

Лабораторная работа №1 – 3  
по дисциплине «Сравнительный анализ языков программирования»  
**Работа с STL C++, C# и Java. Исследование  
сложности алгоритма**

Факультет	ПМИ
Группа	ПМ 33, ПМ 32
Студенты	Жиляков А. Антоненко В. Губский Е.
Преподаватели	Рояк М. Э. Ступаков И. М.
Вариант	5

# I Указания к выполнению работы

Разработайте программу для решения поставленной задачи на языке C++, удовлетворяющую следующим требованиям:

1. программа должна использовать для ввода-вывода потоки **STL**,
2. программа не должна содержать собственных реализаций стандартных алгоритмов и структур данных, а использовать существующие в **STL**,
3. размер каждой подпрограммы не должен превышать 10 строк.

Протестируйте разработанную программу. Исследуйте асимптотические свойства разработанной программы на системе тестов с возрастающей размерностью.

Разработайте программы на языках C# и Java с использованием аналогичных структур данных и сравните результаты работы.

(<http://dispace.edu.nstu.ru/didesk/course/show/5626/5>)

## I.1 Вариант задания

Задан большой текст (книга). Для наиболее часто встречающегося слова найти  $N$  наиболее часто встречающихся пар слов.

# II Размышления

Для поиска наиболее популярного слова (или пары слов) удобно использовать контейнер `std::map` — ассоциативный массив с красно-чёрным деревом под кашушоном. Ключ суть слово, значение — частота. (Аналоги в Java и C# — `TreeMap` и `SortedDictionary`)

Через  $n$  обозначим количество слов в нашей входной книжке. Тогда для занесения в контейнер данных потребуется  $O(n \log m)$  времени (поиск в дереве осуществляется за логарифм), где  $m \leq n$  — количество слов без повторений.

Положим  $m = n$  (в такой книжке каждое слово уникально), чтобы проверить на практике логарифмическую сложность — в действительности  $m \ll n$  и алгоритм работает практически за линейное время. Будем увеличивать длину книжки в 4 раза при начальном  $n = 250000$ . Тогда отношение времени работы программы суть

$$\frac{4n \log_2 4n}{n \log_2 n} = 4 \left( 1 + \frac{2}{\log_2 n} \right) \searrow 4, n \rightarrow \infty, \quad (1)$$

то есть предполагается, что программа станет работать медленнее в 4+ раз, при этом при  $n \rightarrow \infty$  отношение (1) должно становиться всё меньше и меньше.

# III Тестирование

## III.1 Тестирование на реальной книжке

Для имплементации на всех языках результат один и тот же, во что и хотелось бы верить.

частота	пара слов
65	martin eden
44	and martin
42	martin was
40	martin had
27	that martin
25	martin said
20	but martin
19	to martin
18	martin did
18	martin and
17	you martin
13	when martin
13	martin answered
13	as martin
12	martin went

Таблица 1: “Martin Eden,” Jack London. Число слов — 138 754

## III.2 Тестирование сложности и проверка «размышлений»

Результаты тестирования при  $m = n$ .

$i$	$n$	$t$ , сек	$t_i/t_{i-1}$
1	250 000	.742	
2	1 000 000	3.855	5.2
3	4 000 000	19.738	5.12
4	16 000 000	94.566	4.79

Таблица 2: C++

$i$	$n$	$t$ , сек	$t_i/t_{i-1}$
1	250 000	5.01	
2	1 000 000	26.53	5.2
3	4 000 000	101.15	3.8

Таблица 3: C#

$i$	$n$	$t$ , сек	$t_i/t_{i-1}$
1	250 000	5.32	
2	1 000 000	27.01	5.08
3	4 000 000	108.47	4.02

Таблица 4: Java

## IV Приложение

Компилируем исходники с помощью MSVC:

```
cl -EHsc -Ox testGen.cpp user.cpp,
```

генерируем тест из  $n$  разных слов ( $n$  подаётся на стандартный ввод):

```
testGen > dummy.txt
```

и наша программа готова к работе:

```
user < dummy.txt » out.txt.
```

## IV.1 testGen.cpp

Генератор больших тестов — общий для всех языков.

```
1 #include <random>
2 #include <algorithm>
3 #include <iterator>
4 #include <iostream>
5
6 int main() {
7     int n;
8     std::cin >> n;
9     std::vector<int> v(n);
10    for (int i = 0; i < n; ++i)
11        v[i] = i;
12    std::random_device rd;
13    std::mt19937 g(rd());
14    std::shuffle(v.begin(), v.end(), g);
15    std::copy(v.begin(), v.end(), std::ostream_iterator<int>(std::cout, " "));
16    std::cout << "\n";
17    return 0;
18 }
```

