

Universidade Federal de Minas Gerais

Programação e Desenvolvimento de Software II - (PDSII)

Alunos: André Luiz Alves Costa e Bruno Henrique Evangelista Pereira

Link do repositório: https://github.com/andrelac963/TP_maquina_busca

Trabalho Prático: Máquina de Busca

Introdução

O trabalho prático teve o objetivo de desenvolver uma máquina de busca feita a partir da linguagem de programação C++, onde em um banco de dados de arquivos txt, o usuário irá dizer o que deseja consultar e o programa desenvolvido retornará em quais arquivos as palavras dadas pelo usuário aparecem.

Ao início do projeto nos reunimos para decidir a estrutura da aplicação, a modelagem das classes que seriam necessárias e como seria implementado cada funcionalidade. Para isso foi necessário fazer uma abstração das funcionalidades aplicando os conhecimentos aprendidos na disciplina de PDS 2 e isso auxiliou na organização geral. E com as ideias bem organizadas podemos partir para a implementação do projeto de fato.

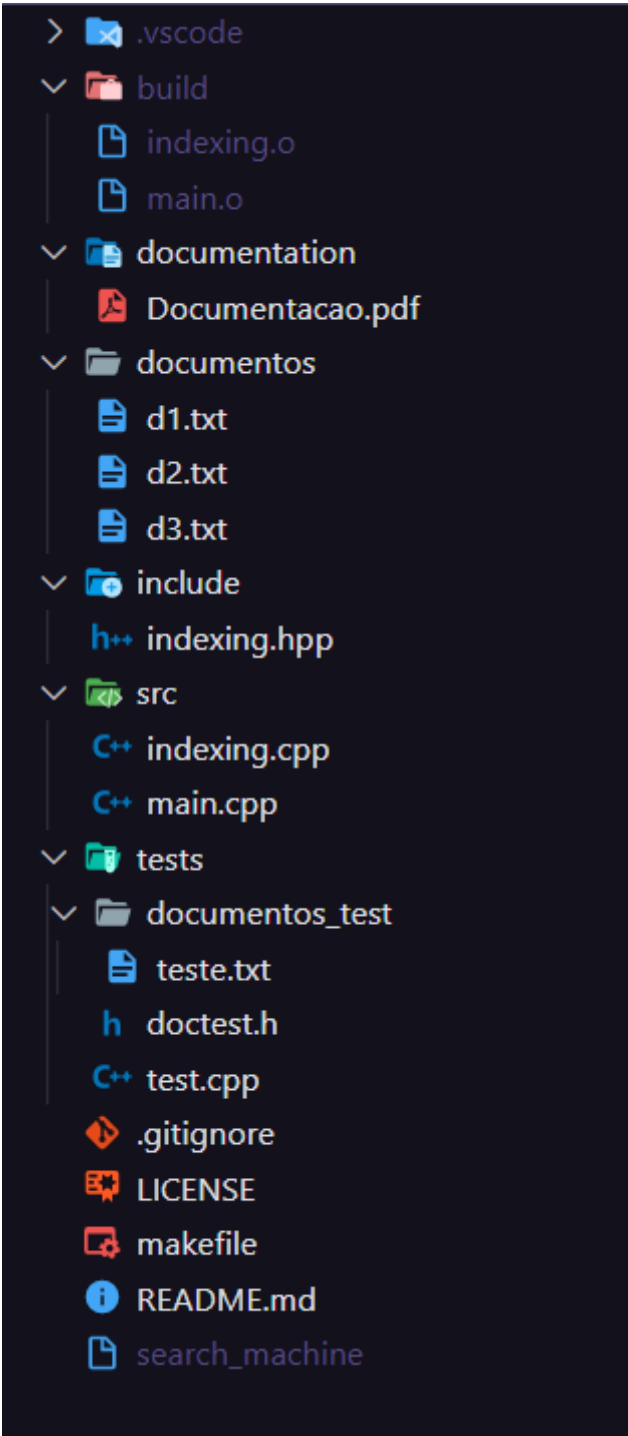
Implementação

Para versionamento do código foi utilizado o github, que funciona baseado em git, e também a utilização de pull requests para que os membros do projeto pudessem validar a implementação dos colegas de trabalho.

Para compilar o código foi utilizado o gcc versão 11.3.0 rodando no WLS com o Ubuntu 22.04.1 LTS. E para facilitar os testes foi feito um makefile com os comandos para compilar todos os arquivos do projeto.

