# A Software Parser for Regular Expressions Using Python
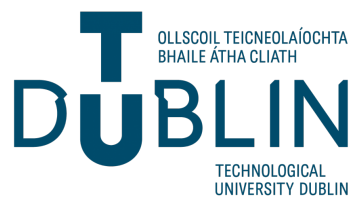
Carl Brady

November 2019

TU Dublin Blanchardstown

# Contents

**Abstract**

This document describes in detail, a parser for regular expressions that is created using python, a programming language. The parser analyses a combination of sentences and returns an output of the bracketed phrasal structure of the sentence. This document will cover the steps taken in creating such a program and an explanation of how each method works and the necessary files needed. To create this system, the author implemented an agile methodology to ensure that there was sufficient time to research the project and test each use case. It also meant that the system was continuously improved and could be reviewed on a regular basis.

# 1   Introduction

Computational Linguistics is a subsidiary of Computer Science and Linguistics. It is used in order to help computers understand human languages. Although, speaking in ones native tongue may be an action a person can do without speaking, for a computer this process can be extremely difficult [3]. We as humans, do not need to actively think when creating regular sentences or simple conversation, this act comes naturally to us. In order for a computer to act in the same way, they too must have a method of accessing this ability to understand the complex structure of sentences that we usually do not consider such as index of all words and their meaning, syntax and also the morphology of the language. Spoken language is even harder for a computer to understand, factors such as pitch, pronunciation of words or accents also must be taken into consideration when creating such a system.

Work in this field can be seen as early as the 1950's with concepts such as machine translation but it was not until the late 1980's / 1990's that real improvements were noticed [5]. The corpus based approach, which has been very successful, so successful it has even had an impact the way language is thought [2], has helped in pushing a new life into the field. In modern times the field has seen expansion in terms of voice recognition and machine learning [6]. Siri[1], Apples virtual assistant is an example of this technology used today. This report describes a program based on a simplified corpus based approach using lexicon, a rule-set and bottom-up-parsing.

# 2   Lexicon & Rules

This section will detail the lexicon used for the making of the software parser. A lexicon is a group of lexical categories or words that essentially describe a language [4]. A lexicon can be used to detail information about languages such as; words, morphological structures, syntax and how words are pronounced. The lexicon the author decide to create was based on words, syntax and morphological structures. The rules that were created for this program are based upon the book 'Natural Language Processing with Python' by Bird et al[4].

## 2.1   Lexicon

The lexicon for this program is created with a group of listings in a .txt file. Every word available from the set of example sentences is combined with a part of speech (POS) tag. The POS tag is based on the NLTK abbreviations for their POS Tagger [4]. For example, each line in the text file had a word and its associated POS tag (eg, DT the). In A.4 the full lexicon.txt file is available.

## 2.2 Rules

The rules used in order to create this program were made in a similar fashion, each line in the rule.txt file was made up of a left hand side argument and a right hand side argument. The left hand side of the argument is the parent to the right hand side argument. For example, each line in the text had the parent specified followed by its children (NP → DT NN). In A.5 the full rules.txt file is available

# 3 System Architecture

This section covers the system architecture, all of the methods used in the program and the overall code design are detailed in this section. The software program is made up of six components; 'main.py', 'parser.py', 'file_manager.py', 'lex.txt', 'rules.txt' and 'tests.txt'. The three text files are used to each store one of the following; lexicon, rules, test cases. Three python files were created, 'main.py' which acts the main file and calls upon methods from other files. 'file_manager.py' handles the opening and reading of all .txt files. 'parser.py' contains the methods for POS tagging, grammar checking, bracketed formatting and creating a tree. A class diagram used in the creation of this system can be seen in A.7.

## 3.1 File Manager

A file manager was used in order to read .txt files into the software program. Each .txt was stored as an array of sentences so that each element in the array could be accesed at a later stage.

```
1    lex_file = open("lex.txt","r") #Open and read file that contains lexicon
2    lex_arr = lex_file.readlines()
3    lex_file.close()
4    return lex_arr
```

As seen in the above code segment, the lexicon file is read using the open() method, which takes two variables , the name of the file and the permission needed (r = read, w = write, etc). Once the file is read, the method readLines() then stores each line of the file as an individual item in an array, which in this case we assign to the variable lex_arr. The file is then closed as it is no longer needed and the array is returned.

