

# Raport Haskell Words Counting

Mikołaj Wielebnowski Krzysztof Kubień

9 stycznia 2021

## 1 Wstęp

W niniejszym raporcie znajduje się opis prac przeprowadzonych w celu stworzenia programów w języku Haskell . Projekt został oparty o poradnik i został stworzony w celu zaliczenia projektu na przedmiot *"Programowanie Funkcyjne"*.

## 2 Opis projektu

Etapy prac:

1. Wstęp.
2. Przygotowanie środowiska Haskell.
3. Wyświetlanie liczb w innej bazie.
4. Odczytywanie liczb z innej bazy.
5. Znajdowanie ciągów znaków używając Data.ByteString.
6. Znajdowanie ciągu znaków używając algorytmu Boyer-Moore-Horspool.
7. Poszukiwanie ciągu znaków używając algorytmu Rabin-Karpa.
8. Rozdzielenie ciągu znaków.
9. Znajdowanie najdłuższej wspólnej sekwencji.
10. Kod fonetyczny.
11. Przeliczanie dystansu edycji.
12. Przeliczanie dystansu Jaro-Winkler pomiędzy dwoma ciągami.
13. Znajdowanie ciągu z jednym dystansem.
14. Naprawianie błędów w zapisie.
15. Podsumowanie.

## 2.1 Wstęp.

Haskell to czysto funkcyjny język programowania nazwany na cześć Haskella Curry'ego. Był on początkowo intensywnie rozwijany wokół ośrodka University of Glasgow, popularny kompilator tego języka to Glasgow Haskell Compiler (GHC) kompilujący szybki kod maszynowy porównywalny w szybkości wykonania do kodów z GCC.



Rysunek 1: Logo Haskell

Programy w projekcie pokazują podstawowe sposoby interpretacji ciągów znaków oraz zmiany na nie liczb, co po rozwinięciu pozwala na modyfikację oraz analizę dużej ilości tekstu. Takie techniki w dzisiejszych czasach są bardzo istotne i przyczyniają się zarówno do analiz big data, jak i technologiami AI i badaniem tekstów takich jak dzieła literackie.

## 2.2 Przygotowanie środowiska Haskell.

Przed rozpoczęciem pracy nad projektem zespół musiał przygotować środowisko programistyczne dla języka Haskell. W tym celu postanowiono posłużyć się edytorem Visual Studio Code. Następnie użyto zakładki extensions w celu dodania środowiska Haskell, oraz poprzez instalację pakietu Chocolatey. Zespół korzystał z kompilatora ghci.

Kolejnym krokiem zespołu było zapoznanie się z treścią opracowania zawartego w pliku WordsCounting.pdf

## 2.3 Wyświetlanie liczb w innej bazie.

Ten program pozwala na wczytanie liczby dziesiętnej  $n$  i otrzymanie jako output liczby w postaci stringu w zadanej bazie  $b$ . Funkcja `showIntAtBase` bierze za bazę wybraną liczbę, a następnie przypisuje ją do odpowiedniego dla bazy znaku z przedziału `0, 1, 2, ...`, `a, b, ...`, `z`. Aby funkcja mogła zadziałać poprawnie wymaga zamiany liczby na ciąg cyfr, tą funkcję pełni tutaj `numToLetter`

```
liczbawinnejbazie.hs > main
1  import Data.Char (intToDigit, chr,ord)
2  import Numeric (showIntAtBase)
   inBase :: (Integral a, Show a) => a -> a -> String
3  n `inBase` b = showIntAtBase b numToLetter n ""
4  numToLetter :: Int -> Char
5  numToLetter n
6      | n < 10 = intToDigit n
7      | otherwise = chr (ord 'a'+ n-10)
8  main :: IO ()
9  main = do
10     putStrLn $ 8 `inBase` 12
11     putStrLn $ 10 `inBase` 12
12     putStrLn $ 12 `inBase` 12
13     putStrLn $ 47 `inBase` 12
```

Rysunek 2: liczbawinnejbazie.hs

Podane dla przykładu liczby 8, 10, 12 oraz 47, po przekonwertowaniu ich na bazę 12 dają wynik:

```
8
a
10
3b
```

Rysunek 3: Wynik programu.