## IV.2 C++

### user.cpp

```
1 #include <iostream>
2 #include <fstream>
3 #include <vector>
4 #include <map>
5 #include <set>
6 #include <string>
7 #include <algorithm>
8 #include <ctime> //замерить время
9 #include <cctype> //isalnum
10
11 using namespace std; //программа коротенькая, так что можно и расслабиться
12
13 //сегмент данных
14
15 vector<string> book; //текст книги
16 string word; //самое популярное слово
17 set<string> words2ignore; //игнорируемые слова (предлоги, артикли и т.д.)
18 unsigned N = 15; //количество выводимых сочетаний слов
19
20
21 void getBook() { //считываем текст книжки из стандартного ввода в контейнер
22     while (cin >> word) {
23         transform(word.begin(), word.end(), word.begin(), ::tolower); //прописные буквы к строчным
24         word.erase(remove_if(word.begin(), word.end(), [](char c){return !isalnum(c);}), word.end()); //избавляемся от лишних символов в слове
25         if (word.length()) book.push_back(word);
26     }
27 }
28
29 void getWord() { //находим самое популярное слово
30     map<string, size_t> freq; //ключ - слово, значение - частота появления
31     size_t i, maxFreq = 0; //максимальная частота появления слова
32     for (i = 0; i < book.size(); ++i)
33         if (!words2ignore.count(book[i]) && ++freq[book[i]] > maxFreq) {
34             word = book[i];
35             maxFreq = freq[word];
36         }
37 }
38
39 int main() {
40     map<string, size_t> freq; //ключ - сочетание слов, значение - частота появления
41     map<string, size_t>::const_iterator iter;
42     vector<pair<size_t, string>> res; //отсортированный (см. далее) по частоте контейнер
43     size_t i;
44     clock_t begTime, endTime;
45     ifstream ignoreFile("ignore.txt");
46     while (ignoreFile >> word) words2ignore.insert(word);
47     getBook();
48
49     begTime = clock();
50     getWord();
51     for (i = 1; i < book.size() - 1; ++i)
52         if (book[i] == word) {
53             ++freq[book[i - 1] + " " + word];
54             ++freq[word + " " + book[i + 1]];
55         }
56     for (iter = freq.begin(); iter != freq.end(); ++iter)
57         res.push_back(make_pair(iter->second, iter->first));
58     sort(res.rbegin(), res.rend());
59     endTime = clock();
60     N = min(N, res.size());
61     for (i = 0; i < N; ++i)
62         cout << res[i].first << ' ' << res[i].second << '\n';
63     cout << "time diff: " << double(endTime - begTime) / CLOCKS_PER_SEC << '\n';
64     cout << "words count: " << book.size() << "\n\n";
65
66     return 0;
67 }
```

## IV.3 C#

### Program.cs

```
1 using System;
2 using System.Collections.Generic;
3 using System.Linq;
4 using System.Text;
5 using System.Threading.Tasks;
6
7 namespace LAB2
8 {
9     class Program
10     {
11         static void Main(string[] args)
12         {
13             Analyzer a = new Analyzer("Martin Eden.txt", "ignore.txt");
14             Console.WriteLine("MaxFreeQ={0}", a.MaxFreeq);
15             int n = 15;
16             for (int i = 0; i < a.Pairs.Count && i < n; ++i)
17                 Console.WriteLine("{0}\t{i}", i + 1, a.Pairs[i]);
18             Console.WriteLine(a.RunTime);
19             Console.ReadKey();
20         }
21     }
22 }
```

