

### Trabalho Prático I de CCF212 - 20 pts

#### Estudo comparativo entre Árvore PATRICIA e Tabela HASH como estruturas para implementar arquivo invertido

**Formação do grupo:** O trabalho deverá ser desenvolvido em **10 grupos** de **5** integrantes e **1 grupo** com **4** integrantes. Sugere-se que os grupos sejam organizados dentro da divisão de turmas práticas (não obrigatório).

**Problema:** Construção de índice invertido para uma máquina de busca dos POCs do curso de Ciência da Computação. Atualmente, os POCs estão disponíveis como arquivos, em formato pdf, organizados no site do curso (<https://ccp.caf.ufv.br/tccs/>) por semestre de apresentação. Ainda não existe nenhum tipo de sistema de busca que permita a aplicação de filtros para recuperação dos trabalhos.

**Fundamentação Teórica:** Máquinas de busca, tais como Google, trabalham com a busca de palavras-chave em textos armazenados na Web. Para que os documentos contendo os termos sejam recuperados, os mesmos precisam ser devidamente indexados à priori. Nesse contexto, são utilizadas estruturas de dados que facilitem a recuperação das informações, como é o caso do uso de arquivos invertidos. Dada uma coleção de documentos, um índice invertido é uma estrutura contendo uma entrada para cada palavra (termo de busca) que aparece em, pelo menos, um dos documentos. Essa entrada associa a cada palavra do texto um ou mais pares do tipo  $\langle qtde, idDoc \rangle$ , onde  $qtde$  corresponde ao número de vezes em que a palavra em questão apareceu em um determinado documento identificado por  $idDoc$ . O nome dado à estrutura indica que houve uma inversão da hierarquia da informação, ou seja, ao invés de uma lista de documentos contendo termos, é obtida uma lista de termos, referenciando documentos. Essa estrutura de índices é comumente implementada com base em *árvores* e tabelas *hash*, pois as mesmas não precisam ser reconstruídas a cada atualização.

Para exemplificar, considere os dois documentos mostrados abaixo<sup>1</sup>:

Texto 1 (arquivo1.txt) = *"Quem casa quer casa. Porem ninguem casa. Ninguem quer casa tambem. Quer apartamento."*

Texto 2 (arquivo2.txt) = *"Ninguem em casa. Todos sairam. Todos. Quer entrar? Quem? Quem?"*

<sup>1</sup> Nos exemplos, os textos não têm acentuação. Além disso, palavras com letras maiúsculas foram transformadas para minúsculas antes da inserção no índice.

Supondo os identificadores 1 e 2 para os textos apresentados, arquivo1.txt e arquivo2.txt, respectivamente, o índice invertido para as palavras contidas nos textos pode ser visualizado na Tabela 1.

**Tabela 1** - Índices invertidos para os textos apresentados

Palavra	<qtde, idDoc>	<qtde, idDoc>	...
apartamento	<1,1>		
casa	<4, 1>	<1, 2>	
em	<1, 2>		
entrar	<1, 2>		
ninguem	<2, 1>	<1, 2>	
porem	<1, 1>		
quem	<1, 1>	<2, 2>	
quer	<3, 1>	<1, 2>	
sairam	<1, 2>		
tambem	<1, 1>		
todos	<2, 2>		

Conforme pode ser observado na Tabela 1, para cada palavra, existe uma lista de pares de números *<qtde, idDoc>*, ordenada pelo campo *idDoc*. Em termos práticos, essa lista poderia ser implementada como uma lista encadeada para cada palavra indexada.

A partir dos textos indexados no arquivo invertido, podem ser realizadas pesquisas, com base em termos de busca, para se calcular a relevância de documentos. Neste trabalho, será utilizado um método de ponderação baseado no cálculo da frequência da ocorrência do(s) termo(s) nos documentos, conhecida como **TF-IDF (Term frequency - Inverse Document Frequency)**. Esse método é baseado na frequência dos termos da consulta em cada documento (TF) da coleção bem como na frequência inversa dos documentos (IDF). O componente IDF estima o quanto um termo ajuda a discriminar os documentos entre relevantes e não relevantes. Um termo que aparece em muitos documentos têm valor de IDF baixo, enquanto um termo que aparece em poucos documentos apresenta IDF alto, sendo um bom discriminador.