## 2.4 Odczytywanie liczb z innej bazy.

W tym programie staramy się osiągnąć odwrotność poprzedniej operacji, czyli przekonwertować liczbę w postaci stringa zapisanej w systemie liczbowym o danej bazie na liczbę dziesiętkową typu int. Program dla zadanego stringa stosuje funkcję `readInt` wczytuje liczbę i konwertuje ją na daną bazę. Aby funkcja zadziałała poprawnie należy zamienić typ zmiennych ze stringów na int (`letterToNum`) i sprawdzić czy otrzymany znak jest zgodny (`isValidDigit`).

```
odczytbazy.hs > ...
1 import Data.Char (ord, digitToInt, isDigit)
2 import Numeric (readInt)
  base :: Num a => String -> a -> [(a, String)]
3 str `base` b = readInt b isValidDigit letterToNum str
4 letterToNum :: Char -> Int
5 letterToNum d
6   | isDigit d = digitToInt d
7   | otherwise = ord d - ord 'a' + 10
8
9 isValidDigit :: Char -> Bool
10 isValidDigit d = letterToNum d >= 0
11 main :: IO ()
12 main = do
13   print $ "8" `base` 12
14   print $ "a" `base` 12
15   print $ "10" `base` 12
16   print $ "3b" `base` 12
```

Rysunek 4: odczytbazy.hs

Po podstawieniu przykładowych liczb 8, a, 10, 3b w bazie 12 otrzymujemy następujący output:

```
[(8, "")]
[(10, "")]
[(12, "")]
[(47, "")]
```

Rysunek 5: Wynik programu.

## 2.5 Znajdowanie ciągów znaków używając Data.ByteString.

W tym programie pobieramy ciąg znaków, przy pomocy jednego z dostępnych algorytmu dzielącego ciągi znaków z biblioteki Data.ByteString, który nadaje się do używania na dużej ilości danych ze względu na swoją efektywność. Program otrzymując dwa string konwertuje je przy pomocy funkcji Char8 na 8 bitowe, następnie następnie przy pomocy algorytmu breakSubstring poszukuje zadanego w pierwszej części (query) inputu ciągu w drugim ciągu. Gdy ciąg zostanie znaleziony program wyrzuca wartości logiczne True lub False, w zależności czy został on znaleziony czy nie.

```
databyte.hs > main
1 import Data.ByteString (breakSubstring)
2 import qualified Data.ByteString.Char8 as C
3 substringFound :: String -> String -> Bool
4 substringFound query str =
5     (not . C.null . snd) $
6     breakSubstring (C.pack query) (C.pack str)
main :: IO ()
7 main = do
8     print $ substringFound "scraf" "swedish scraf mafia"
9     print $ substringFound "flute" "swedish scraf mafia"
```

Rysunek 6: databyte.hs

### 3 Znajdowanie ciągu znaków używając algorytmu Boyer-Moore-Horspool.

Algorytm Boyer'a-Moor'a-Horspool'a to algorytm doskonale nadający się do znajdowania ciągów w innych ciągach. Ten wyjątkowo dobrze sprawdza się dłu-gich ciągów i poszukiwania ich w bardzo dużych ciągach, ponieważ algorytm jest w stanie efektywnie pomijać części tekstu, dla których porównania już zostały wykonane. W tym programie zastosowana jest uproszczona wersja określana ja-ko algorytm Horspool'a, która nie wymaga postprocessingu i jednocześnie traci znacznie na efektywności na dużych ciągach, ale jest nie jest tak „przeciążony” na początku.