### Analyzer.cs

```
1 using System;
2 using System.Collections.Generic;
3 using System.Linq;
4 using System.Text;
5 using System.Threading.Tasks;
6 using System.Diagnostics;
7
8 namespace LAB2
9 {
10     /// <summary>
11     /// анализирует файл
12     /// </summary>
13     public class Analyzer
14     {
15         public long RunTime;
16         public class WordFreeq
17         {
18             public readonly string Word;
19             public readonly uint Freeq;
20
21             public WordFreeq(string Word, uint Freeq)
22             {
23                 this.Word = Word;
24                 this.Freeq = Freeq;
25             }
26             public override string ToString()
27             {
28                 return "[" + Word + ", " + Freeq + "]";
29             }
30         }
31         List<string> book = new List<string>(); //строки книги
32         public HashSet<string> wordsToIgnore = new HashSet<string>(); //слова, которые не учитываются
33         public SortedDictionary<string, uint> freeq = new SortedDictionary<string, uint>(); //набор слов и их встречаемость
34         public List<WordFreeq> Pairs = new List<WordFreeq>(); //пары элементов сортированные по встречаемости
35
36         public KeyValuePair<string, uint> MaxFreeq
37         {
38             get
39             {
40                 KeyValuePair<string, uint> max = new KeyValuePair<string, uint>("", 0);
41                 foreach (var v in freeq)
42                 {
43                     if (max.Value == 0 || v.Value > max.Value) max = v;
44                 }
45                 return max;
46             }
47         }
48
49         public readonly string FileName;
50         public Analyzer(string fileName, string ignoreFile)
51         {
52             LoadExcludeWords(ignoreFile);
53             this.FileName = fileName;
54             foreach (var v in wordsToIgnore)
55                 this.wordsToIgnore.Add(v.ToLower());
56             ReadBook(fileName);
57         }
58         void LoadExcludeWords(string fileName)
59         {
60             wordsToIgnore.Clear();
61             System.IO.StreamReader r = new System.IO.StreamReader(fileName);
62             foreach (var word in r.ReadToEnd().Split(' ', ',', '\n', '\t', '\0', '\r'))
63                 wordsToIgnore.Add(word.ToLower());
64             r.Close();
65         }
66
67         void ReadBook(string fileName)
68         {
69             //открываем файловый поток
70             System.IO.StreamReader r = new System.IO.StreamReader(fileName);
71             //ограничитель
72             //читаем файл по словам в нижнем регистре
73             book.Clear();
74             foreach (var word in r.ReadToEnd().Split(' ', ',', '\n', '\t', '\0', '\r'))
75                 book.Add(word.ToLower());
76         }
77     }
78 }
```

```

76         r.Close();
77         //удаляем все ненужное
78         book.RemoveAll(x => string.IsNullOrEmpty(x) || wordsToIgnore.Contains(x));
79         //определяем сколько раз встретилось каждое слово
80         freeq.Clear();
81         var watch = System.Diagnostics.Stopwatch.StartNew();
82         Console.WriteLine(watch.ElapsedMilliseconds);
83         Console.WriteLine(GC.GetTotalMemory(true));
84
85         unchecked
86         {
87             for (int i = 0; i < book.Count; i++)
88             {
89                 if (freeq.ContainsKey(book[i])) freeq[book[i]]++;
90                 else freeq[book[i]] = 1;
91             }
92         }
93         //берем макс слово
94         Console.WriteLine(watch.ElapsedMilliseconds);
95         string maxWord = MaxFreeq.Key;
96         Console.WriteLine(watch.ElapsedMilliseconds);
97         //вытаскиваем все пары
98         SortedDictionary<string, uint> pairs = new SortedDictionary<string, uint>(); //пары слов и их количества
99         for (int i = 1; i < book.Count - 1; ++i)
100         {
101             //отсеиваем то, что не содержит макс слова
102             if (book[i] != maxWord && book[i + 1] != maxWord) continue;
103             //создаем пару
104             string p = book[i] + " " + book[i + 1];
105             //вставка пары и количества пар
106             if (pairs.ContainsKey(p)) pairs[p]++;
107             else pairs[p] = 1;
108
109             /* p = book[i-1] + " " + book[i];
110             if (pairs.ContainsKey(p)) pairs[p]++;
111             else pairs[p] = 1;*/
112         }
113         //создаем список сортированных пар
114         Pairs.Clear();
115         foreach (var v in pairs)
116             Pairs.Add(new WordFreeq(v.Key, v.Value));
117         Console.WriteLine(watch.ElapsedMilliseconds);
118         //сортируем по убыванию встречаемости
119         Pairs.Sort((a, b) => (int)(b.Freeq - a.Freeq));
120         watch.Stop();
121         Console.WriteLine(GC.GetTotalMemory(true));
122         Console.WriteLine(watch.ElapsedMilliseconds);
123         RunTime = watch.ElapsedMilliseconds;
124     }
125 }
126 }