A estrutura de pastas utilizadas foi da seguinte forma:

- build: para guardar os arquivos de compilação;
- documentos: para a base de dados lida pelo programa;
- include: para os arquivos de cabeçalho;
- src: para os arquivos de implementação.
- tests: para realização dos testes de unidade.



```
> .vscode
  build
    indexing.o
    main.o
  documentation
    Documentacao.pdf
  documentos
    d1.txt
    d2.txt
    d3.txt
  include
    indexing.hpp
  src
    indexing.cpp
    main.cpp
  tests
  documentos_test
    teste.txt
    doctest.h
    test.cpp
  .gitignore
  LICENSE
  makefile
  README.md
  search_machine
```

The image shows a file explorer interface with a dark background. It displays a project structure with several folders and files. The folders are expanded, showing their contents. The files are color-coded by extension: .o for object files, .pdf for PDF documents, .txt for text files, .hpp for C++ header files, .cpp for C++ source files, .h for C header files, .gitignore for Git ignore files, LICENSE for license files, makefile for build files, README.md for documentation, and search_machine for the main program.

Decidimos que só seria necessário implementar uma classe indexing onde conteria as funcionalidades de indexação das palavras da base de dados e de recuperação delas.

```
4  #include <map>
5  #include <set>
6  #include <vector>
7  #include <string>
8  #include <iostream>
9  #include <bits/stdc++.h>
10 #include <dirent.h>
11
12 using namespace std;
13
14 class indexing
15 {
16 private:
17     map<string, set<string>> index;
18
19 public:
20     indexing();
21     ~indexing();
22
23     set<string> read_directory(const string &);
24     void read_files(const string &);
25
26     void insert(string, string);
27     map<string, set<string>> get_index();
28
29     string normalize(string);
30
31     void recovery(vector<string>);
32 };
```

Para guardar as informações do index decidimos por utilizar um map responsável por manter a string dada pelo base de dados e o set de strings que são os nomes dos arquivos txt onde ela ocorre. Como há necessidade de que outros arquivos acessem os métodos da classe somente o container

responsável por guardar os dados fique privado. Partindo para implementação do indexing, além do construtor e destrutor de instâncias foi preciso desenvolver métodos para a leitura dos diretórios em que os .txt se encontravam e para ler os arquivos em si .

```
1  #include "indexing.hpp"
2
3  indexing::indexing()
4  {
5  }
6
7  indexing::~~indexing()
8  {
9  }
```

```
11  set<string> indexing::read_directory(const string &name)
12  {
13      DIR *dirp = opendir(name.c_str());
14      struct dirent *dp;
15      set<string> files;
16
17      while ((dp = readdir(dirp)) != NULL)
18      {
19          if (dp->d_name[0] != '.')
20              files.insert(dp->d_name);
21      }
22      closedir(dirp);
23
24      return files;
25  }
```

Como o nome dos arquivos dentro do diretório ./Documentos não é fixo, foi necessário pesquisar uma forma de saber o nome de todos os arquivos dentro de uma pasta, que é o que foi implementado na imagem acima. Utilizando a biblioteca dirent.h, foi possível percorrer o nome de cada arquivo e guardá-los em um container set.

```

27 void indexing::read_files(const string &name)
28 {
29     set<string> files = read_directory(name);
30     fstream file;
31     string word, filename;
32
33     for (auto it = files.begin(); it != files.end(); it++)
34     {
35         filename = *it;
36         file.open(name + filename);
37         while (file >> word)
38         {
39             this->insert(word, filename);
40         }
41         file.close();
42     }
43 }

```

Agora utilizando a biblioteca bits/stdc++, foi possível abrir cada arquivo com o nome salvo no set e inserir cada palavra com seu respectivo nome de arquivo com a função insert.

```

45 void indexing::insert(string word, string filename)
46 {
47     string normalized_word = this->normalize(word);
48     if(normalized_word != "")
49         this->index[normalized_word].insert(filename);
50 }
51
52 map<string, set<string>> indexing::get_index()
53 {
54     return this->index;
55 }

