

Universidade Federal de Ouro Preto - UFOP
Instituto de Ciências Exatas e Biológicas - ICEB
Departamento de Computação - DECOM
Ciência da Computação

Hashing para pesquisa de documentos

BCC202 - Estrutura de Dados I

Lucas Chagas, Nicolas Mendes, Pedro Moraes
Professor: Pedro Henrique Lopes Silva

Ouro Preto
17 de março de 2023

Sumário

1	Introdução	1
1.1	Especificações do problema	1
1.2	Considerações iniciais	1
1.3	Ferramentas utilizadas	1
1.4	Especificações da máquina	1
1.5	Instruções de compilação e execução	1
2	Desenvolvimento	2
2.1	Funções de hash	2
2.2	Funções auxiliares	6
2.3	Funções para coleta de informações	8
3	Análise do Experimento e Resultados	9
4	Considerações Finais	9

Lista de Códigos Fonte

1	Código da função inicia indice invertido	2
2	Código da função leDocumento	2
3	Código da função inserePalavraChave	3
4	Código da função insereDocumento	3
5	Código da função busca	4
6	Código da função h	4
7	Código da função consulta	5
8	Código da função vaziosTodosDocumentos	5
9	Código da função selecionaNaoVazio	6
10	Código da função imprimeResultadoBusca	6
11	Código da função ordena	6
12	Código da função leOpção	7
13	Código da função imprime	8
14	Código da função imprimeColisoes	8
15	Parte do código da "main" para o cálculo do tempo de execução	8
16	Código da função imprimeMemoria	9

1 Introdução

O objetivo deste trabalho pratico foi implementar uma tabela hash que ao inves de armazenar quais palavras existem em um documento, armazena quais documentos possuem uma certa palavra.

1.1 Especificações do problema

Nos arquivos proporcionados pelo professor, já havia algumas funções e constantes que definem tamanhos máximos e strings específicas, sendo elas: "VAZIO", "N", "D", "ND", "NM", "M", "MAX STR", "h" e "pegarChaves". Tendo esta base foi exigido que implementássemos uma tabela hash que possui alguma palavra como chave e os documentos que possuem tal palavra dentro desta chave, além disso, os métodos desta tabela deverão permitir ao usuário ver todos os elementos da tabela hash ou buscar quais documentos possuem uma sequência de palavras informada pelo usuário.

1.2 Considerações iniciais

Algumas ferramentas foram utilizadas durante a criação deste projeto:

- Ambiente de desenvolvimento do código fonte: Visual Studio Code. ¹
- Linguagem utilizada: C.
- Ambiente de desenvolvimento da documentação: Overleaf L^AT_EX. ²

1.3 Ferramentas utilizadas

Algumas ferramentas foram utilizadas para testar a implementação, como:

- *CLANG*: ferramentas de análise estática do código.
- *Valgrind*: ferramentas de análise dinâmica do código.

1.4 Especificações da máquina

A máquina onde o desenvolvimento e os testes foram realizados possui a seguinte configuração:

- Processador: Ryzen 7 3700x.
- Memória RAM: 8 Gb.
- Sistema Operacional: Ubuntu.

1.5 Instruções de compilação e execução

Para a compilação do projeto, basta digitar:

```
Compilando o projeto
```

```
make
```

Usou-se para a compilação as seguintes opções:

- *-std=99*: para usar-se o padrão ANSI C 99, conforme exigido.
- *-g*: para compilar com informação de depuração e ser usado pelo Valgrind.
- *-Wall*: para mostrar todos os possível *warnings* do código.

Para a execução do programa basta digitar:

```
./exe < arquivoDeEntrada.in
```

¹Vscode está disponível em <https://code.visualstudio.com/>

²Disponível em <https://www.overleaf.com/>

2 Desenvolvimento

As funções de manipulação da hash implementadas no código(localizadas no arquivo "indiceInvertido.c"), consiste em zerar vetores, inserir valores em vetores, buscar elementos em vetores do tipo abstrato "IndiceInvertido", sendo elas:"inicia", "leDocumentos", "vazioTodosDocumentos", "inserePalavraChave","busca", "insereDocumento","h"e "consulta".

