

## ЛР3: Булев индекс

### Задание

Требуется построить поисковый индекс, пригодный для булева поиска, по подготовленному в ЛР1 корпусу документов. Требования к индексу:

- Самостоятельно разработанный, бинарный формат представления данных. Формат необходимо описать в отчёте, в побайтовом представлении.
- Формат должен предполагать расширение, так как в следующих работах он будет меняться под требования новых лабораторных работ.
- Использование текстового представления или готовых баз данных не допускается.
- Кроме обратного индекса, должен быть создан «прямой» индекс, содержащий в себе как минимум заголовки документов и ссылки на них (понадобятся для выполнения ЛР4, при генерации страницы поисковой выдачи).
- Для термов должна быть как минимум понижена капитализация.

В отчёте должно быть отмечено как минимум:

- Выбранное внутренне представление документов после токенизации.
- Выбранный метод сортировки, его достоинства и недостатки для задачи индексации.

Среди результатов и выводов работы нужно указать:

- Количество термов.
- Средняя длина терма. Сравнить со средней длиной токена, вычисленной в ЛР1 по курсу ОТЕЯ. Объяснить причину отличий.
- Скорость индексации: общую, в расчёте на один документ, на килобайт текста.
- Оптимальна ли работа индексации? Что можно ускорить? Каким образом? Чем она ограничена? Что произойдёт, если объём входных данных увеличится в 10 раз, в 100 раз, в 1000 раз?

### Метод решения

1. Выбрать алгоритм индексации
  2. Выбрать бинарный формат представления данных
  3. Реализовать алгоритм
  4. Провести анализ решения
- В качестве алгоритма индексации был выбран алгоритм SPIMI. Этот алгоритм является эффективным с точки зрения масштабирования. SPIMI использует термины, а не их идентификаторы, записывает словарь каждого блока на диск, а затем для нового блока начинает создавать новый словарь. При наличии достаточного объема памяти на диске с помощью алгоритма SPIMI можно проиндексировать коллекцию любого размера.

- Для сохранения индексов в побайтовом представлении был выбран инструмент Boost.Serialization. Данная библиотека позволяет преобразовывать объекты в последовательность байтов, которая может быть сохранена и загружена для восстановления объектов. Доступны различные форматы данных для определения правил генерации последовательностей байтов.

## Результаты выполнения

Процесс индексации коллекции из ~1.8млн документов занимает ~50 минут. В качестве ограничения используемой индексатором RAM был программно указан лимит в размере 100тыс статей на индекс, что соответствует ~550мб.

## Исходный код

```
C indexer.h
1  #ifndef TINDEXER
2  #define TINDEXER
3
4  #include <iostream>
5  #include <unordered_map>
6  #include <vector>
7  #include <map>
8
9
10 class TIndexer
11 {
12 public:
13     TIndexer() {};
14     void BuildIndex(std::string input_articles_path, std::string output_indexes_path, unsigned long long batch_size);
15     void Load(std::string index_path);
16 private:
17     void WriteToIndex(std::string wiki_doc_path);
18     bool InitBuilder(std::string input_articles_path, std::string output_indexes_path, unsigned long long batch_size);
19     void Save(std::string counter);
20     unsigned long long batch_size;
21     std::string input_articles_path;
22     std::string output_indexes_path;
23
24     std::map<std::string, std::vector<unsigned long long>> index;
25 };
26
27 #endif // TINDEXER
```

```

G- indexer.cpp
1  #include <iostream>
2  #include <sstream>
3  #include <chrono>
4  #include <map>
5  #include <filesystem>
6  #include <string>
7  #include <boost/archive/binary_oarchive.hpp>
8  #include <boost/archive/binary_iarchive.hpp>
9  #include <boost/serialization/map.hpp>
10 #include <boost/serialization/vector.hpp>
11 #include <boost/filesystem.hpp>
12 #include <boost/range/iterator_range.hpp>
13
14 #include "nlohmann/json.hpp"
15 #include "indexer.h"
16
17
18 bool FilesystemCreateFolderIdempotent(std::string path) {
19     boost::filesystem::path dir(path);
20
21     if (!(boost::filesystem::exists(dir))) {
22         return boost::filesystem::create_directory(dir);
23     }
24     return true;
25 }
26
27
28 bool TIndexer::InitBuilder(std::string input_articles_path, std::string output_indexes_path, unsigned long long batch_size) {
29     if (!boost::filesystem::exists(input_articles_path)) {
30         std::cout << "input directory not found" << std::endl;
31         return false;
32     }
33
34     if (boost::filesystem::is_empty(input_articles_path)) {
35         std::cout << "input directory is empty" << std::endl;
36         return false;
37     }
38
39     if (!FilesystemCreateFolderIdempotent(output_indexes_path)) {
40         std::cout << "cannot create output indexes dir" << std::endl;
41         return false;
42     }
43
44     this->input_articles_path = input_articles_path;
45     this->output_indexes_path = output_indexes_path;
46     this->batch_size = batch_size;
47     this->index.clear();
48
49     return true;
50 }
51
52