Program podobnie jak w poprzednim poszukuje ciągu z pierwszej części inputu w drugiej. Algorytm w uproszczonej formie zastosowanej tutaj tworzy „okno”, o wielkości równej zadanemu ciągowi i porównuje znak po znaku następnie po-ruszając „okno” z prawej do lewej.

```
boyermoore.hs > bmh'
1 import Data.Map (fromList, (!), findWithDefault)
  indexMap :: (Ord k, Num k, Enum k) => [a] -> Data.Map.Internal.Map k a
2 indexMap xs = fromList $ zip [0..] xs
  revIndexMap :: (Ord k, Num a, Enum a) => [k] -> Data.Map.Internal.Map k a
3 revIndexMap xs = fromList $ zip (reverse xs) [0..]
4 bmh :: Ord a => [a] -> [a] -> Maybe Int
5 bmh pat xs = bmh' (length pat - 1) (reverse pat) xs pat
6 bmh' :: Ord a => Int -> [a] -> [a] -> [a] -> Maybe Int
7 bmh' n [] xs pat = Just (n + 1)
8 bmh' n (p:ps) xs pat
9   | n >= length xs = Nothing
10  | p == (indexMap xs) ! n = bmh' (n - 1) ps xs pat
11  | otherwise       = bmh' (n + findWithDefault
12    (length pat) (sMap ! n) pMap)
13    (reverse pat) xs pat
14  where sMap = indexMap xs
15        pMap = revIndexMap pat
16 main :: IO ()
17 main = print $ bmh "wor" "Hello World"
```

Rysunek 7: boyermoore.hs

W tej implementacji wynik programu pokazuje ostatnie miejsce przed roz-poczęciem pierwszego wystąpienia zadanego stringu stąd dla przykładu wynik to:

```
Just 6
```

Rysunek 8: Wynik programu.hs

## 4 Poszukiwanie ciągu znaków używając algorytmu Rabin-Karpa.

Kolejnym algorytmem, który stosuje się do poszukiwania do poszukiwania ciągu w ciągu jest algorytm Rabina-Karpa, ten algorytm porównuje „okienko” z tekstem, ale nie znak po znaku, a przez numer o bazie 26+, identyfikujący zestaw znaków. Algorytm ten nie jest zbyt szybki dla pojedynczego przeszukania, jednak zaczyna być bardzo sprawny, gdy zapytań jest więcej, ponieważ po pierwszym przejściu przez tekst jest w stanie szybko znajdować kolejne.

Określenia haystack (stóg siana) oraz needle (igła) stosowane są w angielskim w celu określenia odpowiednio przeszukiwanego tekstu i zapytania.

Po zainstalowaniu modułu cabala stringsearch, algorytm Rabina-Karpa można zaimplementować z biblioteki ByteString. Program zamienia litery z naszego zapytania używając miejsca na którym znajduje się znak jako potęgę bazy i funkcji ord do zamiany znaku na liczbę jej odpowiadającą. I po raz kolejny program porusza „oknem”, tym razem lewej do prawej, porównując wartości identyfikujące znaki i w ten sposób znajdując dopasowania.

```
rabinkarp.hs > ...
1 {-# LANGUAGE OverloadedStrings #-}
2 import Data.ByteString.Search.KarpRabin (indicesOfAny)
3 import qualified Data.ByteString as BS
4 main = do
5     .....let needles = ["preparing to go away"
6     .....|.....|.....|....., "is some letter of recommendation"]
7     .....haystack <- BS.readFile "big.txt"
8     print $ indicesOfAny needles haystack
```

Rysunek 9: rabinkarp.hs

Po wykonaniu programu otrzymujemy listę, która składa się z numeru określającego położenie i liczbę w nawiasie określającą numer zapytania.

```
[(3738968, [1]), (5632846, [0]), (5714386, [0])]
```

Rysunek 10: Wynik programu.