As funções auxiliares do código são responsáveis por imprimir, ler e ordenar vetores de variados tipos, sendo elas: "selecionaNaoVazio", "imprimeResultadoBusca", "ordena", "leOpcao"e "imprime".

2.1 Funções de hash

- "inicia: esta função recebe uma tabela hash do TAD "IndiceInvertido" coloca todas as posições do vetor em um estado vazio,zerando todos os campos "n" e igualando todos os campos "chave" a constante "VAZIO" definida anteriormente no programa.

```
1 void inicia(IndiceInvertido dic)
2 {
3     int i;
4     for (i = 0; i < M; i++)
5     {
6         strcpy(dic[i].chave, VAZIO);
7         dic[i].n = 0;
8     }
9 }
```

Código 1: Código da função inicia indice invertido

- "leDocumento": esta função recebe uma tabela hash do TAD "IndiceInvertido" e a quantidade de elementos que serão colocados na tabela. Com manipulação de strings, a função separa todas as chaves e documentos através de tokenização, após a separação, com o auxílio das funções: "insereDocumento" e "inserePalavraChave", a função insere a chave e seus respectivos documentos em uma posição da tabela.Quanto ao gasto de memória da função,isto é utilizado pelo vetor "str",com tamanho "MAX-STR"(20001) e também pelo vetor "nomeDocumento",com o tamanho "D"(51).

```
1 void leDocumento(IndiceInvertido dic, int n, int *colisoas)
2 {
3     for (int i = 0; i < n; i++)
4     {
5
6         char str[MAX_STR];
7         fgets(str, MAX_STR, stdin);
8
9
10        str[strcspn(str, "\n")] = 0;
11
12
13        char *aux = strtok(str, " ");
14
15
16        char nomeDocumento[D];
17        int qtdPalavras = 0;
18
19
20        while (aux != NULL)
21        {
22            if (qtdPalavras == 0)
23            {
24                strcpy(nomeDocumento, aux);
25                qtdPalavras++;
26            }
27        }
28    }
```

```

27         else
28         {
29
30             inserePalavraChave(dic, aux, colisoes);
31             insereDocumento(dic, aux, nomeDocumento);
32             qtdPalavras++;
33         }
34
35         aux = strtok(NULL, " ");
36     }
37 }
38 }

```

Código 2: Código da função leDocumento

- "inserePalavraChave": esta função recebe uma tabela hash do TAD "IndiceInvertido" e uma "Chave". Com auxílio da função "busca" a função verifica se a chave já existe na tabela e onde há um local válido para a inserção da palavra chave caso ela não esteja lá ainda.

```

1 bool inserePalavraChave(IndiceInvertido dic, Chave chave, int *colisoes)
2 {
3
4     int pos = busca(dic, chave);
5     if (pos != -1)
6     {
7         return false;
8     }
9
10    int j = 0;
11    int ini = h(chave);
12    while ((strcmp(dic[(ini + j) % M].chave, VAZIO) != 0) && (j < M))
13    {
14        j++;
15        *colisoes = *colisoes + 1;
16    }
17
18    if (j < M)
19    {
20        strcpy(dic[(ini + j) % M].chave, chave);
21        return true;
22    }
23
24    return false;
25 }

```

Código 3: Código da função inserePalavraChave

- "insereDocumento": esta função recebe uma tabela hash do TAD "IndiceInvertido", uma "Chave" e uma string de documentos. com auxílio da função "busca", a função acha a posição da chave recebida por parâmetro, verifica se a quantidade de documentos naquela chave já atingiu seu limite e caso não tenha atingido, insere os documentos no campo da chave e incrementa a quantidade de documentos no campo "n".

```

1 bool insereDocumento(IndiceInvertido dic, Chave chave, NomeDocumento
2     doc)
3 {
4     int pos = busca(dic, chave);
5
6     if (dic[pos].n == ND)
7     {
8         return false;

```

```

9      }
10
11      strcpy(dic[pos].documentos[dic[pos].n], doc);
12      dic[pos].n++;
13      memoria++;
14      return true;
15 }