Dada uma consulta com  $q$  termos,  $t_1, t_2, \dots, t_q$ , a relevância de um documento  $i$ ,  $r(i)$ , é computada como:

$r(i) = \frac{1}{n_i} \sum_{j=1}^q w_j^i$	$n_i$ = número de termos distintos do documento $i$ $w_{j,i}$ = peso do termo $t_j$ no documento $i$
$w_j^i = f_{j,i} \frac{\log(N)}{d_j}$	$f_{j,i}$ = número de ocorrências do termo $t_j$ no documento $i$ $d_j$ = número de documentos na coleção que contém o termo $t_j$ $N$ = número de documentos na coleção. Se o termo $t_j$ não aparece no documento $i$ , $f_{j,i} = 0$ * O log é na base 2!

Para exemplificar, considere uma consulta com os termos “quer” e “todos”<sup>2</sup>:

- dois termos ( $q = 2$ )
- dois documentos:  $N = 2$ 
  - o documento 1 tem 7 termos ( $n_1 = 7$ )
  - o documento 2 tem 8 termos ( $n_2 = 8$ )
- O número de ocorrências do primeiro termo (quer), no documento 1 é 3 e no documento 2 é 1, logo  $f_{1,1} = 3$  e  $f_{1,2} = 1$ . Além disto  $d_1 = 2$ .
  - $w_{1,1} = 3 * \log_2 2 / 2 = 1.5$
  - $w_{1,2} = 1 * \log_2 2 / 2 = 0.5$
- O número de ocorrências do segundo termo (todos), no documento 1 é 0 e no documento 2 é 2, logo  $f_{2,1} = 0$ ,  $f_{2,2} = 2$  e  $d_2 = 2$ 
  - $w_{2,1} = 0 * \log_2 2 / 2 = 0$
  - $w_{2,2} = 2 * \log_2 2 / 2 = 1$
- Logo, as relevâncias dos documentos para esta consulta são:
  - $r(1) = 1/7 * (1.5 + 0) = 0.21$
  - $r(2) = 1/8 * (0.5 + 2) = 0.31$

A partir das relevâncias calculadas para cada documento, o método retornará os documentos ordenados da seguinte forma:

Texto 2 (arquivo2.txt)  
Texto 1 (arquivo1.txt)

### Tarefas:

1. **Receber, como entrada,** o arquivo **entrada.txt**, com  $N$  arquivos com textos (POCs) cujas palavras serão indexadas. O arquivo deverá conter o seguinte formato:

$N$   
arquivo1.txt  
arquivo2.txt  
arquivo3.txt  
.....  
arquivoN.txt

A primeira linha deste arquivo contém um número  $N$  que representa o número de documentos da coleção. Cada linha a seguir contém o nome do arquivo que contém um dos documentos da coleção. No exemplo acima, há  $N$  documentos que estão armazenados nos arquivos arquivo1.txt, arquivo2.txt, ....., arquivoN.txt. Você pode assumir que estes arquivos, caso existam (você precisa testar), estarão no diretório corrente de execução.

O sistema deverá processar cada um dos arquivos, lendo palavra após palavra e construindo os índices invertidos, um em cada uma das estruturas (PATRICIA e HASH). Ele deverá também associar a cada documento um *idDoc* único e associar, em memória, este identificador com o nome do documento.

Para simplificar e reduzir o volume de informações manipuladas, podem ser indexados apenas Título e Resumo dos trabalhos (inclua também palavras-chave, se houver). Sugere-se que os textos a serem indexados sejam copiados manualmente e armazenados localmente em um formato de manipulação mais fácil (txt, csv etc). Também para simplificar, sugere-se que sejam utilizados os trabalhos em inglês para evitar a manipulação de acentos.

2. **Criar um Índice Invertido** para cada um dos TADs: **PATRICIA e HASH**. No caso do TAD PATRICIA, você deve adaptar os algoritmos fornecidos em sala de aula para permitir o armazenamento de palavras. A solução mais comum é inserir mais um campo de comparação em cada nó, ou seja, além do campo de índice (que avança  $x$  posições na palavra) será necessário também ter um campo com o caracter que está sendo comparado naquela posição para se decidir o caminho a seguir (esquerda ou direita). A decisão de se colocar, no nó interno, o menor ou o maior caracter de comparação e se os iguais ficarão à esquerda ou à direita deste nó, deverá levar em conta o melhor uso de memória e a

diferença de tamanho entre as palavras. Essa decisão é um componente importante no trabalho prático e precisa ser tomada com atenção! No caso do TAD HASH, deve ser usado o endereçamento com encadeamento, ou seja, hash usando listas encadeadas para tratar colisões. No caso da tabela HASH, o tamanho da tabela deve ser definido considerando-se um fator de carga definido a priori.