## 5 Rozdzielenie ciągu znaków.

Problem ze użytecznymi danymi jest taki, że trzeba je w jakiś sposób rozdzielić i do tego idealna jest funkcja `splitOn`, która pochodzi z modułu `cabala.split`. Dzieli ona dane na osobne ciągi jak i usuwa separator. Ten program ma za zadanie wczytać dane z pliku `input.txt`, a następnie podzielić zawartość na różne sposoby poprzez linie, spacje, przecinki lub dowolne znaki odpowiednio dla każdej linijki z pliku.

Program jest banalnie prosty wystarczy wczytać plik i użyć `splitOn` z odpowiednim separatorem podanym jako pierwszy argument, gdzie wyjątkiem jest nowa linia, ponieważ tak zaimportowany plik Haskell dzieli automatycznie.

```
splitting.hs
1 import Data.List.Split (splitOn)
2 main = do
3   ... input <- readFile "input.txt"
4   ... let ls = lines input
5   ... print $ ls
6   ... let ws = words $ ls !! 2
7   ... print ws
8   ... let cs = splitOn "," $ ls !! 3
9   ... print cs
10  ... let ds = splitOn "an" $ ls !! 4
11  ... print ds
```

Rysunek 11: `splitting.hs`

W pierwszym przykładzie dzielimy po linijce:

```
first line
second line
words are split by space
comma,separated,values
or any delimiter you want
```

Rysunek 12: Wynik pierwszego przykładu

W drugim po spacji

```
words,are,split,by,space
```

Rysunek 13: Wynik drugiego przykładu

W trzecim po przecinku:

```
comma,separated,values
```

Rysunek 14: Wynik trzeciego przykładu

W czwartym po dowolnym znaku w tym wypadku „e”:

```
or any d,limit,r you want
```

Rysunek 15: Wynik czwartego przykładu



## 6 Znajdowanie najdłuższej wspólnej sekwencji.

Jednym ze sposobów na znajdowanie podobieństw w ciągach znaków jest znajdowanie ich najdłuższej wspólnej sekwencji. Sposób ten jest użyteczny w znajdowaniu mutacji danych. Algorytm ten jest zaimplementowany z zapamiętaniem rekursywnych przywołań. Jeśli dwa pierwsze składniki listy są takie same, wtedy najdłuższą wspólną sekwencją jest funkcja lcs zastosowana do pozostałej części listy. W przeciwnym przypadku najdłuższą wspólną sekwencją jest ta, która jest najdłuższy z możliwych. Do usprawnienia programu została użyta funkcja memorize aby zapamiętać wcześniej przeliczone wartości.

```
common.hs
1 import qualified Data.MemoCombinators as Memo
2 memoize :: (String -> String -> r) -> String -> String -> r
3 memoize = Memo.memo2
4 ... (Memo.list Memo.char) (Memo.list Memo.char)
5 lcs :: String -> String -> String
6 lcs = memoize lcs'
7 where lcs' xs@(x:xs) ys@(y:ys)
8     | x == y = x : lcs xs ys
9     | otherwise = longer (lcs xs' ys) (lcs xs ys')
10     lcs' _ _ = []
11     longer as bs
12     | length as > length bs = as
13     | otherwise = bs ...
14 main :: IO ()
15 main = do
16     let xs = "find the lights"
17         ys = "there are four lights"
18     print $ lcs xs ys
```

Rysunek 16: common.hs

Po wykonaniu programu otrzymujemy następujący wynik:

**the lights**

Rysunek 17: Wynik programu

## 7 Kod fonetyczny.

Kiedy mamy do czynienia z korpusem angielskich słów to możemy je podzielić na fonetyczny kod żeby zobaczyć jak podobnie brzmią. Kod fonetyczny działa dla dowolnego alfabetycznego ciągu znaków. Będziemy korzystać z pakietu `Text.PhoneticCode`. Program wykonuje manipulacje ciągiem bazując na heurystyce wzorców języka angielskiego.

```
phonetic.hs
1 import Text.PhoneticCode.Soundex (soundexNARA, ... soundexSimple)
2 import Text.PhoneticCode.Phonix (phonix)
3 ws = ["haskell", "hackle", "haggle", "hassle"]
4 main :: IO ()
5 main = do
6     ... print $ map soundexNARA ws
7     ... print $ map soundexSimple ws
8     ... print $ map phonix ws
```

