

# Finding Transposons in DNA, A Python Approach

Louis Chatta

*University of California, Santa Cruz: BME 160*

## Abstract

Jumping genes, or Transposons, can be found all throughout many biological domains. They are DNA segments that are constantly moving around a genome. These types of DNA segments have a wide range of use in the genome, from creating new mutations, to developing anti-biotic resistance in bacteria, to even DNA fingerprinting in genetic testing. Finding transposons manually can be daunting. It is difficult to search a sequence of DNA that may be thousands of lines long to find all IR(inverted repeats) or DR(direct repeats). This paper outlines a new method of finding all IRs and DRs in a genetic sequence from a FASTA file

## 1 Introduction

Transposons exist in two different forms Replicative(Retrotransposons, Class I), which copy DNA segments, and Non-replicative (Class II), which cuts the DNA segments. Both classes 'paste' into their target DNA. The Class I first transcribe the DNA into RNA and then use reverse transcriptase to make a DNA copy of the RNA to insert in a new location. The Class II use the Transpose enzyme, which binds with the IRs(Inverted Repeats), which then cleaves the phosphodiester bond at each end of the segment, then the nicked hairpins attack the new target DNA, which is then incorporated through the use of DNA Polymerase and DNA Ligase. This process will leave short direct repeats before and after the inverted repeats.[2]

There are many different subclasses of the Class I and Class II transposons. In Class I, there exists TEs(Transposable Elements with Long Term Repeats), SINES(Short Interspersed Transposable Elements) and LINES(Long Interspersed Transposable Elements). The LINES and SINES lack long term repeats, and are encoded with RNA polymerase I and II respectively. The

TEs with long term repeats, encode reverse transcriptase, similar to retroviruses like HIV. Combined, Class I transposons represent 30 % of the human genome, while Class II comprising of only 2 % of the human genome.[2] Another type of TE, is the autonomous TE. Autonomous TEs can move by themselves, while non-autonomous TEs need another TE to move. That is because dependent TEs lack transposase or reverse transcriptase.

Overall many transposons are useless or harmful, but we can use them as a tool of fingerprinting DNA. They provide a good method of detailing inheritance. Finding transposons in the same sequence at the same positions gives direct evidence that two organisms have a common ancestor. An example of this is discovery that Whales and Hippos share a common ancestor. See figure [1].

Finding all the transposons in a sequence can be tremendously useful, but quite difficult. For example, the Drosophila genome has 175,000,000 base pairs. That would take 1.8 years to just read the sequence figure [2]. Taking a programmatic approach to this problem, would dramatically reduce the computation time. This paper will discuss a python program that can take in any amount of sequences from a FASTA file and process each sequence and output all the possible Transposons by finding all IRs and DRs in a given sequence.

## 2 Methods

### 2.1 Program Overview

This program, findRepeats.py, reads in FASTA file with  $N$  amount of sequences, and outputs to the STDOUT all sequences with DRs and IRs, which can be investigated further to be true Transposons. The program runs fast on short sequences, but begins to slow down on larger ones.

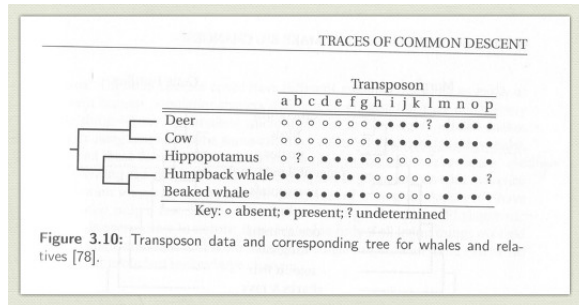


Figure 1: This is a diagram of the relationship between whales and hippos. Based on Transposons. [1]

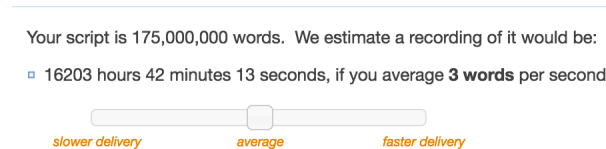


Figure 2: This is an estimated of how long it would take a human to read the Drosophila genome. [3]

## 2.2 Program Architecture

### 2.2.1 Overview

The program is split into the following components

```

1 # previously written library, that reads a
2 # sequence from a FASTA file
3 from sequenceAnalysis import NucParams,
4 FastReader, ProteinParam
5
6 # handy utility class
7 class Util():
8
9 # handles all command line arguments
10 class CommandLine():
11
12 # handles all of the IRs and DRs
13 class RepeatFinder():
14
15 # handles all the flow
16 main():

