```csharp
 1  using System;
 2  using System.Collections;
 3  using System.Collections.Generic;
 4  using Microsoft.SqlServer.Server;
 5  using System.Data.SqlTypes;
 6  using System.IO;
 7  using System.Text;
 8  using System.Text.RegularExpressions;
 9
10  namespace CLRUDF
11  {
12      public class SpeechParser
13      {
14          public static readonly string DEBUG_ROOT = "C:\\Projects\\CIS 612 Lab2\
                \";
15          public static readonly string LOG_FILENAME = "ErrorLog.txt";
16          public static readonly string TSV_FILENAME =
                "Speech_Inverted_Index.tsv";
17          public static readonly string CONNECTION_STRING = "Data
                Source=DESKTOP-335I8BU;Initial Catalog=SPEECHES;Integrated
                Security=True";
18          public static readonly string STOP_FILE = "stop.txt";
19          public static readonly char[] JUNK_CHARS = {
20              '`', '~', '!', '@', '#', '$', '%', '^', '&', '*',
21              '(', ')', '_', '=', '+', '[', ']', '{', '}',
22              '\\', '|', '\'', '\"', ';', ':', ',', '<', '.', '>', '/', '?'};
23
24          public static List<Speech> GetSpeeches(string filename)
25          {
26              // This will be written in a very, very file specific manner, if
                  the format of the file ever
27              //  changes, this code WILL have to change with it
28              string[] split = { "NEW ADDRESS=>" , "FIRSTNAME=", "LASTNAME=",
                  "MONTH=", "DAY=", "YEAR=", "WEBLINK=",
29                  "FILELINK=", "SPEECH=" };
30              List<Speech> speeches = new List<Speech>();
31              try
32              {
33                  string line = null;
34                  StreamReader reader = new StreamReader(filename);
35                  while (!reader.EndOfStream)
36                  {
37                      line = reader.ReadLine();
38                      string[] parts = line.Split(split,
                          StringSplitOptions.None);
39                      int day = -1, year = -1;
40                      for (int i = 1; i < parts.Length; i += 9)
41                      {
42                          // Give a default value to day/year if they are bad
```

```csharp
43                      if (!int.TryParse(parts[i + 4], out day))
44                      {
45                          day = -1;
46                      }
47                      if (!int.TryParse(parts[i + 5], out year))
48                      {
49                          year = -1;
50                      }
51                      // Skip the first slot, its added by split and contains ⮡
                            nothing
52                      speeches.Add(new Speech
53                      {
54                          Firstname = parts[i + 1],
55                          Lastname = parts[i + 2],
56                          Month = parts[i + 3],
57                          Day = day,
58                          Year = year,
59                          Weblink = parts[i + 6],
60                          Filelink = parts[i + 7],
61                          Text = parts[i + 8]
62                      });
63                  }
64              }
65              reader.Close();
66          }
67          catch (Exception ex)
68          {
69              Console.WriteLine(ex.ToString());
70          }
71          return speeches;
72      }
73
74      public static Dictionary<string, int> CreateSpeechIndex(List<Speech>    ⮡
           speeches)
75      {
76          // Stores the number of times each term comes up across all        ⮡
               speeches
77          Dictionary<string, int> termFrequency = new Dictionary<string, int> ⮡
               ();
78
79          // Read in the stop words
80          HashSet<string> stopWords = new HashSet<string>();
81          try
82          {
83              StreamReader reader = new StreamReader(DEBUG_ROOT + STOP_FILE);
84              while (!reader.EndOfStream)
85              {
86                  stopWords.Add(reader.ReadLine());
87              }
```

```
88                  reader.Close();
89              }
90          catch (Exception ex)
91          {
92              throw ex;
93          }
94
95          // holds the current speech
96          StringBuilder sb = new StringBuilder();
97          string currentSpeech = "";
98
99          // Iterate through each speech
100         foreach (var speech in speeches)
101         {
102             // Remove characters we do not care about
103             foreach (var c in speech.Text.ToLower())
104             {
105                 if (!Array.Exists(JUNK_CHARS, j => j == c))
106                 {
107                     sb.Append(c);
108                 }
109             }
110
111             currentSpeech = sb.ToString();
112
113             // Remove the stop words
114             foreach (var word in stopWords)
115             {
116                 // Not the best way to do this
117                 currentSpeech = currentSpeech.SafeReplace(word, "", true);
118             }
119
120             // Remove the excess spaces in the string
121             RegexOptions options = RegexOptions.None;
122             Regex regex = new Regex("[ ]{2,}", options);
123             currentSpeech = regex.Replace(currentSpeech, " ");
124
125             // Split the string into the query space
126             string[] words = currentSpeech.Split(' ');
127
128             // Update the index
129             foreach (var word in words)
130             {
131                 // update the term frequency
132                 if (termFrequency.ContainsKey(word))
133                 {
134                     termFrequency[word]++;
135                 }
136                 else
```

```
137                    {
138                        termFrequency[word] = 1;
139                    }
140
141                    // Update the document frequency
142                    speech.UpdateFrequency(word);
143                }
144
145                // Clear out the current speech
146                currentSpeech = "";
147                // Clear out the string builder
148                sb.Clear();
149            }
150
151            return termFrequency;
152        }
153
154        [SqlFunction(FillRowMethodName = "FillRow")]
155        public static IEnumerable InitMethod(string filename)
156        {
157            // Attempt to read in the speeches from file
158            List<Speech> speeches = GetSpeeches(filename);
159            // Hold the inverted index
160            List<Tuple<string, int>> termFrequency = new List<Tuple<string,
                 int>>();
161
162            try
163            {
164                // CreateSpeechIndex returns the term frequency
165                // Document frequency can be found as an attribute of the
                     Speech class and will be
166                //  up-to-date when CreateSpeechIndex() exits
167                foreach (var term in CreateSpeechIndex(speeches))
168                {
169                    termFrequency.Add(new Tuple<string, int>(term.Key,
                         term.Value));
170                }
171
172                // TODO: Write out the speeches to the database
173                /*foreach (var speech in speeches)
174                {
175
176                }*/
177
178                // Write the term frequency out to a tsv file
179                StreamWriter output = new StreamWriter(File.OpenWrite
                     (DEBUG_ROOT + TSV_FILENAME));
180                foreach (var term in termFrequency)
181                {
```

```
182                        output.WriteLine($"{term.Item1}\t{term.Item2}");
183                    }
184                output.Flush();
185                output.Close();
186            }
187            catch (Exception ex)
188            {
189                StreamWriter log = new StreamWriter(File.OpenWrite(DEBUG_ROOT + ⇗
                       LOG_FILENAME));
190                log.WriteLine(ex.ToString());
191                log.Flush();
192                log.Close();
193            }
194
195            // Return the speeches to the database, so it can handle updating  ⇗
                  the database
196            return termFrequency;
197        }
198
199        public static void FillRow(Object obj, out SqlString term, out SqlInt32 ⇗
               frequency)
200        {
201            Tuple<string, int> pair = obj as Tuple<string, int>;
202            term = new SqlString(pair.Item1);
203            frequency =  new SqlInt32(pair.Item2);
204        }
205    }
206 }
207
```