# Japanese Scanner/Parser/Translator Project
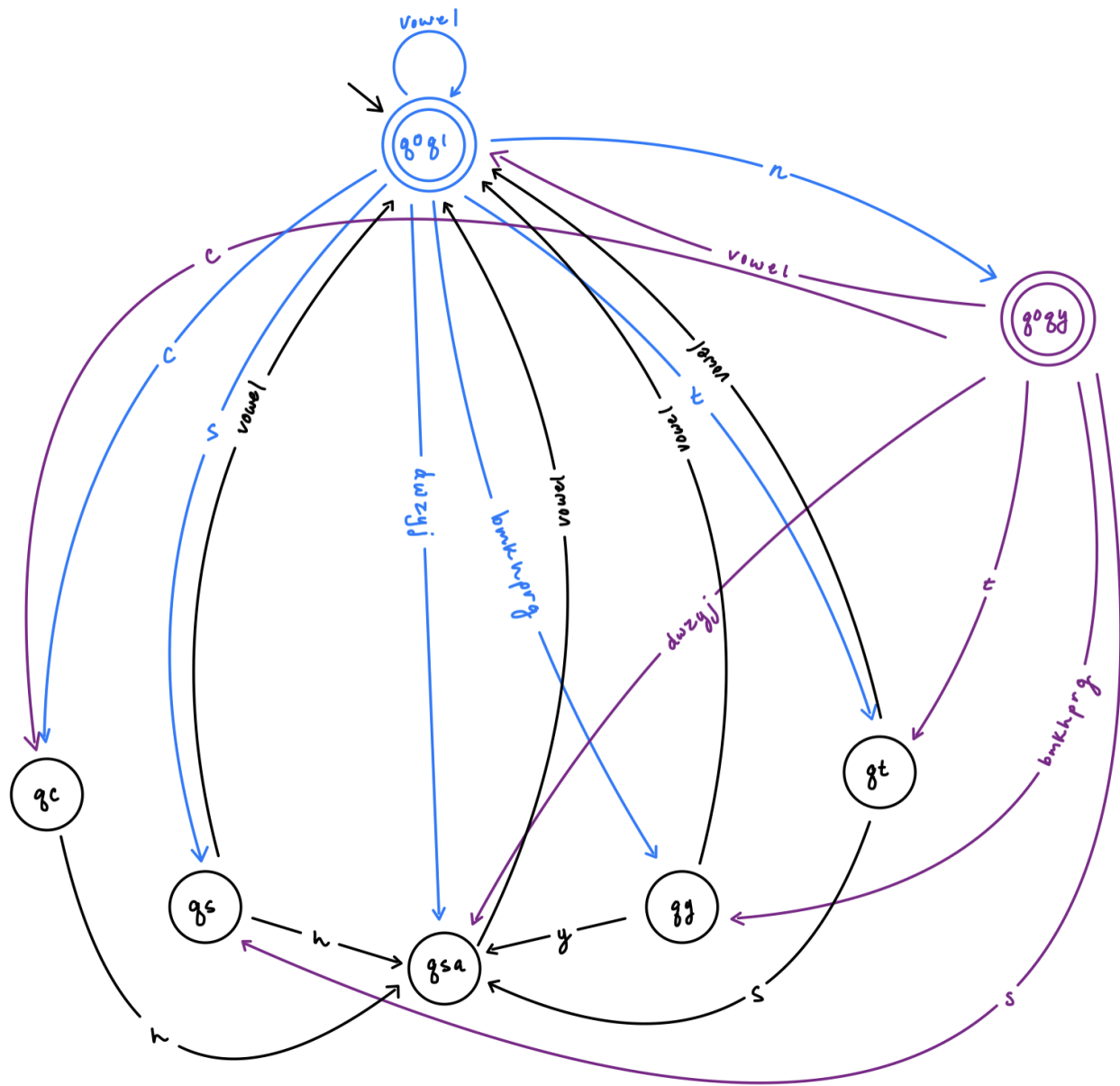
Group 23
Leouel Guanzon
Marco Flores
John Foster