```

Listing 1: overview

### 2.2.2 SequenceAnalysis.py

This is a program written from a previous lab. Its functionality is to read in FASTA files and generate readable sequences. Here is how it used:

```

1 from sequenceAnalysis import NucParams,
2 FastReader, ProteinParam
3
4 def main():
5     for head, seq in orfReader.readFasta():
6         nucParams = NucParams('')

```

```

6     nucParams.addSequence(seq)
7     nucParams.buildNucleotide()
8     #nucParams.codons now has the sequence

```

Listing 2: sequenceAnalysis.py

### 2.2.3 Util Class Overview

This class is used a utility for the whole program. It contains two functions, one that returns the inverted repeat of a sequence and one that returns the complement strand. The code is below:

```

1 class Util():
2
3     # Function that returns the inverted repeat
4     # of a sequence
5     def getInvertedRepeat(self, string):
6         return self.getComplementStrand(string[::-1])
7
8     # Function that returns the complement of a
9     # given strand
10    def getComplementStrand(self, bases):
11        reverseBase = []
12        for base in bases:
13            if base == 'A':
14                reverseBase.append('T')
15            elif base == 'T':
16                reverseBase.append('A')
17            elif base == 'C':
18                reverseBase.append('G')
19            elif base == 'G':
20                reverseBase.append('C')
21            else:
22                reverseBase.append(base.upper())
23
24        return reverseBase
25
26    ...
27
28    A nice function that returns the reverse
29    strand
30    ...
31
32    def getReverseStrand(bases):
33        reverseBase = []
34        for base in bases:
35            if base == 'A':
36                reverseBase.append('T')
37            if base == 'T':
38                reverseBase.append('A')
39            if base == 'C':
40                reverseBase.append('G')
41            if base == 'G':
42                reverseBase.append('C')
43
44        return reverseBase[::-1]

```

Listing 3: Util Class

### 2.2.4 RepeatFinder Class Overview

This class is used for all the heavy lifting of finding all of the IR and DR in the sequence. It takes in the sequence that was generated from the sequenceAnalysis.py and then finds all the possible sub-sequences and from

those sub-sequences it finds all of IR and DR. Here is the overview of the class:

```

1 class RepeatFinder():
2
3     def __init__(self, seq, header, revSeq=False):
4         # FASTA Sequence
5         self.seq = seq
6
7         # Reverse Sequence
8         self.revSeq = revSeq
9
10        # Header of FASTA
11        self.header = header
12
13        # All the sub-sequences
14        self.subStrings = []
15
16        # An instance of utility class
17        self.util = Util()
18
19
20    # Function that creates all of the possible
21    # contiguous sub-sequences
22    def getAllSubstrings(self, string):
23
24    # Function that finds all the Direct Repeat
25    # sequences and prints to the STDOUT
26    def findAllRepeats(self):
27
28    # Function that finds all the Inverted
29    # Repeat sequences and prints to the STDOUT
30    def findAllInversionRepeats(self):
31
32    # Function that builds a sequence from the
33    # original sequence based on the start and
34    # end positions
35    def buildSequence(self, startRepeat,
36                      endRepeat, startPos, endPos):
37
38    # prints all of the transposons
39    def printTransposons(self):

```

Listing 4: RepeatFinder Class

## 2.2.5 RepeatFinder.getAllSubstrings(self,string)

This function is the one that is the most computationally complex. It loops through every element in a string and computes all possible combinations from that position to the end of the string, which runs at  $O(n^2)$ . In an attempt to optimize this algorithm, I added a line, (6), to skip all sequences smaller than 1. After running a few test cases, it appears that this algorithm is the bottleneck of the program, under larger files. In the future, I hope to optimize this algorithm further. The code is below:

```

1 def getAllSubstrings(self, string):
2     length = len(string)
3     alist = []
4
5     # loop through every element in the
6     # string
7     for i in range(length):

```

```

7
8         # loop through every possible combination
9         # of the current index to the end of the
10        string
11        for j in range(i, length):
12            if len(string[i:j]) > 7 or len(string
13            [i:j]) < 1: continue
14            newString = string[i:j + 1]
15            if newString not in alist:
16                alist.append((newString, i, j+1))
17            self.subStrings = alist
18        return alist

```

Listing 5: getAllSubstrings(self)