```

## IV.4 Java

### Main.java

```

1  import java.io.*;
2  import java.nio.file.Files;
3  import java.nio.file.Paths;
4  import java.util.*;
5
6  class WordFreeq {
7      public String word;
8      public Integer freeq;
9
10     public WordFreeq( String w, int f ) {
11         word = w;
12         freeq = f;
13     }
14
15     public String toString() {
16         return "[" + word + ", " + freeq + "]";
17     }
18 }
19
20 class KeyValuePair<A, B> {
21     public A first;
22     public B second;
23 }
24
25 public class Main {
26     //Fields
27     public static ArrayList<String> words;
28
29     public static String uniqueWord;
30     public static int wordsCount;
31     public static int uniqueCount;
32
33     public static HashSet<String> wordsToExclude;
34     public static TreeMap<String, Integer> freeq;
35     public static ArrayList<WordFreeq> pairs;
36
37
38     //Methods
39     public static void readExcludeWords( String path ) throws IOException {
40         byte[] bytes = Files.readAllBytes( Paths.get( path ) );
41         String allWords = new String( bytes );
42         wordsToExclude = new HashSet<String>( Arrays.asList( allWords.split( "[\\n]" ) ) );
43     }
44
45     public static void readText( String path ) throws IOException {

```

```

46     byte[] bytes = Files.readAllBytes( Paths.get(path) );
47     String fullText = new String( bytes );
48     words = new ArrayList<String>( Arrays.asList( fullText.split( "[\n]" ) ) );
49 }
50
51 public static void main(String[] args) throws IOException {
52     readExcludeWords( "ignore.txt" );
53     for ( String s : wordsToExclude ) {
54         wordsToExclude.add( s.toLowerCase() );
55     }
56     readText( "Martin Eden.txt" );
57
58     words.removeAll( wordsToExclude );
59
60     long begin = System.currentTimeMillis();
61
62     //Частота встречи слов
63     for ( int i = 0, n = words.size(); i < n; ++i ) {
64         if ( freeq.containsKey( words.get(i) ) ) {
65             Integer k = freeq.get( words.get(i) );
66             freeq.put( words.get( i ), k );
67         }
68         else {
69             freeq.put( words.get( i ), 1 );
70         }
71     }
72
73     String maxWord = "";
74
75     TreeMap<String, Integer> pairsTree = new TreeMap<String, Integer>();
76     for ( int i = 1, n = words.size() - 1; i < n; ++i ) {
77         if ( words.get(i) != maxWord && words.get(i+1) != maxWord )
78             continue;
79
80         String p = words.get(i) + " " + words.get(i+1);
81         if ( pairsTree.containsKey( p ) ) {
82             Integer k = pairsTree.get( p );
83             pairsTree.put( p, k );
84         }
85         else
86             pairsTree.put( p, 1 );
87     }
88
89     pairs.clear();
90     for ( Map.Entry<String, Integer> entry : pairsTree.entrySet() ) {
91         pairs.add( new WordFreeq( entry.getKey(), entry.getValue() ) );
92     }
93
94     Collections.sort( pairs, new Comparator<WordFreeq>() {
95         public int compare( WordFreeq s1, WordFreeq s2 ) {
96             return s2.freeq - s1.freeq;
97         }
98     });
99
100     long end = System.currentTimeMillis() - begin;
101
102     //writeAnswer();
103     System.out.println( "Time " + end );
104
105     //writer.close();
106 }
107 }

```

## V Делаем выводы

Как видно из таблиц (2 – 4), асимптотика, предсказанная в секции (II), работает всегда — независимо от языка программирования (другого мы и не ожидали).

Так же стоит отметить, что C++, будучи языком компилируемым, справляется с задачей быстрее — особенно это заметно на больших тестах (с тестом на 16 000 000 он справился в 5 раз быстрее, чем Java и C# справились с тестом в 4 раза меньшего размера). Это можно объяснить тем, что Java и C# — языки с виртуальной машиной. Интерпретация байткода для машин в исполняемый код занимает своё время. Этим же можно оправдать то, что оба языка показали себя примерно одинаково.

Было замечено, что в случае использования неупорядоченных ассоциативных контейнеров (хэш-таблиц), — `std::unordered_map` и `Dictionary` в STL C++ и C# соответственно, — на небольших тестах C# обогнал C++ приблизительно в 2 раза. Что наводит на мысль, что мейкрософтовская реализация хэш-таблиц в MSVC для STL несколько хуже их же реализации для C#.