## 3.2 POS Tagging

The POS tagger was implemented by creating a method pos_tagger(x,tests,lex_arr), this method takes three arguments. 'X' is the 'test select' input from the user (further explained in section 3.6), 'tests' is the array of test cases and 'lex_arr' is the lexicon returned as an array as seen in 3.1. The selected test case was then split using the .split() method and stored in an array 'test'.

```
1    for word in test:
2        for lexi in lex_arr:
3            lex = lex_arr[j].split()
4            if test[i] == lex[1]:
5                sentence += lex_arr[j]
6            j+= 1
7        j = 0
8        i += 1
```

As can be seen above, the POS tagger checks each word in the test case sentence against every lexical entry's right hand side argument in an attempt to find a match. If a match is found, the matching lexicon entry is then appended to the string 'sentence', this continues until all words in a test sentence have been checked, the method then returns 'sentence' on completion.

## 3.3 Parser

A parser was implemented based on the bottom up parsing method. The method grammar_check(pos, rules) takes two arguments, the first is the sentence that was created by the POS tagger and the second is the rules that have been stored in an array from the getRules() method in file_manager.py.

```
1    pos_split = pos.split()
2    pos_NP = pos_split[0] + ' ' + pos_split[2] # gets the POS tag for the first
     two words and adds to string
3    pos_VP = pos_split[4] # gets the POS tag for the third word and adds to
     string
4    pos_NP2 = pos_split[6] + ' '+pos_split[8] + ' ' + pos_split[10] # gets the
     POS tag for the last three words and adds to string
```

Firstly, the method splits the POS tagged sentence. Three strings are then created based on the split and as the POS tag only appears every second word in the array only the even indexes are required. As it is known there are three phrases (NP, VP, NP) that essentially make up the correct sentence. The method first splits them accordingly, the first two tags are created as a string, etc.

```
1    if pos_split[1] in 'a' and pos_split[3] == 'men': # checks for 'a men' in a
     sentence
2        print('\'A men\' is grammatically incorrect.')
3        return False
4    elif pos_split[1] in 'a' and pos_split[3] == 'women': # checks for 'a women'
      in a sentence
5        print('\'A women\' is grammatically incorrect')
6        return False
```

The parser / grammar checker also checks if there is any immediate grammar mistakes that can rule this sentence out from being a coherent sentence and if so, the method returns false.

```
1    while i<3: # Checking the if the childern of NP that has a parent of S are
     correct
2        rule = rules[i].split('->')
3        x = rule[1].split() # rule[1] is only split here again to fix the
     formatting and remove extra spacing
4        y = x[0]+' '+x[1] # here we assign formatted x to y
5        if pos_NP == y:
6            checker += rule[0] # rule[0] is the left hand side of the matching
     rule
7        i += 1
8    i=4
```

Above we can see how a loop is created in order to ensure the children of NP are correct. As we know, all of the rules that cover the first Noun Phrase(NP) grouping in the sentence are covered in the first three lines of the rule set, it does not make sense to examine each rule in the array, therefore the bound is set as i < 3. The program then checks if the right hand side of the rule matches the tags from the test sentence and if they match it appends the right hand side argument to a string. This is done for all three possible parent tags (NP,VP,NP), until it is possible whether to determine if the sentence is correct or not, if the sentence is correct the method returns true.

## 3.4 Bracketed Phrasal Formatting

The bracketed phrasal method is only called if the grammar_check returns true. Once this happens, the bracketed phrasal takes the POS tagged sentences and uses .split. Once the sentence has been split, it is recreated with the added syntax required to create a bracketed phrasal structure.

```
while i<3:
    br_sentence += '[' + split[i] + ' ' + split[i+1] + ']'
    i+=2
br_sentence += ']][VP[VP'
```

As seen in the above segment, the sentence is rejoined using a loop for each phrase contained within a sentence.

## 3.5 Tree Structure

In order to create a tree structure like display on the command line, the program imports treelib, a package that assists creating trees with python.

```
split = sentence.split()
tree = Tree()
tree.create_node("S1", "s1")
tree.create_node("S", "s", parent="s1")
```

Once again we can see the sentence is split, from here we can add nodes to the tree with .createnode() this method takes two arguments, how the node will be displayed in the tree and its ID. As the first node in the tree does not have a parent there is none listed, for every other node in the tree a parent is listed.