```

Código 4: Código da função insereDocumento

- "busca": esta função recebe uma tabela hash do TAD "IndiceInvertido" e uma "Chave". Usando da função "h" que retorna o local ideal para a chave estar, a função verifica se a posição ideal já está ocupada por outra chave e procura por outra posição disponível ou se a chave já está inserida na tabela. O retorno da função é a posição válida ou a posição da chave caso ela já esteja na tabela.

```

1  int busca(IndiceInvertido dic, Chave chave)
2  {
3      int j = 0;
4      int ini = h(chave);
5
6
7      while ((strcmp(dic[(ini + j) % M].chave, VAZIO) != 0) && (strcmp(dic
      [(ini + j) % M].chave, chave) != 0) && (j < M))
8      {
9          j++;
10     }
11
12     if (strcmp(dic[(ini + j) % M].chave, chave) == 0)
13     {
14         return (ini + j) % M;
15     }
16
17     return -1;
18 }

```

Código 5: Código da função busca

- "h": função implementada para calcular a posição ideal de uma dada chave como parâmetro baseado em pesos definidos dentro da função.

```

1  int h(char * chave) {
2      float p[] = {0.8326030060567271, 0.3224428884580177,
3                  0.6964223353369197, 0.1966079596929834,
4                  0.8949283476433433, 0.4587297824155836,
5                  0.5100785238948532, 0.05356055934904358,
6                  0.9157270141062215, 0.7278472432221632};
7
8      int tamP = 10;
9      unsigned int soma = 0;
10     for (int i=0; i<strlen ( chave ); i++)
11         soma += (unsigned int) chave [i] * p[i % tamP];
12     return soma % M;
13 }

```

Código 6: Código da função h

- "consulta": esta função recebe uma tabela hash do TAD "IndiceInvertido", um vetor de "Chave", a quantidade de chaves no vetor, e um vetor de documentos passado por referência. Em primeiro lugar, a função busca a primeira chave do vetor na tabela hash, caso ela ache, ela coloca todos os documentos da primeira chave no vetor de documentos passado por referência (resultados). Após este passo, a próxima chave do vetor é buscada, e quando achada, os documentos do vetor "resultado" são comparados com os da chave atual, aqueles documentos que pertencerem aos 2

vetores permaneceram no vetor "resultado", e aquelas não pertencerem aos 2, serão substituídos por "VAZIO". este processo é repetido para todas as chaves, fazendo com que no final da execução o vetor "resultado" contenha apenas os documentos que estão presentes em todas as chaves.

```

1  void consulta(IndiceInvertido dic, Chave *chaves, int n,
2      NomeDocumento *resultados)
3  {
4      for (int i = 0; i < n; i++)
5      {
6          int pos = 0;
7          for (int j = 0; j < M; j++)
8          {
9              if (i == 0)
10             {
11                 if (strcmp(dic[j].chave, chaves[i]) == 0)
12                 {
13                     for (int k = 0; k < dic[j].n; k++)
14                     {
15                         strcpy(resultados[pos], dic[j].documentos[k]);
16                         pos++;
17                     }
18                 }
19             }
20             else
21             {
22                 if (strcmp(dic[j].chave, chaves[i]) == 0)
23                 {
24                     for (int l = 0; l < ND; l++)
25                     {
26                         int docExistente = 0;
27                         for (int k = 0; k < dic[j].n; k++)
28                         {
29                             if (strcmp(resultados[l], dic[j].documentos[k]) == 0)
30                             {
31                                 {
32                                     docExistente = 1;
33                                 }
34                             }
35                             if (!docExistente)
36                             {
37                                 strcpy(resultados[l], VAZIO);
38                             }
39                         }
40                     }
41                 }
42             }
43         }
44     }

```

Código 7: Código da função consulta

- "vazioTodosDocumentos": esta função recebe um vetor do TAD "NomeDocumentos" e o preenche com a string "VAZIO".

```

1  void vazioTodosDocumentos(NomeDocumento *documentos)
2  {
3      for (int i = 0; i < ND; i++)
4      {
5          strcpy(documentos[i], VAZIO);
6      }