Rysunek 18: phonetic.hs

Po wykonaniu otrzymujemy wynik:

```
H240,H240,H240,H240][H240,H240,H240,H240][H82,H2,H2,H8]
```

Rysunek 19: wynik programu

## 8 Przeliczanie dystansu edycji.

Dystans edycji lub dystans Levenshteina jest minimalną liczbą prostych operacji na ciągu wymaganej do konwersji jednego ciągu znaków na drugi.

$$\text{lev}_{a,b}(i,j) = \begin{cases} \max(i,j) & \text{if } \min(i,j) = 0, \\ \min \begin{cases} \text{lev}_{a,b}(i-1,j) + 1 \\ \text{lev}_{a,b}(i,j-1) + 1 \\ \text{lev}_{a,b}(i-1,j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

Rysunek 20: Formuła

Kod będzie bezpośrednią interpretacją powyższej formuły. Algorytm rekurencyjnie próbuje usuwać, dodawać oraz zamieniać na wszelki sposób do czasu aż znajdzie minimalny dystans z jednego ciągu do drugiego.

```
editdist.hs
1 import qualified Data.MemoCombinators as Memo
2 lev :: Eq a => [a] -> [a] -> Int
3 lev a b = levM (length a) (length b)
4 where levM = memoize lev'
5     lev' i j
6     | min i j == 0 = max i j
7     | otherwise = minimum
8     [ ( 1 + levM (i-1) j )
9     , ( 1 + levM i (j-1) )
10    , ( ind i j + levM (i-1) (j-1) ) ]
11 ind i j
12 | a !! (i-1) == b !! (j-1) = 0
13 | otherwise = 1
14 memoize = Memo.memo2 (Memo.integral) (Memo.integral)
15 main = print $ lev "mercury" "sylveste" |
```

Rysunek 21: editdist.hs

Po wykonaniu trzymujemy poniższy wynik:

8

Rysunek 22: Wynik programu

## 9 Przeliczanie dystansu Jaro-Winkler pomiędzy dwoma ciągami.

Dystans Jaro-Winkler'a określa podobieństwo pomiędzy ciągami za pomocą liczby z zakresu 0 i 1, gdzie 0 to brak podobieństwa a 1 to pełne podobieństwo. Algorytm wyraża się za pomocą matematycznego równania przedstawionego poniżej:

The Jaro distance  $d_j$  of two given strings  $s_1$  and  $s_2$  is

$$d_j = \begin{cases} 0 & \text{if } m = 0 \\ \frac{1}{3} \left( \frac{m}{|s_1|} + \frac{m}{|s_2|} + \frac{m-t}{m} \right) & \text{otherwise} \end{cases}$$

where:

- $m$  is the number of matching characters (see below);
- $t$  is half the number of transpositions (see below).

Rysunek 23: Formuła

Gdzie  $s_1$  i  $s_2$  to dwa ciągi,  $m$  to liczba identycznych znaków w maksymalnie połowie długości ciągu dłuższego,  $t$  to połowa ilości identycznych znaków, ale o różnych pozycjach.

```
jarowinkler.hs > jaro
1 import Data.List (elemIndices)
2 jaro :: Eq a => [a] -> [a] -> Double
3 jaro s1 s2
4   | m == 0   = 0.0
5   | otherwise = (1/3) * (m/ls1 + m/ls2 + (m-t)/m)
6   where ls1 = toDouble $ length s1
7         ls2 = toDouble $ length s2
8         m' = matching s1 s2 d
9         d = fromIntegral $
10            max (length s1) (length s2) `div` 2 - 1
11        m = toDouble m'
12        t = toDouble $ (m' - matching s1 s2 0) `div` 2
13 toDouble :: Integral a => a -> Double
14 toDouble n = (fromIntegral n) :: Double
15 matching :: Eq a => [a] -> [a] -> Int -> Int
16 matching s1 s2 d = length $ filter
17   (\(c,i) -> not (null (matches s2 c i d)))
18   (zip s1 [0..])
19 matches :: Eq a => [a] -> a -> Int -> [Int]
20 matches str c i d = filter (<= d) $
21   map (dist i) (elemIndices c str)
22   where dist a b = abs $ a - b
main = IO ()
23 main = do
24   print $ jaro "marisa" "magical"
25   print $ jaro "haskell" "hackage"
```