## 3.6 Main File

The main file simply calls on all of the methods mentioned above.

```
while True:
        x = input('\n\nInput a test number, type \'help\' to see all test cases
    or \'exit\' to exit the program\n')
        if x in 'help':
            print('-----------------------------------------------------------')
            print_tests()

        elif x in 'exit':
            return
```

Once we run this program we can see the program prompts a user to input a test number or type 'help' to see the test cases or 'exit' to exit the program.



```
Input a test number, type 'help' to see all test cases or 'exit' to exit the program
```
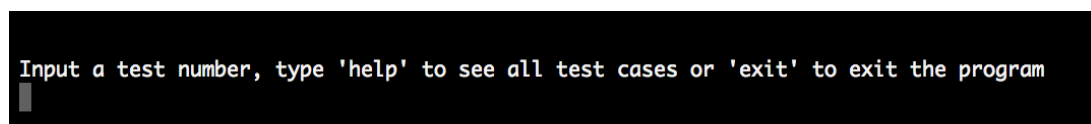
Figure 1: Initial prompt

If the user inputs help, the 32 possible test cases are shown to the user and the user is once again shown the initial prompt.

Figure 2: Help on screen information

An example of how a correct and incorrect test sentence is displayed can be seen in A.8:

# 4    Conclusion

The results of this software parser for regular expressions showed a 100% accuracy rate when tested against all thirty two test cases. The program and report give insight into how much detail is required to simply determine the validity of a sentence regardless of its context. The natural language toolkit for python is a great mechanism for understanding how to create corpus based systems for language processing. To conclude, the author feels that given more time, a more efficient algorithm could have been designed that takes advantage of the computational efficiency of a hash map and also its key value pairings which could be used to store left hand and right hand side arguments.

# References

[1] Siri apple voice assistant.

[2] Samina Dazdarevic, Amela Lukac-Zoranic, and ahreta Fijuljanin. Corpus-based research for language teaching. *Corpus Linguistics in Language Teaching.*

[3] Michael J. Garbade. A simple introduction to natural language processing, Oct 2018.

[4] Edward Loper and Steven Bird. *NLTK: the Natural Language Toolkit.* S.n., 2009.

[5] Lenhart Schubert. Computational linguistics, Feb 2014.

[6] Jiawei Yao. Automated sentiment analysis of text data with nltk. *Journal of Physics: Conference Series*, 1187(5):052020, 2019.

# A Appendix

## A.1 main.py

```
1 from parser import pos_tagger, grammar_check, bracketed_sentence,
      incorrect_sentence_parser, build_tree
2 from file_manager import get_tests, get_lex, get_rules, print_tests
3
4
5 def main():
6     rules = get_rules()
7     lex_arr = get_lex()
8     tests = get_tests()
9
10    while True:
11        x = input('\n\nInput a test number, type \'help\' to see all test cases
    or \'exit\' to exit the program\n')
12        if x in 'help':
13            print('----------------------------------------------------------')
14            print_tests()
15
16        elif x in 'exit':
17            return
18
19        elif x.isdigit():
20            print('\n\nRESULTS\n
    ----------------------------------------------------------')
21            x = int(x)
22            pos = pos_tagger(x,tests,lex_arr)
23
24            if grammar_check(pos, rules) == True:
25                print('Verdict: Sentence is grammatically correct\n')
26                bracketed_sentence(pos)
27                build_tree(pos)
28            else:
29                test = tests[x].split()
30                print('Verdit: Sentence is not grammatically correct\n')
31                incorrect_sentence_parser(test)
32        else:
33            print('\n----------------------------------------------------------
    ')
34            print('\nPLEASE INPUT A CORRECT COMMAND!\n')
35        print('----------------------------------------------------------')
36
37
38
39
40
41
42
43 if __name__== "__main__":
44   main()
```