# 1 – DFA

```cpp
1  #include<iostream>
2  #include<fstream>
3  #include<string>
4  using namespace std;
5
6  /* Look for all **'s and complete them */
7
8  //================================================
9  // File scanner.cpp written by: Group Number: 23
10 //================================================
11
12 //Done by: Leouel Guanzon and John Foster
13 //Tables are moved here so functions can access them.
14 enum tokentype {ERROR, WORD1, WORD2, PERIOD, VERB, VERBNEG, VERBPAST, VERBPASTNEG, IS, WAS, OBJECT, SUBJECT, DEST
   INATION, PRONOUN, CONNECTOR, EOFM};
15
16 string tokenName[16] = {"ERROR", "WORD1", "WORD2", "PERIOD", "VERB", "VERBNEG", "VERBPAST", "VERBPASTNEG", "IS",
   "WAS", "OBJECT", "SUBJECT", "DESTINATION", "PRONOUN", "CONNECTOR", "EOFM"};
17
18 //Array is used for simplicity
19 string reservedWords[38] = {"masu", "VERB", "masen", "VERBNEG", "mashita", "VERBPAST",
20 "masendeshita", "VERBPASTNEG", "desu", "IS", "deshita", "WAS",
21 "o", "OBJECT", "wa", "SUBJECT", "ni", "DESTINATION",
22 "watashi", "PRONOUN", "anata", "PRONOUN", "kare", "PRONOUN",
23 "kanojo", "PRONOUN", "sore", "PRONOUN", "mata", "CONNECTOR",
24 "soshite", "CONNECTOR", "shikashi", "CONNECTOR",
25 "dakara", "CONNECTOR", "eofm", "EOFM"};
26
27 // --------- Two DFAs --------------------------------
28 // WORD DFA
29 // Done by: Leouel Guanzon
30 // RE:
31 // (b|g|h|k|m|p|r) (y ((a|e|i|o|u|E|I)^+ | (a|e|i|o|u|E|I)^+ n) | (a|e|i|o|u|E|I)^+ | (a|e|i|o|u|E|I)^+ n)
32 // (a|e|i|o|u|E|I)^+ | (a|e|o|i|u|E|I)^+ n
33 // (d|j|w|y|z) ( (a|e|i|o|u|E|I)^+ | (a|e|i|o|u|E|I)^+ n)
34 // s (h((a|e|i|o|u|E|I)^+ | (a|e|i|o|u|E|I)^+ n)) | ((a|e|i|o|u|I|E)^+ | (a|e|i|o|u|I|E)^+ n))
35 // t (s((a|e|i|o|u|I|E)^+ | (a|e|i|o|u|I|E)^+ n)) | ((a|e|i|o|u|I|E)^+ | (a|e|i|o|u|I|E)^+ n)
36 // ch ((a|e|i|o|u|I|E)^+ | (a|e|i|o|u|I|E)^+ n)
37
38 //Done by: Marco Flores
39 bool isVowel(char c){
40   return (c=='a'|| c=='e'|| c=='i'|| c=='o'|| c=='u'|| c=='I'|| c=='E');
41 }
42
43 //Done by: Marco Flores
44 bool isConsonant1(char c){
45     return(c=='b'|| c=='g'|| c=='h'|| c=='k'|| c=='m'|| c=='p' || c=='r');
46 }
47
48 //Done by: Marco Flores
49 bool isConsonant2(char c){
50     return(c=='d'|| c=='j'|| c=='w'|| c=='y'|| c=='z');
51 }
52
53 //Done by: Leouel Guanzon and Marco Flores
54 bool word (string s)
55 {
56   int state = 0;
57   int charpos = 0;
58
59   /* replace the following todo the word dfa */
60   while (s[charpos] != '\0')
61     {
62       /* States:
63       * q0q1 = 0
64       * qsa = 1 = consonant
65      * qy = 2 = pair
66      * qs = 3 = s
67      * qt = 4 = t
68      * qc = 5 = c
69      * q0qy = 6 = q1
70      */
71
72      // q0q1 ==(d|j|w|y|z)==> qsa
73      if (state == 0 && isConsonant2(s[charpos]))
```