```

```
7 }
```

Código 8: Código da função vazioTodosDocumentos

2.2 Funções auxiliares

- "selecionaNaoVazio": esta função recebe 2 vetores do TAD "NomeDocumento" e um inteiro representando quantidade. Um dos vetores contém informação e o outro o vetor irá receber apenas os valores válidos(aqueles que não estão preenchidos com "VAZIO").

```
1 void selecionaNaoVazio(NomeDocumento *origem, NomeDocumento *destino,
2 int *qtd)
3 {
4     int pos = 0;
5     for (int i = 0; i < ND; i++)
6     {
7         if (strcmp(origem[i], VAZIO) != 0)
8         {
9             strcpy(destino[pos], origem[i]);
10            pos++;
11        }
12    }
13    *qtd = pos;
14 }
```

Código 9: Código da função selecionaNaoVazio

- "imprimeResultadoBusca": esta função recebe um vetor do TAD "NomeDocumento" e um inteiro representando quantidade de elementos. Quando a quantidade de elementos for 0 a função imprime "none" no terminal, caso contrário a função imprime todo o vetor no terminal.

```
1 void imprimeResultadoBusca(NomeDocumento *documentos, int tamNaoVazio)
2 {
3
4     if (tamNaoVazio == 0)
5     {
6         printf("none\n");
7         return;
8     }
9     else
10    {
11        for (int i = 0; i < tamNaoVazio; i++)
12        {
13            printf("%s\n", documentos[i]);
14        }
15        return;
16    }
17 }
```

Código 10: Código da função imprimeResultadoBusca

- "ordena": esta função é um quicksort adaptado para vetores de strings do TAD "NomeDocumento", que ordena um vetor recebido por parâmetro de maneira crescente, tendo como escolha do pivô, para a ordenação, o elemento inicial do vetor. Por conta do comportamento da tabela hash, a função também verifica os endereços de para evitar erros de overlap. Como o algoritmo em questão é um quicksort, ele acaba gastando $\log(n)$ em espaço, sendo "n" um vetor de documentos com um máximo de 100 documentos sendo cada documento um vetor tipo string com 51 casas.

```
1 void ordena(NomeDocumento *documentos, int inicio, int final)
2 {
3     int pivo, esquerda, direita;
```



```

4     NomeDocumento aux;
5     pivo = inicio;
6     esquerda = inicio;
7     direita = final;
8
9     while (esquerda <= direita)
10    {
11        while ((esquerda < final) && (strcmp(documentos[esquerda],
12            documentos[pivo]) < 0))
13        {
14            esquerda++;
15        }
16        while ((direita > inicio) && (strcmp(documentos[direita],
17            documentos[pivo]) > 0))
18        {
19            direita--;
20        }
21        if (esquerda <= direita)
22        {
23            if (&aux != &documentos[esquerda])
24            {
25                strcpy(aux, documentos[esquerda]);
26            }
27            if (&documentos[esquerda] != &documentos[direita])
28            {
29                strcpy(documentos[esquerda], documentos[direita]);
30            }
31            if (&documentos[direita] != &aux)
32            {
33                strcpy(documentos[direita], aux);
34            }
35            esquerda++;
36            direita--;
37        }
38    }
39    if (direita > inicio)
40    {
41        ordena(documentos, inicio, direita);
42    }
43    if (esquerda < final)
44    {
45        ordena(documentos, esquerda, final);
46    }
47 }

```

Código 11: Código da função ordena

- "leOpcao": esta função recebe uma variável tipo char por referência, um vetor de chaves e uma variável que representa a quantidade de chaves passada por referência. A função extrai o comando do usuário(I ou B) e armazena no char passado por parâmetro caso haja chaves. O custo de memória desta função cai sobre a criação do vetor "str" com tamanho "MAX STR"(20001).

```

1 void leOpcao(char *c, Chave *chaves, int *quantChaves)
2 {
3
4     char str[MAX_STR];
5     fgets(str, MAX_STR, stdin);
6
7
8     str[strcspn(str, "\n")] = 0;
9
10

```

```

11     char *aux = strtok(str, " ");
12
13
14     *c = aux[0];
15
16     aux = strtok(NULL, " ");
17
18     int posChaves = 0;
19
20     while (aux != NULL)
21     {
22         strcpy(chaves[posChaves], aux);
23         aux = strtok(NULL, " ");
24         posChaves++;
25     }
26
27     *quantChaves = posChaves;
28 }