## A.2 file_manager.py

```
1
2 def get_tests():
3     t=''
4     tests_file = open("tests.txt","r") #Open and read file that contains tests
    sentnces
5     tests = tests_file.readlines() #store each line as an individual test
    sentence
6     tests_file.close()
7     i = 0
8     return tests
9
```

```
10  def get_lex():
11
12      lex_file = open("lex.txt","r") #Open and read file that contains lexicon
13      lex_arr = lex_file.readlines()
14      lex_file.close()
15      return lex_arr
16
17  def get_rules():
18
19      rule_file = open("rules.txt","r") #Open and read file that contains rules
20      rule_arr = rule_file.readlines()
21      rule_file.close()
22      return rule_arr
23
24  def get_grammar():
25      file = open("grammar.txt","r") #Open and read file that contains rules
26      file_str = file.read()
27      grammar = CFG.fromstring(file_str) #CFG = Read File
28      file.close()
29      return grammar
30
31
32  def print_tests():
33      t=''
34      tests_file = open("tests.txt","r") #Open and read file that contains tests
            sentnces
35      tests = tests_file.readlines() #store each line as an individual test
            sentence
36      tests_file.close()
37      i = 0
38      for t in tests:
39          print(str(i) +': ' + t)
40          i += 1
```

## A.3  parser.py

```
1  from treelib import Node, Tree
2
3  #pos_tagger(x,tests,lex_arr) x = test number, tests = list of test sentences,
        lex_arr = list of lexicon rules
4  #This method compares each word of the test sentence against the lexicon file
        and allocates the related tag
5  def pos_tagger(x,tests,lex_arr):
6      sentence= ''
7      j = 0
8      i = 0
9      test = tests[x].split()
10      print('\nTest Sentence: ' + tests[x])
11      for word in test:
12          for lexi in lex_arr:
13              lex = lex_arr[j].split()
14              if test[i] == lex[1]:
15                  sentence += lex_arr[j]
16              j+= 1
17          j = 0
18          i += 1
19      return sentence
20
21
22  # grammar_check(pos,rules) is a boolean method that takes in the tagged pos
        sentence
23  # and a list of rules and compares the two and returns true if the sentence is
        correct
24  def grammar_check(pos, rules):
25      checker = ''
26      i = 1
27      pos_split = pos.split()
```