73,45-51          Top

guanz004@empress:~...

```
 74             state = 1;
 75         // q0q1 ==(b|g|h|k|m|p|r)==> qy
 76         else if (state == 0 && isConsonant1(s[charpos]))
 77             state = 2;
 78         // q0q1 ==s==> qs
 79         else if (state == 0 && s[charpos] == 's')
 80             state = 3;
 81         // q0q1 ==t==> qt
 82         else if (state == 0 && s[charpos] == 't')
 83             state = 4;
 84         // q0q1 ==c==> qc
 85         else if (state == 0 && s[charpos] == 'c')
 86             state = 5;
 87         // q0q1 ==n==> q0qy
 88         else if (state == 0 && s[charpos] == 'n')
 89             state = 6;
 90
 91         // (a|e|i|o|u|E|I) (a|e|i|o|u|E|I)^*
 92         // q0q1 ==(a|e|i|o|u|I|E)==> q0q1
 93         // qsa ==(a|e|i|o|u|I|E)==> q0q1
 94         // qy ==(a|e|i|o|u|I|E)==> q0q1
 95         // qs ==(a|e|i|o|u|I|E)==> q0q1
 96         // qt ==(a|e|i|o|u|I|E)==> q0q1
 97         // qc ==(a|e|i|o|u|I|E)==> q0q1
 98         // q0qy ==(a|e|i|o|u|I|E)==> q0q1
 99         else if ((state == 0||state ==  1||state == 2||state == 3||state == 4||state == 6) && isVowel(s[charpos]))
100             state = 0;
101
102         // pair followed by 'y'
103         // qy ==y==> qsa
104         else if (state == 2 && s[charpos] == 'y')
105             state = 1;
106         // from state 3 || 5
107         // followed by 'h'
108         // qs ==h==> qsa || qc ==h==> qsa
109         else if ((state == 3 || state == 5) && s[charpos] == 'h')
110             state = 1;
111         // from state 4
112         // followed by 's'
113         // qt ==s==> qt
114         else if (state == 4 && s[charpos] == 's')
115             state = 1;
116
117         // from sate 5
118         // followed by 'h'
119         // qc ==h==> qsa
120         //else if (state == 5 && s[charpos] == 'h')
121         //   state = 1;
122
123         // q0qy ==(d|j|w|y|z)==> qsa
124         else if (state == 6 && isConsonant2(s[charpos]))
125             state = 1;
126         // q0qy ==(b|g|h|k|m|p|r)==> qy
127         else if (state == 6 && isConsonant1(s[charpos]))
128             state = 2;
129         // q0qy ==s==> qs
130         else if (state == 6 && s[charpos] == 's')
131             state = 3;
132         // q0qy ==t==> qt
133         else if (state == 6 && s[charpos] == 't')
134             state = 4;
135         // q0qy ==c==> qc
136         else if (state == 6 && s[charpos] == 'c')
137             state = 5;
138
139         else
140             return ERROR;
141         charpos++;
142     }//end of while
143
144      // where did I end up????
145     if (state == 0)
146     {
147         return WORD1; //scanner() function will overwrite to WORD2 if string ends in 'I' or 'E'
                                                              147,89-95        30%
```

