

# Report of Assignment 7

Haodong Liao & Liyuan Zhang

July 12, 2019

## 1 Preface

This was a  $\text{\LaTeX}$  report written by Haodong Liao and Liyuan Zhang, students of UESTC who participated in the summer school of UCPH. More codes of this summer school were at this [repository](#).

## 2 Introduction

A string in F# was a sequence of characters<sup>1</sup>, and some of the string functions were similar to list functions, after the practicing of the list function, it was time to move on.

In this assignment, we worked with simple text processing and met the requirements of analyzing and generating some text according to statistics. In detail, we used pattern-matching and explicit type definition to create our own recursive functions, besides, we used powerful List-library functions to read, convert, analyze and generate target text.

## 3 Analysis and Design

### 3.1 Text Processing of character

#### 3.1.1 ReadText

```
1 readText: filename:string -> string
2 // readFile.fsx
3 let filename = "readFile.fsx"
4 let line =
5     try
6         let reader = System.IO.File.OpenText filename
7         reader.ReadToEnd ()
8     with
9         _ -> "" // The file cannot be read, so we return an empty string
10 printfn "%s" line
```

In this sub-assignment, we needed to convert *readFile.fsx* into a more flexible function which reads the content of any text file.

Referred to the input processing code in assignment 6, there were two branches needed to be handled, one was to set the entry point of our code, the other was to change the fixed argument to a flexible one so that the function could read the content of any text file. The key part of the pseudocode was shown as following:

```
1 [<EntryPoint>]
2 let main (input argument) =
3     if (argument correct)
4         call readText function with responding name of file
5     else
6         display error message
```

#### 3.1.2 ConvertText

```
1 convertText: src:string -> string
```

---

<sup>1</sup><https://docs.microsoft.com/en-us/dotnet/fsharp/language-reference/lists>

The goal of this sub-assignment was to convert letters of a string to lower case and removes all characters except a..z. According to [1], there was a *ToLower()* function *returns a copy of the string where each letter has been converted to lower case*, so we could use it to finish the 'convert' part and moved on to 'remove' part. The key part of the pseudocode is shown as following:

```
1 let convertText inputString =
2   let lowerStr = inputString.ToLower()
3   if (element of lowerStr was in the range of a..z)
4     keep the current element
5   else remove the current element
```

### 3.1.3 Get Histogram

```
1 histogram: src:string -> int list
```

We needed to count occurrences of each lower-case letter of the English alphabet in a string and returned a list in this sub-assignment. According to [1], there was a *String.filter* function similar to *List.filter* function, we used it to filter the specific letter of the string and count the length of it to get the occurrences. The key part of the pseudocode was shown as following:

```
1 let countchar (src:string)(letter:char) =
2   if letter was in the range of a..z then
3     filter the current letter
4     handle the rest string
5   else
6     returned an empty list
7 let histogram (src:string): int list = (countchar src 'a')
```

### 3.1.4 Generate Random String

```
1 randomString: hist: int list -> len:int -> string
```

The requirement of this sub-assignment was to change the given program and generate a string of a given length and contained random characters distributed according to a given histogram which used our own functions. The flow of the program was shown as following:

```
1 step1 - Read the target text file and saved it
2 step2 - Converted the text to lower case & removed characters except a..z
3 step3 - Got the histogram of the text
4 step4 - Generated a new random string according to the histogram
```

## 3.2 Text Processing of pair of characters

### 3.2.1 Get Cooccurrence

In this sub-assignment, we needed to count occurrences of each pair of lower-case letter of the English alphabet in a string and return a list of lists. Our solution was to use double recursion and process the letter in each pair separately. The key part of the pseudocode was shown as following:

```
1 // count the occurrence of each pair
2 let countnum =
3   if the length of src > 2 then
4     if the first two char = charpair then
5       1 + sum recursively
6     else
7       ignore current and sum recursively
8   else
9     return 0
10 // handle the second letter
11 let countcharpair =
12   if the second letter <= 'z' then
13     sum the occurrence and process the left string recursively
```

```

14     else return an empty list
15     // process the first letter
16     let countchar =
17         if the first letter <= 'z' then
18             process the second letter and process the left string recursively
19         else return an empty list
20     // call function from here
21     let cooccurrence (src:string) = countchar src 'a'