```

```

53 void TIndexer::WriteToIndex(std::string wiki_doc_path) {
54     std::ifstream wiki_article(wiki_doc_path);
55     nlohmann::json wiki_doc = nlohmann::json::parse(wiki_article);
56     std::string doc_id_str = wiki_doc["id"];
57
58     unsigned long long wiki_doc_id = std::stoi(doc_id_str);
59
60     std::string tokenized_text = wiki_doc["text"];
61     std::istringstream token_stream(tokenized_text);
62     std::string token;
63
64     while(token_stream >> token) {
65         auto index_iterator = this->index.find(token);
66         if (index_iterator == this->index.end()) {
67             this->index[token] = std::vector<unsigned long long>{wiki_doc_id};
68         } else {
69             if (index_iterator->second.back() != wiki_doc_id) {
70                 index_iterator->second.emplace_back(wiki_doc_id);
71             }
72         }
73     }
74 }
75
76
77 void TIndexer::BuildIndex(std::string input_articles_path, std::string output_indexes_path, unsigned long long batch_size) {
78     if (!this->InitBuilder(input_articles_path, output_indexes_path, batch_size)) {
79         return;
80     }
81
82     // statistics variables
83     unsigned long long articles_counter = 0;
84     unsigned long long articles_counter_general = 0;
85     std::chrono::steady_clock::time_point time_clock_begin = std::chrono::steady_clock::now();
86     std::chrono::steady_clock::time_point time_clock_end;
87
88
89     for (auto& entry : boost::make_iterator_range(boost::filesystem::directory_iterator(this->input_articles_path), {})) {
90         articles_counter++;
91         articles_counter_general++;
92
93         if (articles_counter < this->batch_size) {
94             this->WriteToIndex(entry.path().string());
95         } else {
96             for (auto& arr : this->index) {
97                 std::sort(arr.second.begin(), arr.second.end());
98             }
99             this->Save(std::to_string(articles_counter_general));
100
101             time_clock_end = std::chrono::steady_clock::now();
102             std::cout << "exec time for batch = " << std::chrono::duration_cast<std::chrono::minutes>(time_clock_end - time_clock_begin).count() << "[m]" << std::endl;

```

```

103
104         articles_counter = 0;
105         time_clock_begin = std::chrono::steady_clock::now();
106         this->index.clear();
107         this->WriteToIndex(entry.path().string());
108     }
109 }
110 if (this->index.size() > 0) {
111     this->Save(std::to_string(articles_counter_general));
112 }
113 }
114
115
116 void TIndexer::Save(std::string postfix)
117 {
118     std::cout << this->output_indexes_path << std::endl;
119
120     std::string index_path = this->output_indexes_path + "/wiki_articles_" + postfix + ".index";
121     std::ofstream f(index_path, std::ios::binary);
122     if (f.fail()) {
123         std::cout << "error" << std::endl;
124         return;
125     }
126
127     boost::archive::binary_oarchive oa(f);
128
129     auto indexValue = this->index;
130     oa << indexValue;
131 }
132
133 void TIndexer::Load(std::string index_path) {
134     this->index.clear();
135     std::ifstream f(index_path, std::ios::binary);
136     if (f.fail()) {
137         std::cout << "error while load data in index" << std::endl;
138         return;
139     }
140     boost::archive::binary_iarchive ia(f);
141     ia >> this->index;
142 }

```

## Выводы:

В результате выполнения данной лабораторной работы нам удалось проиндексировать ~1.8млн документов российской википедии за 50 мин. Процесс оказался достаточно трудоемким в следствие особенностей языка с++. Работать с json оказалось

затруднительно. Понравилось использовать библиотеку boost как для операций с файловой системой, так и для сериализации/десериализации данных.