree4CS([seq1[seqindex:-1],seq2[index:-1]])
        start,end=st.lcsEndStart()
        length=end-start
        startindex=seqindex+start
        seqindex=seqindex+end

        st1 = STree.STree4CS([seq2[index:-1]])
        start1=st1.find(seq1[startindex:seqindex])+index
        if start1>=index and length>0:
            print("start",startindex,"length",length,"seqindex",seqindex)
            print("start1",start1,"length",length,"index",start1+length)
            resultAlignList.append([startindex,start1,length])
            index=start1+length
        else:
            index=index+1
    return resultAlignList
```

Figure 3: Common Sequence Picking Code

3. Pick out common strings in S1 and S2 with sequence S2 and
   Tree-S1 to form two common substring sets;

4. Aligns the remaining unmatched substrings between S1&S2, all substrings are concatenated to form the aligned sequence. And then calculate the score between two sequence. The aligned score id defined by sore(X, Y)=\sum_{i=0}^{n-1}(d(X[i],Y[i]);

```python
def calculate_score_align(seq1,seq2):
    resultList=find_commom_substrings(list(samples[0]),list(samples[1]))
    score=0
    for i in resultList[1:-1]:
        start1=i[0]
        start2=i[1]
        length=i[2]
        str1="".join(seq1[start1:start1+length])
        str2="".join(seq2[start2:start2+length])
        a,b,c=align(str1,str2)
        print("seq1",b,"\n seq2",c,"\n score",a)
        score=score+a
    return score
```

Figure 4: Alignment and Score Calculation

But during experiment, we find the sequence is too long, more than 20000 char, which takes a lot of time to computing, so we improve the algorithm using the binary tree based on the suffix tree. The main idea is that we use suffix tree to find the longest common sequence in the two sequence, and use the sequence aforementioned to split origin sequence to construct binary tree. For the binary tree, each tree node is restore the sub common sequence or two uncommon sequence, for uncommon sequence, we use need Needleman-Wunsch, a kind of dynamic algorithm to calculate the alignment sequence score and get the

```
def createTree(seq1,seq2):
    """
    create a binary tree using digui algorithm, and first find longest comom sequence
    and using the LCS to split original sequence.
    exit requirement: then the sequence length less than 200;
    """
    if len(seq1)==0 and len(seq2)==0:
        return None
    elif len(seq1)>200 and len(seq2)>200:
        try:
            #print("seq1",len(seq1),"seq2",len(seq2))
            st=STree.STree4CS([seq1,seq2])
            #print("FindLSC")
            start,end=st.lcsEndStart()
            st1 = STree.STree4CS([seq2])
            start1=st1.find(seq1[start:end])
            #print("construct PairsequenceItem")
            pair=PairsequenceItem(seq1[start:end],seq2[start1:start1+end-start],start,start1,common=True)
            root=Node(data=pair)
            #print("construct left node")
            root.left=createTree(seq1[0:start],seq2[0:start1])
            #print("construct right node")
            root.right=createTree(seq1[end:-1],seq2[start1+end-start:-1])
            #print("return node")
            return root
        except :
            print("error seq1",len(seq1),"seq2",len(seq2))
    else:
        #print("uncommon sequence:\tseq1",seq1,"seq2",seq2)
        pa =PairsequenceItem(seq1,seq2,0,0,common=False)
        return Node(data=pa)
```

Figure 5: Binary Tree Construction

two alignment sequence. Then we using a middle order algorithm
to go over the whole tree to get the alignment sequence and the
alignment cost. The proposed methods takes much little time and
memory to get the results.

```
def inorderTraversal(root):       #先把根节点入栈，如果左子树一直不为空，就一直入栈，直到
                                  ##把所有左节点入栈，然后pop栈顶元素，指针指向栈顶元素的右子树
    """
    :type root: TreeNode
    :rtype: List[int]
    :return: ordered sequence aligned segment list
    """
    stack=[]
    res=[]
    if not root:
        return []
    while root or stack:
        while root:
            stack.append(root)
            root=root.left
        if stack:
            a=stack.pop()
            root=a.right
            res.append(a.data)
    return res
```

Figure 6: Binary Tree Reverse

