

Московский авиационный институт  
(национальный исследовательский университет)

Факультет информационных технологий и прикладной  
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №5 по курсу «Информационный поиск»

Студент: И. А. Мариничев  
Преподаватель: А. А. Кухтичев  
Группа: М8О-408Б-19  
Дата:  
Оценка:  
Подпись:

Москва, 2023

## Лабораторная работа №5 «Поиск цитат, координатный индекс»

В этом задании необходимо расширить язык запросов булева поиска новым элементом – поиском цитат. Синтаксис этого элемента следующий:

- [ «что где когда» ] – кавычки, включают режим цитатного поиска для терминов внутри кавычек. Этому запросу удовлетворяют документы, содержащие в себе все термины что, где и когда, причём они должны встретиться внутри документа ровно в этой последовательности, без каких либо вкраплений других терминов.
- [ «что где когда» / 5 ] – аналогично предыдущему пункту, но допускаются вкрапления других терминов так, чтобы расстояние от первого термина цитаты до последнего не превышало бы 5.

Новый элемент может комбинироваться с другими стандартными средствами булева поиска, например:

- [ «что где когда» && другъ ]
- [ «что где когда» || квн ]
- [ «что где когда» && !«хрустальная сова» ]

Для реализации цитатного поиска нужно использовать координатный индекс, т.е. для каждого вхождения термина в документ построить и сохранить список позиций внутри документа, где этот термин встречался.

В отчёте нужно описать формат координатного индекса. Привести статистические данные:

- Размер получившегося индекса.
- Время построения индекса.
- Общее количество позиций. Среднее количество позиций на термин и на пару термин-документ.
- Скорость индексации (кб входных данных в секунду)
- Время выполнения поисковых запросов.
- Примеры долго выполняющихся запросов.

Кроме того, нужно привести примеры запросов и результаты их выполнения. В выводах должны быть указаны недостатки работы, приведены примеры их решения. Что можно сделать, чтобы ускорить «долгие» запросы?

# 1 Описание

По выбранному корпусу документов строится координатный индекс. Построение индекса занимает 379853 ms (6.33088333 min), сохранение занимает 1.14806e+06 ms (19.13433333 min). Загрузка требует 96074.5 ms (1.601241667 min).

Координатный индекс имеет следующее представление:

```
<Word 0> <Document Id 0> <Coordinates Size 0> <Coordinates 0>
...
<Word 0> <Document Id N> <Coordinates Size N> <Coordinates N>
...
<Word M> <Document Id 0> <Coordinates Size 0> <Coordinates 0>
...
<Word M> <Document Id N> <Coordinates Size N> <Coordinates N>
```

Алгоритм, учитывающий вкрапления других слов в цитатном поиске, был изменён и имеет другой синтаксис. Общий вид запросов: “word1 [/k1] word2 [/k2] ... wordN”. Алгоритм не ищет цитату так, чтобы можно было указать расстояние именно между первым и последним словом. Вместо этого предлагается искать цитату с вкраплениями между каждой последовательной парой слов. Если расстояние не указано, то оно подразумевается равным единице. Т. е. запрос “whale /3 fish dish” будет подразумевать, что между словами “whale” и “fish” может быть до 2 других слова, а после “fish” обязано идти слово “dish”. Цитата так же участвует в булевом поиске и может компоноваться с другими выражениями. В данном случае цитата представляется как отдельный терм с особым алгоритмом поиска документов, где она встречается.

Ниже приведены несколько запросов и результаты по ним.

Запрос	Время (в ms)	Количество найденных файлов
"fellowship of the ring"	390.703	36
"fellowship /5 of /5 the /5 ring"	21.9868	38
"fellowship of the ring"&& gandalf	3.6879	26
"fellowship of the ring"   grinch	6.5681	54
"fellowship of the ring"&& !"peter jackson"	524.674	8