```

A função insert, que recebe a palavra e o nome do arquivo onde ela foi encontrada, primeiro normaliza a palavra e depois a insere no index.

```

57 string indexing::normalize(string word)
58 {
59     string normalized_word = "", aux;
60
61     string comAcentos = "ÃÄÅĀĂāăǻȃĖĘĚēėëîİıİïİİÖŎöŏōõūŪŮůúûüçç";
62     string semAcentos = "AAAAAAaaaaaaEEEEeeeeIIIIIiiiIIOOOOoooooUUUuuuuCc";
63
64     for (int i = 0; i < word.length(); i++)
65     {
66         // Verificando se o caractere é maiúsculo
67         aux = tolower(word[i]);
68
69         // Verificando se o caractere tem acento
70         for (int j = 0; j < comAcentos.length(); j += 2)
71         {
72             if (word[i] << word[i + 1] == comAcentos[j] << comAcentos[j + 1])
73             {
74                 aux = semAcentos[j / 2];
75                 aux = tolower(aux[0]);
76                 break;
77             }
78         }
79
80         // Verificando se o caractere é um caracter especial
81         if (isalpha(aux[0]))
82         {
83             normalized_word += aux;
84         }
85     }
86
87     return normalized_word;
88 }

```

Para permitir uma pesquisa mais precisa utilizamos o método `normalize` que recebe a palavra que será normalizada, e passando caractere por caractere dela, torna letras maiúsculas em minúsculas, através do `tolower()`, retira caracteres especiais, como pontuação, exclamação, vírgula, etc, com a função `isalpha()`.

Fora isso temos que percorrer a string comAcentos comparando se ao caractere da palavra é acentuado, e se sim, substituí-lo por um caractere da string semAcentos. Uma dificuldade na implementação dessa função foi que mesmo que os caracteres acentuados pareçam visualmente ter o mesmo tamanho que um sem acentos, eles ocupam duas casas dentro de um array. Por isso foi necessário fazer a comparação de duas posições juntas da string word com a string comAcentos, e a posição de semAcentos dividido por dois, já que o laço estava sendo percorrido de dois em dois.

```

89 void indexing::recovery(vector<string> query)
90 {
91     set<string> normalized_query;
92     string aux;
93
94     // Normalizando a query
95     for (int i = 0; i < query.size(); i++)
96     {
97         aux = this->normalize(query[i]);
98         if (aux != "")
99         {
100             normalized_query.insert(aux);
101         }
102     }
103
104     // Contando a quantidade de vezes que cada arquivo aparece na query
105     map<string, int> relevant_files;
106
107     for (auto it = this->index.begin(); it != this->index.end(); it++)
108     {
109         if (normalized_query.find(it->first) != normalized_query.end())
110         {
111             for (auto it2 = it->second.begin(); it2 != it->second.end(); it2++)
112             {
113                 relevant_files[*it2]++;
114             }
115         }
116     }

```

Agora para a função de recovery, responsável por retornar os arquivos com a pesquisa feita pelo usuário, precisamos primeiramente normalizar os parâmetros de pesquisa recebidos.

Feito isso, podemos então criar um map responsável por contar quantas vezes a pesquisa do usuário está presente no index. Dentro de um loop que percorre as chaves do index, o programa compara se a palavra do index está presente no container normalized_query, se sim então ele irá percorrer o container set do index, que guarda os arquivos onde a palavra ocorre, e colocar o arquivo no container relevant_files, onde cada vez que ele aparecer na pesquisa irá ser somado mais um ao seu valor.

```

118 if (relevant_files.size() == 0)
119 {
120     cout << "\nNenhum documento encontrado!" << endl;
121 }
122 else
123 {
124
125     // Ordenando os arquivos por ocorrências na query
126     set<pair<int, string>> ordered_files;
127
128     for (auto it = relevant_files.begin(); it != relevant_files.end(); it++)
129     {
130         ordered_files.insert(make_pair(it->second, it->first));
131     }
132
133     cout << "\nDocumentos encontrados:" << endl;
134
135     for (auto it = ordered_files.rbegin(); it != ordered_files.rend(); it++)
136     {
137         cout << it->second << " - " << it->first << " ocorrências" << endl;
138     }
139 }

```

Para que fosse possível ordenar os arquivos por ordem de relevância, foi necessário transferi-los para um container set, onde o valor de int, que representa a quantidade de vezes que aquele arquivo foi contabilizado na pesquisa, seria o parâmetro de ordenação do set. E no fim, imprimindo o container ordered_files teremos esse resultado no terminal:

```
./search_machine
O que deseja buscar?
casa apartamento quem

Documentos encontrados:
d1.txt - 3 ocorrências
d3.txt - 2 ocorrências
d2.txt - 2 ocorrências

Deseja fazer outra pesquisa? (s/n)
█
```