Rysunek 24: jarowinkler.hs

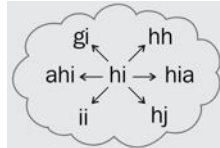
Powyższy program po zaimplementowaniu algorytmu wyrzuca w wyniku dwie wartości dla przykładu marisa i magical oraz haskell i hackage, przypisując im poniższe wartości:

```
0.746031746031746
0.7142857142857142
```

Rysunek 25: Wynik programu

## 10 Znajdowanie ciągu z jednym dystansem.

Ten program znajduje ciągi, które różnią się od zadanego ciągu o jedną długość edycji, czyli szuka wszystkich elementów, na których trzeba wykonać tylko jedną akcję, aby przekształcić go do zadanego ciągu, tworząc rodzinę dla wybranego słowa jak na poniższym rysunku:



Rysunek 26: Rysunek

```

oneedit.hs > ...
1 import Data.Char (toLower)
2 import Data.List (group, sort)
3 edits1 :: String -> [String]
4
5 edits1 word = unique $
6     deletes ++ transposes ++ replaces ++ inserts
7   where splits = [ (take i word', drop i word') |
8       i <- [0..length word'] ]
9     deletes = [ a ++ (tail b) |
10        (a,b) <- splits, (not.null) b ]
11    transposes = [ a ++ [b!i1] ++ [head b] ++ (drop 2 b) |
12        (a,b) <- splits, length b > 1 ]
13    replaces = [ a ++ [c] ++ (drop 1 b)
14        | (a,b) <- splits
15          , c <- alphabet
16            , (not.null) b ]
17    inserts = [ a ++ [c] ++ b
18        | (a,b) <- splits
19          , c <- alphabet ]
20    alphabet = ['a'..'z']
21    word' = map toLower word
22 unique :: [String] -> [String]
23 unique = map head.group.sort
24 main = IO ()
25 main = print $ edits1 "hi"

```

Rysunek 27: oneedit.hs

Dla wybranego przykładu "hi" z dużej bazy słów 1, 2 i 3 literowych otrzymujemy poniższy wynik:

```

["ahi","ai","bhi","bi","chi","ci","dhi","di","ehi","ei","fhi","fi","ghi","gi","h","ha","hal","hb","hbi","hc","hci","hd","hdi","he","hei","hf","hfi","hg","hgi","hh","hhi","hi","hia","hib","hic","hid","hie","hif","hig","hih","hii","hij","hik","hil","hiam","hin","hio","hip","hiq","hir","his","hit","hiu","hiv","hiw","hix","hly","hiz","hj","hji","hk","kki","kl","kll","km","kml","kn","knl","ko","kol","kp","kpl","kq","kql","kr","krl","ks","ksl","kt","ktl","ku","kul","kv","kvl","kw","kwl","kx","kxl","ky","kyl","kz","kzl","li","lhi","lii","lji","ljj","khi","ki","lhi","li","lmi","mi","nhi","ni","ohi","ol","phi","pl","qhi","qi","rhi","ri","shi","si","thi","ti","uhi","ui","vhi","vi","whi","wi","xhi","xi","yhi","yi","zhi","zi"]