## 2 Исходный код

Ниже приведены основные методы реализующие работу с координатным индексом:

```
1 // index.cpp
2 void Index::SaveCoordinateIndex(std::string &outputFile)
3 {
4     std::wofstream wFileOut((outputFile + "_coordinate").c_str());
5     if (wFileOut)
6     {
7         for (const auto &it : coordinateIndex)
8         {
9             wFileOut << it.first << L" ";
10            wFileOut << it.second.size() << L" ";
11            for (size_t i = 0; i < it.second.size(); ++i)
12            {
13                wFileOut << it.second[i] - 1 << L" ";
14            }
15            wFileOut << L"\n";
16        }
17    }
18    wFileOut.close();
19 }
20
21 void Index::LoadCoordinateIndex(std::string &inputFile)
22 {
23     std::wifstream wFileIn((inputFile + "_coordinate").c_str());
24     std::wstring word, docID;
25     uint32_t size;
26     while (!wFileIn.eof())
27     {
28         std::vector<uint32_t> coordinates;
29         wFileIn >> word >> docID >> size;
30
31         uint32_t coordinate;
32         coordinates.reserve(size);
33         for (size_t i = 0; i < size; ++i)
34         {
35             wFileIn >> coordinate;
36             coordinates.push_back(coordinate);
37         }
38
39         coordinateIndex.insert(std::make_pair((word + L" " + docID), coordinates));
40         word.clear();
41     }
42     wFileIn.close();
43 }
```

Ниже приведены основные методы реализующие работу с цитатным поиском:

```
1 // query.cpp
2 void Query::ProcessingQuote(std::wstring &quote)
3 {
4     std::wstring word = L"";
5     std::wstring wordPrevious;
6     std::stack<std::vector<uint32_t>> result;
7
8     size_t i = 0;
9     size_t k = 1;
10    while (i <= quote.size())
11    {
12        if ((i == quote.size()) or (quote[i] == SPACE))
13        {
14            if ((word.size() != 0) and (result.size() == 0))
15            {
16                wordPrevious = word;
17                result.push(GetDocIndices(word));
18                word = L"";
19            }
20            if (word.size())
21            {
22                std::vector<uint32_t> postings1 = result.top();
23                result.pop();
24                std::vector<uint32_t> postings2 = GetDocIndices(word);
25
26                result.push(IntersectionForQuote(wordPrevious, word, postings1,
27                    postings2, k));
28                k = 1;
29                wordPrevious = word;
30                word = L"";
31            }
32        }
33        else if (quote[i] == SLASH)
34        {
35            k = 0;
36            i++;
37            while ('0' <= quote[i] and quote[i] <= '9')
38            {
39                k *= 10;
40                k += (size_t)(quote[i] - '0');
41                i++;
42            }
43            i--;
44        }
45        else
46        {
47            word += quote[i];
48        }
49    }
```

```

48     i++;
49 }
50
51 auto result_ptr = std::make_shared<std::vector<uint32_t>>(
52     result.size() == 1 ? result.top() : std::vector<uint32_t>());
53 operands.push(result_ptr);
54 }
55
56 std::vector<uint32_t> Query::IntersectionForQuote(std::wstring word1, std::wstring
    word2,
57
58                                     std::vector<uint32_t> &l, std::vector<
    uint32_t> &r,
59                                     size_t k)
60 {
61     std::set<uint32_t> result;
62     size_t i = 0, j = 0;
63     while (i < l.size() && j < r.size())
64     {
65         if (l[i] == r[j])
66         {
67             std::vector<uint32_t> coordinates1 = index.coordinateIndex[word1 + L' ' +
                std::to_wstring(l[i])];
68             std::vector<uint32_t> coordinates2 = index.coordinateIndex[word2 + L' ' +
                std::to_wstring(r[j])];
69
70             size_t ii = 0, jj = 0;
71             while (ii < coordinates1.size())
72             {
73                 while (jj < coordinates2.size())
74                 {
75                     if ((coordinates2[jj] - coordinates1[ii]) > 0 and
76                         (coordinates2[jj] - coordinates1[ii]) <= k)
77                     {
78                         result.insert(l[i]);
79                     }
80                     else if (coordinates2[jj] > coordinates1[ii])
81                     {
82                         break;
83                     }
84                     jj++;
85                 }
86                 ii++;
87             }
88             i++;
89             j++;
90         }
91         else if (l[i] < r[j])
92         {
93             i++;

```

```

93     }
94     else
95     {
96         j++;
97     }
98 }
99
100 return std::vector<uint32_t>(result.begin(), result.end());
101 }

```

### 3 Выводы

Выполнив пятую лабораторную работу по курсу «Информационный поиск», я построил свой цитатный поиск, который использует координатный индекс. Координатный индекс занимает значительно больше памяти в отличие от прямого и обратного, но позволяет сохранить информацию о расположении слов в документах друг относительно друга.



## Список литературы

- [1] Маннинг, Рагхаван, Шютце *Введение в информационный поиск* — Издательский дом «Вильямс», 2011. Перевод с английского: доктор физ.-мат. наук Д. А. Ключина — 528 с. (ISBN 978-5-8459-1623-4 (рус.))