```python
def handleAll():
    """
        hande all a pair of sequence
    """
    data = read_from_file_with_enter(data_path)
    keys,samples= statistics_length(data)
    size=len(samples)
    for i in range(0,size):
        for j in range(i+1,size):
            print("handle sequence i=",i,"\t j=",j,"\t",keys[i],"\t",keys[j])
            memomybefore=psutil.Process(os.getpid()).memory_info().rss / 1024
            start = clock()
            strA,strB,score=handleSingle(list(samples[i]),list(samples[j]))
            memomyAfter=psutil.Process(os.getpid()).memory_info().rss / 1024
            finish = clock()
            analyse="\t score:"+str(score)+"\t memory:"+str(memomyAfter-memomybefore)+"\t time:"+str(finish-start)
            saveToFile(keys[i]+"VS"+keys[j],analyse,strA,strB)
            print("save to file:",keys[i],"VS",keys[j],"result:",analyse)
```

Figure 7: Comparison between each other

## Results

[Results]：Here in these section, we will first introduce the experiment environments such as the computer parameters, the data description and the result of the code for 35*34 experiments including spend time, used memory of suffix tree, alignment results and the final alignment, as well as the comparison between BLAST and FASTA algorithm.

Computer settings: processor: Inter(R) Core(TM) i7-9750H CPU @2.60GHz 2.59GHz; memory: 16.0GB; GPU: NVIDIA GeForce GTX 1650。

Data description: the data is provided by teacher, thirty strains of SARS cov-2, which is contained in fasta format file showed in figure 8, and the data about thirty strains of SARS cov-2 distribution is showed in the figure 9, from the figure 9 ,the average count of the dna sequence is about 30000.