9

```
28      pos_NP = pos_split[0] + ' ' + pos_split[2] # gets the POS tag for the first
        two words and adds to string
29      pos_VP = pos_split[4] # gets the POS tag for the third word and adds to
        string
30      pos_NP2 = pos_split[6] + ' '+pos_split[8] + ' ' + pos_split[10] # gets the
        POS tag for the last three words and adds to string
31
32
33      if pos_split[1] in 'a' and pos_split[3] == 'men': # checks for 'a men' in a
        sentence
34          print('\'A men\' is grammatically incorrect.')
35          return False
36      elif pos_split[1] in 'a' and pos_split[3] == 'women': # checks for 'a women'
         in a sentence
37          print('\'A women\' is grammatically incorrect')
38          return False
39
40
41      while i<3: # Checking the if the childern of NP that has a parent of S are
        correct
42          rule = rules[i].split('->')
43          x = rule[1].split() # rule[1] is only split here again to fix the
        formatting and remove extra spacing
44          y = x[0]+' '+x[1] # here we assign formatted x to y
45          if pos_NP == y:
46              checker += rule[0] # rule[0] is the left hand side of the matching
        rule
47          i += 1
48      i=4
49      while i<6: # Checking if the childern of VP the has a parent of VP are
        correct
50          rule = rules[i].split('->')
51          x = rule[1].split()
52          y = x[0]
53          if pos_VP == y:
54              checker += rule[0]
55          i += 1
56
57      rule = rules[6].split('-> ') #checks if the childern of NP that has a parent
         of VP are correct
58      x = rule[1].split()
59      y = x[0]+' '+x[1]+' '+x[2]
60      if pos_NP2 == y:
61          checker += rule[0]
62      i += 1
63
64      str_split = checker.split()
65      if str_split[2] == 'NP': # all POS tags have been analysed and if the
        sentence is correct 'checker' should be NP VPzNP or NPp VPp NP
66          if str_split[0] == 'NP' and str_split[1] == 'VPz':
67              return True
68          elif str_split[0] == 'NPp' and str_split[1] == 'VPp':
69              return True
70          else:
71              return False
72      else:
73          return False
74
75
76 # bracketed_sentence(sentence) takes in a test sentence that has been proven to
        be correct and creates a bracketed phrasal structure
77 def bracketed_sentence(sentence):
78      i = 0
79      br_sentence = '[S1[S[NP['
80      split = sentence.split()
81      while i<3:
82          br_sentence += '[' + split[i] + ' ' + split[i+1] + ']'
83          i+=2
84      br_sentence += ']][VP[VP'
```

```
85      while i<5:
86          br_sentence += '[' + split[i] + ' ' + split[i+1] + ']'
87          i+=2
88      br_sentence += '][NP['
89      while i<11:
90          br_sentence += '[' + split[i] + ' ' + split[i+1] + ']'
91          i+=2
92      br_sentence += ']]]'
93      print(br_sentence)

94
95  # build_tree also takes a correct sentnce as an arg and outputs a tree to the
        cmd line
96  def build_tree(sentence):
97      print('\n\nParser Tree Plot\n')
98      split = sentence.split()
99      tree = Tree()
100     tree.create_node("S1", "s1")
101     tree.create_node("S", "s", parent="s1")
102     tree.create_node("NP", "np", parent="s")
103     tree.create_node("VP", "vp", parent="s")
104     tree.create_node("VP", "vp2", parent = "vp")
105     tree.create_node("_NP", "np2", parent = "vp")
106     tree.create_node(split[0],"dt" , parent="np")
107     tree.create_node(split[2], "nn", parent="np")
108     tree.create_node(split[4], "vb", parent="vp2")
109     tree.create_node(split[6], "dt2", parent="np2")
110     tree.create_node(split[8], "jj", parent="np2")
111     tree.create_node(split[10], "nn1", parent="np2")
112     tree.create_node(split[1],"1" , parent="dt")
113     tree.create_node(split[3], "2", parent="nn")
114     tree.create_node(split[5], "3", parent="vb")
115     tree.create_node(split[7], "4", parent="dt2")
116     tree.create_node(split[9], "5", parent="jj")
117     tree.create_node(split[11], "6", parent="nn1")
118     tree.show()

119

120

121
122 # incorrect_sentence_parser(sentence) takes an incorrect sentence and creates a
        bracketed phrasal structure for as long as the sentence is valid
123 def incorrect_sentence_parser(sentence):
124     br_str = '[S1[S[NP['
125     if repr(sentence[0]) == "'the'":
126         br_str += "DT the]"

127
128     elif repr(sentence[0]) == "'a'":
129         br_str += "DT a]"

130
131     else:
132         return

133
134     if repr(sentence[1]) in "'men'":
135         if repr(sentence[0]) in ("'a'"):
136             return
137         else:
138             x = repr(sentence[1])
139             br_str += "[NNP "+ x[1:-1] +"]"

140
141     elif repr(sentence[1]) in "'man'":
142         if repr(sentence[0]) in ("'a'"):
143             return
144         else:
145                 x = repr(sentence[1])
146                 br_str += "[NNP "+ x[1:-1]+"]"

147
148     else:
149         return

150
151     if repr(sentence[2]) == "'bites'":
```

```
152        if repr(sentence[1]) not in ("'man'" , "'woman'"):
153            print('Sentence incorrect but correct as far as ' + repr(sentence
    [1]))
154            br_str += "]]"
155            print('\n' + br_str)
156            return
157
158    elif repr(sentence[2]) == "'bite'":
159        if repr(sentence[1]) not in ("'men'" , "'women'"):
160            print('Sentence incorrect but correct as far as ' + repr(sentence
    [1]))
161            br_str += "]]]"
162            print('\n' + br_str)
163            return
164
165    elif repr(sentence[2]) == "'like'":
166        if repr(sentence[1]) not in ("'men'" , "'women'"):
167            print('Sentence incorrect but correct as far as ' + repr(sentence
    [1]))
168            br_str += "]]]"
169            print('\n' + br_str)
170            return
171
172    elif repr(sentence[2]) == "'likes'":
173        if repr(sentence[1]) not in ("'man'" , "'woman'"):
174            print('Sentence incorrect but correct as far as ' + repr(sentence
    [1]))
175            br_str += "]]]"
176            print('\n' + br_str)
177            return
178
179
180    br_str += "]]"
181    return
```

## A.4   Lexicon

DT the
DT a
NN man
NN woman
NNS men
NNS women
VBZ bites
VBP bite
VBZ likes
VBP like
JJ green
NN dog