## 2.2.6 RepeatFinder.findAllRepeats()

This function is one that loops through all of the possible sub-sequences and tries to find all the Direct Repeats in a sequence. Here is some psuedo-code to help you understand. Loop through all of the sub-sequences, if it has not been seen before, add it to a list of seen sequences. If we encounter a sub-sequence that has been seen before, then it is a repeat! . Then we print the sequence from the first repeat to the second repeat. Here is the actual code:

```

1 def findAllRepeats(self):
2     seenRepeats = []
3     for subString, startPos, endPos in self
4     .subStrings:
5         foundRepeat = ([item for item in
6         seenRepeats if item[0] == subString])
7         if not foundRepeat:
8             seenRepeats.append((subString,
9             startPos, endPos))
10        else:
11            print(' ', self.
12            buildSequence((foundRepeat[0])[0],
13            subString, (foundRepeat[0])[1], endPos))

```

Listing 6: findAllRepeats(self)

## 2.2.7 RepeatFinder.findAllInvertedRepeats()

This Function follows the same structure as the previous one, with only one addition. Instead of checking if we encounter the same sequence again, we check if we encounter the Inverted Repeat of the original sequence. Here is the code:

```

1 def findAllInversionRepeats(self):
2     seenRepeats = []
3     for subString, startPos, endPos in self
4     .subStrings:
5         invertedRepeat = self.util.
6         getInvertedRepeat(subString)
7         if subString == invertedRepeat:
8             continue
9         foundRepeat = ([item for item in
10        seenRepeats if item[0] == invertedRepeat])
11
12        if foundRepeat:

```

```

9         print(' ', self.
    buildSequence((foundRepeat[0])[0],
    subString,(foundRepeat[0])[1],endPos))
10     else:
11         seenRepeats.append((subString,
    startPos,endPos))

```

Listing 7: findAllInvertedRepeats(self)

### 2.2.8 RepeatFinder.buildSequence()

This function takes in the positions of the start and end of the new sequence, and the start and end repeats of the sequence. It then splices the original sequence according to the parameters, and returns it to the caller. Here is the code:

```

1  def buildSequence(self, startRepeat,
    endRepeat, startPos, endPos):
2
3      startString = ''.join(startRepeat)
4      endString = ''.join(endRepeat)
5
6      innerSequence = ''.join((self.seq[
    startPos+len(startRepeat):endPos-len(
    endRepeat)]))
7
8      return '{} {} {}'.format(startString,
    innerSequence, endString)

```

Listing 8: buildSequence(self)

## 3 Discussion

After implementing all the functions and the classes defined in the methods, I went ahead and ran the program on multiple FASTA files. One test file, and one encoding the AJ009736 locus on the Drosophila Melanogaster genome [4]. The program returned all the IRs and DRs found in all the sequences.

On the test file, the program ran instantly and outputted every IR and DR for all the sequences. A sample run can be seen at figure[3]. The output is in the following format:

- [sequence separated by spaces]
- 'All Direct Repeats'
- [DR] [sequence] [DR]
- 'All Inverted Repeats'
- [IR] [sequence] [IR]

All of the output looked good, and was in good format. However, there were a few things that could have been done better. For example when running the Drosophila Melanogaster, AJ009736 FASTA file, the speed was much slower. The output completed in 5 mins, as opposed to instantly before. I

```

A A T T A C A T C A T A A
Direct Repeats
AT TAC AT
CA T CA
CAT CAT
AT TACATC AT
TA CATCA TA
AA TTACATCAT AA
Inverted Repeats
AAT ATT
AA TT
TTA CATCA TAA

```

Figure 3: This is the sample output of the sequence: AAT-TACATCATAA

suppose this is due to the slow nature of the algorithm, *RepeatFinder.getAllSubstrings(self,string)*, which runs at  $O(n^2)$ . In addition to the slow nature of the program, I would maybe think about also outputting the positions of the sequences as well as printing them. I didn't do so in this version, because I wanted to keep the formatting looking nice.

## 4 Conclusion

Finding transposons in genomes can be very useful. From genetic fingerprinting to studying new mutations, finding all the transposons in a genome can really help researchers. Doing so by hand is daunting, but when utilizing the power of computation with the power of Python, a solution emerges. *findRepeats.py*, allows researchers to input any number of sequences and immediately find out all the possible transposons, reducing the number of hours that would have been spent doing so by hand. Moving forward, it would be helpful to optimize the *findAllSubstrings* code, in order to process larger data sets faster. Another approach could be a multi-threaded approach, to increase speed. Overall, this program is a good start in solving the Transposon Finding problem.

## References

- [1] Masato N. "The Evidence For Evolution". In: *Nat'l Acad. of Sciences* 98.13 (2001), pp. 7384-7389.
- [2] *Transposons: Mobile DNA*. URL: <http://users.rcn.com/jkimball.ma.ultranet/BiologyPages/T/Transposons.html>.
- [3] *Words to Time Calculator*. URL: <http://www.edgestudio.com/production/words-to-time-calculator>.

Figure 4: This is one of the outputs of the Drosophila IR Sequence, Transposon