Obs: Não conseguimos implementar a ordenação lexicográfica dos documentos que têm a mesma quantidade de ocorrências, pois o set implementado com pair<int, string> só leva em conta uma das chaves para ordenação, que no caso foi a chave int que representa as ocorrências. Pesquisando encontramos algumas soluções envolvendo bubble sort que vão além dos nossos conhecimentos.

```
7  int main()
8  {
9
10     string option, word, aux;
11     indexing search;
12
13     do
14     {
15         vector<string> query;
16
17         cout << "O que deseja buscar?\n ";
18         getline(cin, word);
19
20         for (int i = 0; i <= word.length(); i++)
21         {
22             if (word[i] == ' ' || i == word.length())
23             {
24                 query.push_back(aux);
25                 aux = "";
26             }
27             else
28             {
29                 aux += word[i];
30             }
31         }
32
33         search.recovery(query);
34
35     do
36     {
37         cout << "\nDeseja fazer outra pesquisa? (s/n)" << endl;
38         cin >> option;
39
40         if (option != "s" && option != "n")
41         {
42             cout << "Opção inválida!" << endl;
43         }
44
45     } while (option != "s" && option != "n");
46
47     if (option == "s")
48     {
49         cout << "\n-----\n"
50             << endl;
51     }
52
53     cin.ignore();
54
55 } while (option == "s");
56
57 return 0;
58 }
```

Por fim, para que o usuário tenha uma interface onde utilizar o programa, no main foi implementado um menu com do/while. O usuário insere a string de pesquisa na variável através do getline, essa string é particionada palavra por

palavra dentro de um vector para facilitar a manipulação, e é chamado a função de recovery passando este vector como parâmetro, retornando o resultado que foi mostrado mais acima. Além disso, o usuário tem a opção de sair do programa ou realizar uma nova consulta.

E com o objetivo de garantir o desempenho correto de todos os métodos do projeto, foi feito uma série de testes de unidade, via doctest.h, colocando diversas formas de entradas para corroboração da corretude do algoritmo desenvolvido.

```
1  #define DOCTEST_CONFIG_IMPLEMENT_WITH_MAIN
2  #include "doctest.h"
3
4  #include "../include/indexing.hpp"
5
6  TEST_CASE("Testa Indexing::Indexing()")
7  {
8      indexing test;
9      CHECK(test.get_index().size() == 0);
10 }
11
12 TEST_CASE("Testa Indexing::read_directory()")
13 {
14     indexing test;
15     set<string> files = test.read_directory("../tests");
16     CHECK(files.size() == 3);
17     CHECK(files.find("test.cpp") != files.end());
18     CHECK(files.find("doctest.h") != files.end());
19     CHECK(files.find("documentos_test") != files.end());
20 }
21
22 TEST_CASE("Testa Indexing::read_files()")
23 {
24     indexing test;
25     test.read_files("../tests/documentos_test/");
26     CHECK(test.get_index().size() == 3);
27 }
28
29 TEST_CASE("Testa Indexing::Insert()")
30 {
31     indexing test;
32     test.insert("teste", "teste.txt");
33     CHECK(test.get_index().size() == 1);
34 }
```

Conclusão

Ao término da aplicação foi possível desenvolver não apenas conhecimentos a respeito da matéria e da linguagem C++, como também o desenvolvimento prático de um projeto.

A priori os fundamentos discutidos em aula foram fundamentais para o desenvolvimento do trabalho como boas práticas e modularização que permitiram uma boa compreensão e organização do código, bem como a programação orientada a objeto e controle de versão que tornaram possível o desenvolvimento do mesmo projeto em paralelo entre os integrantes.

Além dos conceitos discutidos em aula já citados, houve também a busca por informações mais técnicas da ferramenta como por exemplo as funções `isalpha()`, que verifica se o caracter é uma letra, `tolower()`, que transforma letras maiúsculas em minúsculas, além do `getline()`, que foi necessário utilizar a função `cin.ignore()` para que funcionasse mais de uma vez no menu do usuário.

Vale destacar que graças a esse TP descobrimos coisas como uma letra acentuada vale por 2 espaços numa string, por exemplo, a string `comAcento` tem o dobro do tamanho da `semAcento`, presentes no método `normalize`, e como um espaço sozinho gera um caractere inválido, necessitando percorrer esse vetor de duas em duas casas.

Concluindo, foi um projeto para treinar nossos conhecimentos básicos, mas também foi desafiador em certos pontos, o que agregou bastante para o aprendizado na disciplina.