3. **Imprimir os índices invertidos de cada TAD:** imprime as palavras da **PATRICIA** e da **HASH**, em ordem alfabética, uma por linha, seguidas da lista das ocorrências, ordenada pelo índice do documento (pares `<qtde, idDoc>`).
4. **Implementar uma função de busca por POCs com base em termo(s) de busca**, utilizando, **individualmente**, cada um dos índices invertidos para localizar os textos com os POCs em que ele(s) aparece(m). Os textos devem ser retornados de forma ordenada pela sua relevância para a consulta, ou seja, textos que contém um maior número de ocorrências do(s) termo(s) de busca devem aparecer primeiro.
5. **Implementar um menu com as seguintes opções:** a) receber o arquivo de entrada com os textos a serem indexados; b) construir os índices invertidos, a partir dos textos de entrada, usando os TADs PATRICIA e HASH; c) imprimir os índices invertidos, contendo as palavras em ordem alfabética, uma por linha, com suas respectivas listas de ocorrências; d) realizar buscas por um ou mais termo(s) de busca, nos índices construídos, individualmente, apresentando os arquivos ordenados por relevância, também individualmente para cada TAD.
6. **Realizar medições comparativas para as operações de inserção e pesquisa em cada um dos TADs:** a) número de comparações numa operação de inserção de um termo, usando PATRICIA e HASH (operação usada na criação dos arquivos invertidos); e b) número de comparações numa operação de pesquisa de termos de busca, usando PATRICIA e HASH (operação usada na busca).
7. **Elaborar um relatório sintético**, contendo:
  - Capa com título e formação do grupo;
  - Introdução, com o objetivo do trabalho e as principais fontes (referências) dos algoritmos utilizados;
  - Metodologia, documentando como o grupo se organizou em termos da divisão e da execução das tarefas e compartilhamento de código.
  - Resultados do desenvolvimento, apresentando: a) detalhes técnicos de implementação da PATRICIA, para armazenar palavras, e da HASH; e b) resultados dos testes comparativos, devidamente fundamentados na teoria vista em sala de aula.
  - Considerações finais, em que o grupo pode registrar as principais

facilidades, dificuldades e lições aprendidas; e

- Referências bibliográficas utilizadas.

### **Entrega:**

- ✓ Entrega via PVANET Moodle, por **um** dos integrantes do grupo, através de um **único** arquivo compactado (em formato .zip e nomeado com nome do grupo), contendo:
  - o código-fonte do programa em C, com implementações dos TADs e manipulação dos textos de entrada em arquivos separados;
  - o arquivos utilizados como entrada;
  - o arquivo Makefile (com utilitário make) para compilação e geração de código executável automaticamente. Será fornecido um tutorial!
  - o arquivo "leiam.txt" com explicações de uso e execução do programa; e
  - o relatório em formato pdf.
- ✓ Data de entrega: **30/06/25**
- ✓ Data de apresentação/entrevista: **01 a 03/07/25, em cronograma a ser divulgado.**

### **Comentários Gerais:**

- O código-fonte DEVERÁ ser devidamente comentado;
- As implementações relativas a cada TAD devem estar em arquivos separados;
- O código para leitura e carga dos dados devem estar em um arquivo separado;
- As operações referentes à montagem do índice invertido devem estar em um arquivo separado, um para cada TAD (PATRICIA e HASH);
- Atenção quanto ao uso e inicializações de variáveis no programa principal, que podem comprometer o funcionamento do seu código (Encapsular funções sempre que possível!);
- O grupo deverá estar identificado no cabeçalho de TODOS os arquivos do código-fonte;
- Apesar de o trabalho ser em grupo, a nota poderá ser individualmente atribuída, a critério da professora (entrevistas individuais poderão ser realizadas);
- Em caso de plágio entre trabalhos, será atribuída nota **zero** para todos

os envolvidos (dos grupos em questão) e atribuição de conceito **F**. Se houver discussões entre os grupos acerca de soluções para questões específicas dos algoritmos, isso deve estar documentado no relatório e no código-fonte (nos trechos correspondentes);

- Trabalhos entregues **em atraso** ou que **não sejam apresentados** receberão nota ZERO (A participação na entrevista é obrigatória para TODOS os integrantes);
- Durante o desenvolvimento do trabalho, caberá ao grupo propor e construir uma implementação para o problema apresentado. A professora e os monitores não irão resolver erros em código-fonte nem tampouco fornecer detalhes técnicos da solução a ser construída.