guanz004@empress:~...

```
148        }
149        else if (state == 6)  // end in a final state q0 (0) or q0' (6)
150        {
151            return WORD1;
152        }
153        else
154            return ERROR;
155 }
156
157 // PERIOD DFA
158 // Done by: Leouel Guanzon and Marco Flores
159 bool period (string s)
160 {
161     int state = 0;
162     int charpos = 0;
163
164     while(s[charpos] != '\0'){
165         if(s[charpos] == '.' && s[charpos + 1] == '\0'){
166             return PERIOD;
167         }
168         charpos++;
169     }
170
171     return ERROR;
172 }
173
174 // ------ Three  Tables ------------------------------------
175
176 // TABLES Done by: Leouel Guanzon and John Foster
177 // Moved to the top to be use for global scope
178
179 /*
180 // ** Update the tokentype to be WORD1, WORD2, PERIOD, ERROR, EOFM, etc.
181 enum tokentype {ERROR, WORD1, WORD2, PERIOD, VERB, VERBNEG, VERBPAST, VERBPASTNEG, IS, WAS, OBJECT, SUBJECT, DEST
    INATION, PRONOUN, CONNECTOR, EOFM};
182
183 // ** For the display names of tokens - must be in the same order as the tokentype.
184 string tokenName[16] = {"ERROR", "WORD1", "WORD2", "PERIOD", "VERB", "VERBNEG", "VERBPAST", "VERBPASTNEG", "IS",
    "WAS", "OBJECT", "SUBJECT", "DESTINATION", "PRONOUN", "CONNECTOR", "EOFM"};
185
186 string reservedWords[38] = {"masu", "VERB", "masen", "VERBNEG", "mashita", "VERBPAST",
187 "masendeshita", "VERBPASTNEG", "desu", "IS", "deshita", "WAS",
188 "o", "OBJECT", "wa", "SUBJECT", "ni", "DESTINATION",
189 "watashi", "PRONOUN", "anata", "PRONOUN", "kare", "PRONOUN",
190 "kanojo", "PRONOUN", "sore", "PRONOUN", "mata", "CONNECTOR",
191 "soshite", "CONNECTOR", "shikashi", "CONNECTOR",
192 "dakara", "CONNECTOR", "eofm", "EOFM"};
193 */
194
195 // ** Need the reservedwords table to be set up here.
196 // ** Do not require any file input for this. Hard code the table.
197 // ** a.out should work without any additional files.
198
199
200 // ------------ Scanner and Driver ----------------------
201
202 ifstream fin;  // global stream for reading from the input file
203
204 // Scanner processes only one word each time it is called
205 // Gives back the token type and the word itself
206 // ** Done by: Leouel Guanzon and John Foster
207 int scanner(tokentype& tt, string& w)
208 {
209    // ** Grab the next word from the file via fin
210    // 1. If it is eofm, return right now.
211    fin >> w;
212
213    if(w == "eofm")
214    {
215        return EOFM;
216    }
217
218    /*  **
219    2. Call the token functions (word and period)
220       one after another (if-then-else).
```

```
221          Generate a lexical error message if both DFAs failed.
222          Let the tokentype be ERROR in that case.
223     */
224     if(word(w)){
225         if(w[w.length()-1] == 'I' || w[w.length()-1] == 'E')
226         {
227             tt = WORD2;
228         } else
229         {
230             tt = WORD1;
231         }
232
233     /*
234     3. If it was a word,
235        check against the reservedwords list.
236        If not reserved, tokentype is WORD1 or WORD2
237        decided based on the last character.
238     */
239
240
241     /*
242     4. Return the token type & string   (pass by reference)
243     */
244        for(int i = 0; i < 38; i++)
245        {
246            if(reservedWords[i] == w)
247            {
248                for(int j = 0; j <= 16; j++)
249                {
250                    if(tokenName[j] == reservedWords[i+1])
251                    {
252                        tt = static_cast<tokentype>(j);
253                        break;
254                    }
255                }
256            }
257        }
258     }
259     else if(period(w))
260     {
261         tt = PERIOD;
262     }
263     else
264     {
265         tt = ERROR;
266     }
267
268     if(tt == ERROR)
269     {
270         cout << "Lexical error: " << w << " is not a valid token." << endl;
271     }
272
273     return 1;
274
275 }//the end of scanner
276
277
278
279 // The temporary test driver to just call the scanner repeatedly
280 // This will go away after this assignment
281 // DO NOT CHANGE THIS!!!!!!
282 // Done by:  Louis
283 int main()
284 {
285    tokentype thetype;
286    string theword;
287    string filename;
288
289    cout << "Enter the input file name: ";
290    cin >> filename;
291
292    fin.open(filename.c_str());
293
294    // the loop continues until eofm is returned.
295     while (true)
```
                                                          295,15              93%

```
294    // the loop continues until eofm is returned.
295      while (true)
296        {
297          scanner(thetype, theword);  // call the scanner which sets
298                                     // the arguments
299          if (theword == "eofm") break;  // stop now
300
301          cout << "Type is: " << tokenName[thetype] << endl;
302          cout << "Word is: " << theword << endl;
303
304          cout << endl;
305        }
306
307      cout << "End of file is encountered." << endl;
308      fin.close();
309
310 }// end
```