```

### 3.2.2 Markov Model

```

1  fstOrderMarkovModel: cooc: int list list -> len:int -> string

```

We needed to generate a random string of length *len*, whose character pairs were distributed according to a specified cooccurrence histogram *cooc* of the story "Little Claus and Big Claus" in this sub-assignment. The flow of the program was shown as following:

```

1  step1 - Count the occurrence of each pairs (like 'aa' to 'az')
2  step2 - Count the cooccurrence of all pairs ('aa' to 'zz')
3  step3 - Calculate the ratio of each pairs
4  step4 - Generate a random string of length len
5          according to a given cooccurrence and the ratio of each pairs

```

## 4 Program description

Our implementation of **readText** function was as follows:

```

1  let printErrorMessage () =
2      printfn "Program Input should be 'readfile filename'"
3
4  let readText (filename:string) : string =
5      let line =
6          try
7              let reader = System.IO.File.OpenText filename
8              reader.ReadToEnd ()
9          with
10             _ -> "" // The file cannot be read, so we return an empty string
11      printfn "%s" line
12      line
13  [<EntryPoint>]
14  let main (paramList : string []) : int =
15      if paramList.Length <> 2 then
16          printErrorMessage ()
17          0
18      else
19          match paramList.[0] with
20          | "readfile" ->
21              let str = readText paramList.[1]
22              printfn "str : %A" str
23          | _ ->
24              printErrorMessage ()
25      1

```

As it showed, the key to this function was to match the argument list, the first argument should be the command 'readfile', and the second argument was the name of a file.

Implementation of **convertText** function was as follows:

```

1  let rec convertText (src:string) : string =
2      let lowerSRC = src.ToLower()
3      match lowerSRC with
4      | "" -> ""
5      | elm ->

```

```

6         if ('a' <= elm.[0] && elm.[0] <= 'z')
7         then
8             elm.[0].ToString() + (convertText elm.[1..])
9         else
10            convertText elm.[1..]

```

It made us wondering whether there was an expression of string could work as 'elm::rest' to represent a list at the beginning, things did not go well so we changed our mind to use elm[1..] to represent the rest of a string.

Implementation of **histogram** function was as follows:

```

1 let rec countchar (src:string)(letter:char): int list =
2     if letter >= 'a' && letter <= 'z' then
3         ((String.filter (fun char -> char=letter) src).Length)::
4         (countchar src (char (int letter + 1)))
5     else []
6 let histogram (src:string): int list = (countchar src 'a')

```

This function was the key to solving the following problems.

**Assignment 7.4** was a summary of the previous assignments, and the key part to it was as follows:

```

1  /// step1 readfile
2  let str = readText "littleClausAndBigClaus.txt"
3  /// step2 convertfile
4  let convertRes = convertText str
5  /// step3 get histogram
6  let hist = histogram convertRes
7  let alphabet = List.init hist.Length (fun i -> 'a' + char i)
8  printfn "A histogram:\n %A" (List.zip alphabet hist)
9  /// step4 generages a new random string
10 let ranStr = randomString hist convertRes.Length
11 printfn "A random string: %s" ranStr
12 let newHist = histogram ranStr
13 printfn "Resulting histogram:\n %A" (List.zip alphabet newHist)

```

It was important to have a clear thought of processing flow.

Implementation of **cooccurrence** function was as follows:

```

1 let rec countnum (src:string) (charpair:string): int =
2     if src.Length >= 2 then
3         if src.[0..1] = charpair then
4             1+(countnum src.[1..] charpair)
5         else
6             (countnum src.[1..] charpair)
7     else
8         0
9 let rec countcharpair (src:string) (letter1:char) (letter2:char): int list =
10     if letter2 <= 'z' then
11         let charpair = string letter1 + string letter2
12         (countnum src charpair)::(countcharpair src letter1 (letter2+char 1))
13     else []
14 let rec countchar (src:string) (letter1:char): int list list =
15     if letter1 <= 'z' then
16         (countcharpair src letter1 'a')::(countchar src (letter1 + char 1))
17     else []
18 let cooccurrence (src:string): int list list =
19     countchar src 'a'

```

This was a difficult problem for us, and the key to the solution was the idea of 'divide and conquer'. If we could not process the pair of letters at a time, we handle the letter of it separately.