```

Rysunek 28: Wynik programu

## 11 Naprawianie błędów w zapisie.

Ten program korzystając z algorytmu Petera Norvig'a oraz bazy słów poprawnie zapisanych. Algorytm porównuje dane słowa dzieląc je na ciągi i tworzy rodzinę ze słowami z bazy tak jak w poprzednim przykładzie. Algorytm zakłada, że literówki występuje dla pojedynczego lub podwójnego dystansu edycji, a następnie wybiera słowo z rodziny, korzystając z hierarchii opartej na częstości ich używania, zastępując badany ciąg. Ten algorytm jest dość prosty, jednak dzięki temu działa bardzo szybko, a wykonuje zadanie, które zakrawa o uczenie maszynowe.

```
spelling.hs >
1 import Data.Char (isAlpha, isSpace, toLower)
2 import Data.List (group, sort, maximumBy)
3 import Data.Ord (comparing)
4 import Data.Map (fromListWith, Map, member, (!))
5 autofix :: Map String Int -> String -> String
6 autofix m sentence = unwords $
7   map (correct m) (words sentence)
8 getWords :: String -> [String]
9 getWords str = words $
10   filter (\x -> isAlpha x || isSpace x) lower
11   where lower = map toLower str
12 train :: [String] -> Map String Int
13 train = fromListWith (+) . ('zip' repeat 1)
14 edits1 :: String -> [String]
15 edits1 word = unique $
16   deletes ++ transposes ++ replaces ++ inserts
17   where splits = [ (take 1 word', drop 1 word')
18     | l <- [0..length word'] ]
19   deletes = [ a ++ [tail b]
20     | (a,b) <- splits
21     , (not.null) b ]
22   transposes = [ a ++ [b !! 1] ++ [head b] ++ (drop 2 b)
23     | (a,b) <- splits, length b > 1 ]
24   replaces = [ a ++ [c] ++ (drop 1 b)
25     | (a,b) <- splits, c <- alphabet ]
26   inserts = [ a ++ [c] ++ b |
27     (a,b) <- splits, c <- alphabet ]
28   alphabet = ['a'..'z']
29   word' = map toLower word
30 knownEdits2 :: String -> Map String a -> [String]
31 knownEdits2 word m = unique $ [ e2
32   | e1 <- edits1 word
33   , e2 <- edits1 e1
34   , e2 `member` m ]
35 unique :: [String] -> [String]
36 unique = map head . group . sort
37 known :: [String] -> Map String a -> [String]
38 known ws m = filter ('member' m) ws
39 correct :: Map String Int -> String -> String
40 correct m word = maximumBy (comparing (!)) candidates
41   where candidates = head $ filter (not.null)
42     [ (known [word] m)
43     , (known (edits1 word) m)
44     , (knownEdits2 word m)
45     , [word] ]
46 main :: IO ()
47 main = do
48   rawText <- readFile "big.txt"
49   let m = train $ getWords rawText
50       let sentence = "such codez many hsakell very spel so korrekt"
51       print $ autofix m sentence
```

Rysunek 29: spelling.hs

Jako przykład zastosowano „such codez many hsakell very spel so korrekt”, w którym znajdują się literówki w odpowiednich dystansach edycji, dzięki czemu program wyrzuca poprawiony tekst:

```
"such code many hsakell very spell so correct"
```

Rysunek 30: Wynik programu

## 12 Podsumowanie

Podsumowując prace związane z powyższymi programami zespół stwierdził, że język Haskell jest bardzo użyteczny jeśli chodzi o operacje na dużych ciągach znaków, a nawet korekcję błędów. Warto dodać, że strona Facebook używa właśnie języka Haskell do korekcji błędów literowych.

Jednakże zespół napotkał poważny problem w importowaniu niektórych bibliotek. Po długim poszukiwaniu powodu owego błędu odkryto, że winą była niekompletna instalacja pakietu cabal, przez co dodatkowe importowanie bibliotek poleceniem cabal było przeprowadzane ze skutkiem negatywnym.

## Literatura

- [1] *WordsCounting.pdf*
- [2] *<http://norvig.com/spell-correct.html>*
- [3] *<https://hoogle.haskell.org>*