```

Código 12: Código da função leOpção

- "imprime": esta função recebe uma tabela hash do TAD "IndiceInvertido". A função imprime todas as chaves válidas e seus respectivos documentos no terminal.

```

1 void imprime(IndiceInvertido dic)
2 {
3     for (int i = 0; i < M; i++)
4     {
5         if (strcmp(dic[i].chave, VAZIO) != 0)
6         {
7             printf("%s -", dic[i].chave);
8             for (int j = 0; j < dic[i].n; j++)
9             {
10                 printf(" %s", dic[i].documentos[j]);
11             }
12             printf("\n");
13         }
14     }
15 }

```

Código 13: Código da função imprime

2.3 Funções para coleta de informações

- "imprimeColisoas": esta função irá imprimir a quantidade de colisões ocorrida na execução de cada caso de teste no código. A função é estruturada a partir de uma variável inicializada na main e é passada por referência para a função "busca" e cada colisão é incrementado 1 unidade para a variável.

```

1 void imprimeColisoas(int colisoas)
2 {
3     printf("Quantidade de colisoas: %d\n", colisoas);
4 }

```

Código 14: Código da função imprimeColisoas

- "tempoDeExecução": esta funcionalidade implementada na main, irá calcular o tempo de execução de cada caso de teste feito no código. O cálculo do tempo é feito pela subtração do final pelo início, dividido pelo "CLOCKS-PER-SEC".

```

1     double tempo_gasto = 0.0;
2     clock_t begin = clock();

```

```

3
4     clock_t end = clock();
5     tempo_gasto += (double) (end - begin) / CLOCKS_PER_SEC;
6     printf("Tempo gasto %f s\n", tempo_gasto);

```

Código 15: Parte do código da "main" para o cálculo do tempo de execução

- "imprimeMemoria": esta irá imprimir a quantidade de memória gasta em cada caso de teste executado no código, à vista disso, o cálculo do gasto de memória será igual ao tamanho do vetor multiplicado pelos bytes de cada elemento do vetor. Sendo o primeiro vetor "Chaves chave" utilizado na função "leOpcao" e na função "consulta", o vetor "IndiceInvertido dic" utilizado na maioria das funções, o vetor auxiliar para o nome do documento e vetor de string utilizado na função "leDocumento", o vetor "NomeDocumento resultados" utilizado na função de "consulta" e na "selecionaNaoVazios" e o vetor "NomeDocumento naoVazios" utilizado na função de ordenação(ordena) e na função de "imprime".

```

1 void imprimeMemoria()
2 {
3     int memoria = (M * sizeof(Item)) + (D * sizeof(char)) + (MAX_STR *
4                     sizeof(char)) + (ND * sizeof(NomeDocumento)) + (ND *
5                     sizeof(NomeDocumento)) + (M * sizeof(Chave));
6     printf("Memoria gasta: %d bytes\n", memoria);
7 }

```

Código 16: Código da função imprimeMemoria

3 Análise do Experimento e Resultados

Foram realizados vários testes com diversos tipos de input, com isso conseguimos corrigir alguns erros ao longo do processo, tendo entraves principalmente na função "consulta", tendo em vista que há praticamente toda a lógica do trabalho prático nesta função e também na função "imprime", também houve problemas com o método de ordenação, que deveria realizar todo o processo em menos de 1 segundo, e o erro "overlap" apontado pelo valgrind, nos levando a tomar certas ações para evitar tal problema, o método de ordenação usado foi o quicksort devido sua simples implementação e complexidade adequada, e o que foi implementado para evitar o erro de overlap foi verificar dentro do próprio quicksort se as posições analisadas eram válidas. Com isso, conclui-se que o código implementado cumpriu todos os requisitos propostos e também todos os casos de teste, sendo assim mostrado que cada tópico da disciplina foi aplicado de maneira produtiva e eficiente para atingir o melhor funcionamento do código e alcançar a consolidação do conteúdo entre os membros do grupo.

4 Considerações Finais

O código foi capaz de lidar com todos os casos de teste propostos no "runcodes" e não apresentou nenhum erro na verificação do valgrind. O código ficou suficientemente intuitivo, sendo comentado e indentado em toda sua constituição.