```
>MW411947.1 |Severe acute respiratory syndrome coronavirus 2 isolate SARS-CoV-2/human/West Bank/AAS24/2020, complete genome
GGTTTATACCTTCCCAGGTAACAAACCAACCAACTTTCGATCTCTTGTAGATCTGTTCTC
TAAACGAACTTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCACTCACGCA
GTATAATTAATAACTAATTACTGTCGTTGACAGGACAGAGTAACTCGTCTATCTTCTGC
AGGCTGCTTACGGTTTCGTCCGTGTTGCAGCCGATCATCAGCACATCTAGGTTTTGTCCG
GGTGTGACCGAAAGGTAAGATGGAGAGCCTTGTCCCTGGTTTCAACGAGAAAACACACGT
CCAACTTAGTTTGCCTGTTTTACAGGTTCGCGACGTGCTCGTACGTGGCTTTGGAGACTC
CGTGGAGGAGGTCTTATCAGAGGCACGTCAACATCTTAAAGATGGCACTTGTGGCTTAGT
AGAAGTTGAAAAAGGCGTTTTGCCTCAACTTGAACAGCCCTATGTGTTCATCAAACGTTC
GGATGCTCGAACTGCACCTCATGGTCATGTTATGGTTGAGCTGGTAGCAGAACTCGAAGG
CATTCAGTACGGTCGTAGTGGTGAGACACTTGGTGTCCTTGTCCCTCATGTGGGCGAAAT
ACCAGTGGCTTACCGCAAGGTTCTTCTTCGTAAGAACGGTAATAAAGGAGCTGGTGGCCA
TAGTTACGGCGCCGATCTAAAGTCATTTGACTTAGGCGACGAGCTTGGCACTGATCCTTA
TGAAGATTTTCAAGAAAACTGGAACACTAAACATAGCAGTGGTGTTACCCGTGAACTCAT
GCGTGAGCTTAACGGAGGGGCATACACTCGCTATGTCGATAACAACTTCTGTGGCCCTGA
```
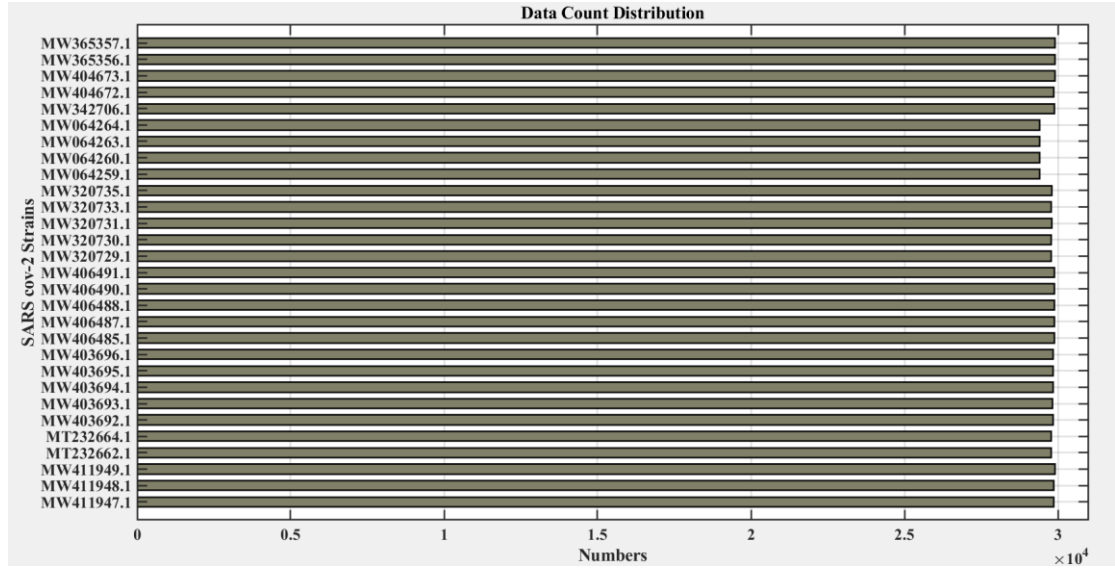
Figure 8: Data Structure

Figure 9: Data Number Distribution

Alignment Score: using score calculation metrics provide in the assignment, sore(X,Y)=\sum_{i=0}^{n-1}(d(X[i],Y[i])); If X[i] and Y[i] both are not gaps, then d(X[i],Y[i]) is defined by the substitution matrix in table 1; If X[i] or Y[i] is a gap, then d(X[i],Y[i])=-4 for opening gap, d(X[i],Y[i])=-1 for extending gap. For each sequence, we are supposed to calculate following 29 sequence, align each pair of them and calculate the aligned sequence score. Here the 30*30 matrix score are showed in the figure 10;

|   | A | C | G | T |
|---|---|---|---|---|
| A | 10 | -5 | 0 | -5 |
| C | -5 | 10 | -5 | 0 |
| G | 0 | -5 | 10 | -5 |
| T | -5 | 0 | -5 | 10 |

Table 1

Figure 10: Aligned Score Matrix

Process Monitor: for each aligned score matrix calculation progress, we use clock() function and psutil module provided by python api to get the information about the code runtime conditions. And the time matrix is showed in the figure 11, and the memory matrix is showed in the figure 12. And the aligned sequence is partial showed in the figure 13, due to the size.



Figure 11: Time Matrix

Figure 12: Memory Matrix



Figure 13: Aligned Sequence

Comparison with BLAST and FASTA algorithm: BLASTN and BLASTP use a rolling window to break down a query sequence into words and word synonyms that form a search set. At least two words

or synonyms in the search set must match a target sequence in the database, for that sequence to be reported in the results. The BLAST algorithm is showed in the figure 14, and the Fasta algorithm mainly has five steps, which is showed as follows:

1. Identify common k-words between I and J;

2. Score diagonals with k-word matches, identify 10 best diagonals;

3. Rescore initial regions with a substitution score matrix;

4. Join initial regions using gaps, penalize for gaps;

5. Perform gynamic programming to find final alignments.



Figure 14: Blast Algorithm workflow

# Conclusion

[Conclusion]：Based on the deep study of the classical suffix tree, we present Binary-Tree based on a suffix tree index structure. Binary-Tree restores all short biology sequences, which divided by Longest common sequence which is calculated using suffix-tree. The Binary-Tree based algorithm solve the "memory bottleneck" that the classical suffix tree exists, and then in order to construct a suffix tree, we sort the short biology sequences by alphabetical order and according to the biggest public prefix of the two adjacent biology sequences confirm the branch point. During constructing the suffix tree, and align all the sequence in one time, through experiment, we can see that the length of sequence data is two long, that we can't align sequence in one time, considering computer requirements, such as the memory, and computing time. So we use suffix tree to find the longest common sequence and use the longest common sequence find aforementioned to split original sequence and construct the binary tree. After constructing the binary tree, we using the middle order algorithm to go over the whole tree to get the alignment sequence and the alignment cost. And for each node, we use Needleman-Wunsch algorithm to align partial sequence and calculate the score using dynamic algorithm.

During solve the task the teacher assigned, I lean something about DNA sequence matching through the internet and am very amazed by the beautiful data structure, which is of so much power to solve the problem in reality. By solving the question, I also improve my coding ability as well as algorithm understanding and constructing ability. Thanks for sincerely teaching.

# References

[1] McCreight, Edward M. "A space-economical suffix tree construction algorithm." Journal of the ACM (JACM) 23, no. 2 (1976): 262-272.

[2] Bergroth, Lasse, Harri Hakonen, and Timo Raita. "A survey of longest common subsequence algorithms." In Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000, pp. 39-48. IEEE, 2000.

[3] Likic, Vladimir. "The Needleman-Wunsch algorithm for sequence alignment." Lecture given at the 7th Melbourne Bioinformatics Course, Bi021 Molecular Science and Biotechnology Institute, University of Melbourne (2008): 1-46.

[4] Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J. and Sayers, E.W., 2010. GenBank. Nucleic acids research, 39(suppl_1), pp.D32-D37.

[5] Aho, Alfred V., and Margaret J. Corasick. "Efficient string matching: an aid to bibliographic search." Communications of the ACM 18, no. 6 (1975): 333-340.

[6] Baeza-Yates, Ricardo A., and Gaston H. Gonnet. "Fast text searching for regular expressions or automaton searching on tries." Journal of the ACM (JACM) 43, no. 6 (1996): 915-936.

[7] Su, Wenhe, Xiangke Liao, Yutong Lu, Quan Zou, and Shaoliang Peng. "Multiple sequence alignment based on a suffix tree and center-star strategy: a linear method for multiple nucleotide sequence alignment on spark parallel framework." Journal of Computational Biology 24, no. 12 (2017): 1230-1242.