                                                    310,7        Bot

3 – Original Scanner test results



```
guanz004@empress:~/CS4  /CS421Progs/ScannerFiles

[guanz004@empress ScannerFiles]$ script scannerOutput.txt
Script started, file is scannerOutput.txt
[guanz004@empress ScannerFiles]$ g++ scanner.cpp -o scanner.out
[guanz004@empress ScannerFiles]$ ./scanner.out
Enter the input file name: scannertest1
Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD1
Word is: rika

Type is: IS
Word is: desu

Type is: PERIOD
Word is: .

Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD1
Word is: sensei

Type is: IS
Word is: desu

Type is: PERIOD
Word is: .

Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD1
Word is: ryouri

Type is: OBJECT
Word is: o

Type is: WORD2
Word is: yarI

Type is: VERB
Word is: masu

Type is: PERIOD
Word is: .

Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD1
Word is: gohan

Type is: OBJECT
Word is: o

Type is: WORD1
Word is: seito

Type is: DESTINATION
Word is: ni

Type is: WORD2
Word is: agE
```

guanz004@empress:~...

```
Type is: VERBPAST
Word is: mashita

Type is: PERIOD
Word is: .

Type is: CONNECTOR
Word is: shikashi

Type is: WORD1
Word is: seito

Type is: SUBJECT
Word is: wa

Type is: WORD2
Word is: yorokobI

Type is: VERBPASTNEG
Word is: masendeshita

Type is: PERIOD
Word is: .

Type is: CONNECTOR
Word is: dakara

Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD1
Word is: kanashii

Type is: WAS
Word is: deshita

Type is: PERIOD
Word is: .

Type is: CONNECTOR
Word is: soshite

Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD1
Word is: toire

Type is: DESTINATION
Word is: ni

Type is: WORD2
Word is: ikI

Type is: VERBPAST
Word is: mashita

Type is: PERIOD
Word is: .

Type is: PRONOUN
Word is: watashi

Type is: SUBJECT
Word is: wa

Type is: WORD2
Word is: nakI

Type is: VERBPAST
```

Word is: mashita

Type is: PERIOD
Word is: .

End of file is encountered.
[guanz004@empress ScannerFiles]$ ./scanner.out
Enter the input file name: scannertest2
Type is: WORD1
Word is: daigaku