**Assignment 7.6** was a summary of all the previous assignments, and the key part of the implementation of it was as follows:

```

1 let fstOrderMarkovModel (cooc:int list list) (len:int): string =
2     let strLengthlist =
3         List.map (fun x -> (float x)) (List.map (List.sum) cooc)

```

```

4     let strLengthSum = List.sum strLengthlist
5     let ratiolist = List.map (fun x -> x/strLengthSum) strLengthlist
6     let randomStringAlt (hist: int list): string = randomString hist
7         (int((float (List.sum (hist))/strLengthSum)*(float len)))
8     let strlist = List.map (randomStringAlt) cooc
9     let composestr (acc: string) (elm: string): string = acc + elm
10    List.fold composestr "" strlist
11    printfn "test: \n %A" (fstOrderMarkovModel (cooccurrence a) (a.Length))
12    printfn "present cooc : \n %A"
13    (cooccurrence (fstOrderMarkovModel (cooccurrence a) (a.Length)))

```

This was another difficult assignment for us, the main challenge was to arrange the text properly so the occurrence was similar. We tried to generate a new character according to the letter before the current letter, but things just did not work well for reason we did not figure out, then we tried with the weight of the pair based on the proportion of the sum of the single letters to the sum of the occurrences of all the letters. In this way, when there was a given length, the length of each character list (like 'aa'..'az') would be its weight timed the given length, such as:  $Length_{aa'to'az'} = \frac{Occurrence_{aa'to'az'}}{Occurrence_{aa'to'zz'}} * (GivenLength)$ , so that we could generate a random string which had a similar distribution but different content.

## 5 Evaluation

The testing environment was macOS Mojave 10.14.5 system with iTerm and Microsoft (R) F# Compiler version 4.1.

We combined all the assignment into a single program, compiled and tested it with parameter `readfile littleClausAndBigClaus.txt`. Output of it included *content of target file*, *converted string of original string*, *histogram of original string*, *random string of single character*, *histogram of random single character*, *cooc of original string*, *cooc of random single character*, *random string of character pairs* and *cooc of random character pairs*. Part of results (because it was too long) were shown as following (the complete results could be seen at [result](#)):

```

→ fsharp --nologo Assig2.fsx && mono Assig2.exe readfile littleClausAndBigClaus.txt
str :
In a village there lived two men who had the self-same name. Both were named Claus. But

```

Figure 1: Testing of readText function

```

cooc of original string:
[[0; 35; 64; 77; 0; 12; 31; 7; 84; 0; 26; 64; 28; 228; 0; 16; 0; 100; 147; 141;
  122; 36; 26; 3; 43; 3];
[19; 0; 0; 0; 64; 0; 0; 0; 44; 0; 0; 12; 0; 0; 29; 0; 0; 8; 0; 0; 61; 0; 0; 0;
  16; 0];

```

Figure 2: Testing of fstOrderMarkovModel function

```

cooc of random character pairs:
[[148; 23; 21; 42; 214; 22; 13; 137; 95; 1; 24; 74; 25; 31; 124; 14; 0; 33; 48;
  90; 30; 3; 35; 3; 27; 0];
[19; 8; 4; 14; 19; 4; 6; 21; 12; 0; 2; 9; 2; 11; 20; 5; 1; 6; 28; 26; 9; 1; 13;
  0; 5; 0];

```

Figure 3: Cooccurrence of fstOrderMarkovModel function

## 6 Conclusion

In this assignment, we worked with simple text processing and met the requirements of analyzing and generating target text according to statistics. It was time-consuming of thinking about the recursive function and avoiding to use things like 'for loop' when we were stuck, it was challenging but this assignment made us get more familiar with F# and functional programming indeed. It was our pleasure to have teachers like Prof.Sporring and we believed that we had a good time. Many thanks!

## References

[1] Jon Spurring. *Learning to Program with F#*. 2019.