## A.5   Rules

S -> NP VP
NP -> DT NN
NPp -> DT NNS
VP -> NP VP
VPp -> VBP
VPz -> VBZ

NP -> DT JJ NN

## A.6 Test Cases

the men bite the green dog
the men bites the green dog
the men like the green dog
the men likes the green dog
a men bite the green dog
a men bites the green dog
a men like the green dog
a men likes the green dog
the man like the green dog
the man likes the green dog
the man bite the green dog
the man bites the green dog
a man like the green dog
a man likes the green dog
a man bite the green dog
a man bites the green dog
the women bite the green dog
the women bites the green dog
the women like the green dog
the women likes the green dog
a women bite the green dog
a women bites the green dog
a women like the green dog
a women likes the green dog
the woman bite the green dog
the woman bites the green dog
the woman like the green dog
the woman likes the green dog
a women bite the green dog
a woman bites the green dog
a woman like the green dog
a woman likes the green dog

## A.7 Class Diagram

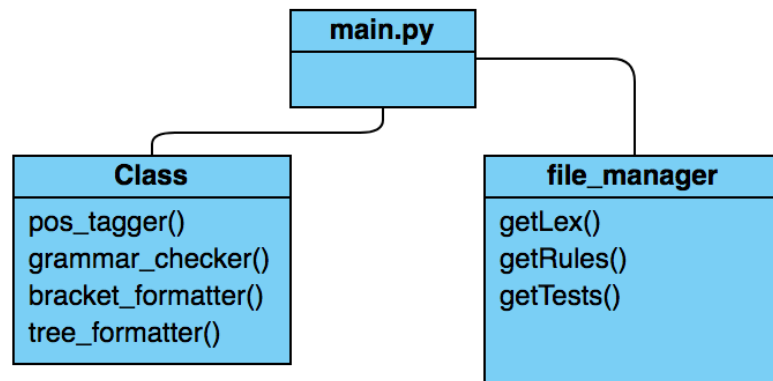The class diagram used for designing the software is shown below

Figure 3: Class diagram

## A.8   Outputs

```
RESULTS
-----------------------------------------------------------

Test Sentence: the men bite the green dog

Verdict: Sentence is grammatically correct

[S1[S[NP[[DT the][NNS men]]][VP[VP[VBP bite]][NP[[DT the][JJ green][NN dog]]]]]


Parser Tree Plot

S1
└── S
    ├── NP
    │   ├── DT
    │   │   └── the
    │   └── NNS
    │       └── men
    └── VP
        ├── VP
        │   └── VBP
        │       └── bite
        └── _NP
            ├── DT
            │   └── the
            ├── JJ
            │   └── green
            └── NN
                └── dog


-----------------------------------------------------------
```

Figure 4: Output displayed for correct test case

```
RESULTS
-----------------------------------------------------------

Test Sentence: the men bites the green dog

Verdict: Sentence is not grammatically correct

Sentence incorrect but correct as far as 'men'

[S1[S[NP[DT the][NNP men]]]]
-----------------------------------------------------------
```
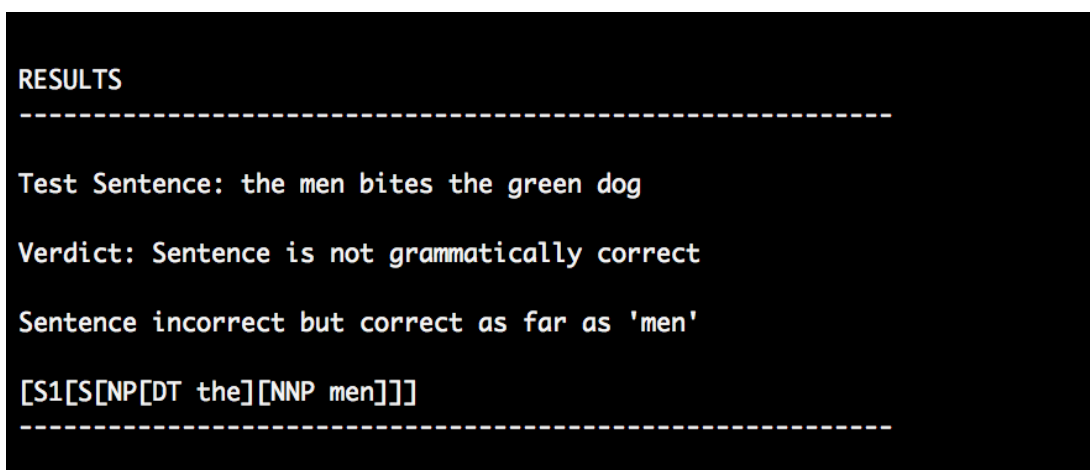
Figure 5: Output displayed for incorrect test case