Lexical error: college is not a valid token.
Type is: ERROR
Word is: college

Type is: WORD1
Word is: kurasu

Lexical error: class is not a valid token.
Type is: ERROR
Word is: class

Type is: WORD1
Word is: hon

Lexical error: book is not a valid token.
Type is: ERROR
Word is: book

Type is: WORD1
Word is: tesuto

Lexical error: test is not a valid token.
Type is: ERROR
Word is: test

Type is: WORD1
Word is: ie

Lexical error: home* is not a valid token.
Type is: ERROR
Word is: home*

Type is: WORD1
Word is: isu

Lexical error: chair is not a valid token.
Type is: ERROR
Word is: chair

Type is: WORD1
Word is: seito

Lexical error: student is not a valid token.
Type is: ERROR
Word is: student

Type is: WORD1
Word is: sensei

Lexical error: teacher is not a valid token.
Type is: ERROR
Word is: teacher

Type is: WORD1
Word is: tomodachi

Lexical error: friend is not a valid token.
Type is: ERROR
Word is: friend

Type is: WORD1
Word is: jidoosha

Lexical error: car is not a valid token.
Type is: ERROR

guanz004@empress:~...

```
Word is: car

Type is: WORD1
Word is: gyuunyuu

Lexical error: milk is not a valid token.
Type is: ERROR
Word is: milk

Type is: WORD1
Word is: sukiyaki

Type is: WORD1
Word is: tenpura

Type is: WORD1
Word is: sushi

Type is: WORD1
Word is: biiru

Lexical error: beer is not a valid token.
Type is: ERROR
Word is: beer

Type is: WORD1
Word is: sake

Type is: WORD1
Word is: tokyo

Type is: WORD1
Word is: kyuushuu

Lexical error: Osaka is not a valid token.
Type is: ERROR
Word is: Osaka

Type is: WORD1
Word is: choucho

Lexical error: butterfly is not a valid token.
Type is: ERROR
Word is: butterfly

Type is: WORD1
Word is: an

Type is: WORD1
Word is: idea

Type is: WORD1
Word is: yasashii

Lexical error: easy is not a valid token.
Type is: ERROR
Word is: easy

Type is: WORD1
Word is: muzukashii

Lexical error: difficult is not a valid token.
Type is: ERROR
Word is: difficult

Type is: WORD1
Word is: ureshii

Lexical error: pleased is not a valid token.
Type is: ERROR
Word is: pleased

Type is: WORD1
Word is: shiawase

Lexical error: happy is not a valid token.
```

```
Type is: ERROR
Word is: happy

Type is: WORD1
Word is: kanashii

Lexical error: sad is not a valid token.
Type is: ERROR
Word is: sad

Type is: WORD1
Word is: omoi

Lexical error: heavy is not a valid token.
Type is: ERROR
Word is: heavy

Type is: WORD1
Word is: oishii

Lexical error: delicious is not a valid token.
Type is: ERROR
Word is: delicious

Lexical error: tennen is not a valid token.
Type is: ERROR
Word is: tennen

Lexical error: natural is not a valid token.
Type is: ERROR
Word is: natural

Type is: WORD2
Word is: nakI

Lexical error: cry is not a valid token.
Type is: ERROR
Word is: cry

Type is: WORD2
Word is: ikI

Lexical error: go* is not a valid token.
Type is: ERROR
Word is: go*

Type is: WORD2
Word is: tabE

Lexical error: eat is not a valid token.
Type is: ERROR
Word is: eat

Type is: WORD2
Word is: ukE

Lexical error: take* is not a valid token.
Type is: ERROR
Word is: take*

Type is: WORD2
Word is: kakI

Lexical error: write is not a valid token.
Type is: ERROR
Word is: write

Type is: WORD2
Word is: yomI

Lexical error: read is not a valid token.
Type is: ERROR
Word is: read

Type is: WORD2
Word is: nomI
```

```
Lexical error: drink is not a valid token.
Type is: ERROR
Word is: drink

Type is: WORD2
Word is: agE

Lexical error: give is not a valid token.
Type is: ERROR
Word is: give

Type is: WORD2
Word is: moraI

Lexical error: receive is not a valid token.
Type is: ERROR
Word is: receive

Type is: WORD2
Word is: butsI

Lexical error: hit is not a valid token.
Type is: ERROR
Word is: hit

Type is: WORD2
Word is: kerI

Lexical error: kick is not a valid token.
Type is: ERROR
Word is: kick

Type is: WORD2
Word is: shaberI

Lexical error: talk is not a valid token.
Type is: ERROR
Word is: talk

End of file is encountered.
[guanz004@empress ScannerFiles]$ exit
exit
Script done, file is scannerOutput.txt
[guanz004@empress ScannerFiles]$
```

guanz004@empress:~...

## 4 – Factored Rules

1 <story> ::= <s> { <s> } // stay in the loop as long as a possible start
                  // of <s> is the next_token  (note it can be CONNECTOR or WORD1 or PRONOUN)


2 <s> ::= [CONNECTOR] <noun> SUBJECT  <verb> <tense> PERIOD
3 <s> ::= [CONNECTOR] <noun> SUBJECT  <noun> <be>    PERIOD
4 <s> ::= [CONNECTOR] <noun> SUBJECT  <noun> DESTINATION  <verb> <tense> PERIOD
5 <s> ::= [CONNECTOR] <noun> SUBJECT  <noun> OBJECT <verb> <tense> PERIOD
6 <s> ::= [CONNECTOR] <noun> SUBJECT  <noun> OBJECT <noun> DESTINATION <verb> <tense> PERIOD


  **// Refer to Left Factoring file to make these into**
  // one rule until things start to differ.

**<s> ::= [CONNECTOR #getEword# #gen(CONNECTOR)#]**
               **<noun> #getEword# SUBJECT #gen(ACTOR)#**
               **<after subject>**

**<after subject> ::= <verb> #getEword# #gen(ACTION)#**
               **<tense> #gen(TENSE)# PERIOD |**
               **<noun> #getEword#**
               **<after noun>**

**<after noun> ::= <be> #gen(DESCRIPTION)# #gen(TENSE)# PERIOD |**
             **DESTINATION #gen(TO)#**
             **<verb> #getEword# #gen(ACTION)#**
             **<tense> #gen(TENSE)# PERIOD |**
             **OBJECT #gen(OBJECT)#**
             **<after object>**

**<after object> ::= <verb> #getEword# #gen(ACTION)#**
              **<tense> #gen(TENSE)# PERIOD |**
              **<noun> #getEword# DESTINATION #gen(TO)#**
              **<verb> #getEword# #gen(ACTION)#**
              **<tense> #gen(TENSE)# PERIOD**

7 <noun> ::= WORD1 | PRONOUN

8 <verb> ::= WORD2

9 <be> ::=  IS | WAS

10 <tense> := VERBPAST  | VERBPASTNEG | VERB | VERBNEG

5 – Updated Parser code for Translation