

**Universidade de São Paulo
Instituto de Ciências Matemáticas e de Computação
Departamento de Ciências de Computação
Disciplina de Organização de Arquivos**

docente

Profª. Dra. Cristina Dutra de Aguiar Ciferri (cdac@icmc.usp.br)

alunos PAE

Turma A. Viviana Elizabeth Romero Noguera (viviana.noguera@usp.br)

Turma B. Guilherme Muzzi da Rocha (guilherme.muzzi.rocha@usp.br)

monitores

Turmas A e B. Matheus Carvalho Raimundo (mcarvalhor@usp.br)

Turma B. Gabriel Alfonso Nascimento Salgueiro (gabrielsalgueiro@usp.br)

Quarto Trabalho Prático

Este trabalho tem como objetivo indexar arquivos de dados usando um índice simples ou linear.

*O trabalho deve ser feito em **dupla** ou **individualmente**. A solução deve ser proposta exclusivamente pelo aluno com base nos conhecimentos adquiridos nas aulas. Consulte as notas de aula e o livro texto quando necessário. A dupla deve escolher um dos trabalhos anteriores a ser utilizado para desenvolver o trabalho prático 4. Isso deve ser informado no início do código do trabalho prático 4 como um comentário.*

Descrição de páginas de disco

No trabalho será usado o conceito de páginas de disco. Cada página de disco tem o tamanho fixo de 32.000 bytes. O conceito de página de disco é um conceito lógico, ou seja, deve ser garantido via programação, de forma que cada página de disco contenha, no máximo, o tamanho fixo especificado.

Descrição do arquivo de índice simples ou linear

Descrição do Registro de Cabeçalho. O registro de cabeçalho deve conter o seguinte campo:

- *status*: indica a consistência do arquivo de índice, devido à queda de energia, travamento do programa, etc. Pode assumir os valores 0, para indicar que o arquivo de dados está inconsistente, ou 1, para indicar que o arquivo de dados está consistente. Ao se abrir um arquivo para escrita, seu status deve ser 0 e, ao finalizar o uso desse arquivo, seu status deve ser 1 – tamanho: *string* de 1 byte.

Representação Gráfica do Registro de Cabeçalho. O tamanho do registro de cabeçalho deve ser de 1 byte, representado da seguinte forma:

1 byte
<i>status</i>

Página de disco. O registro de cabeçalho deve ocupar uma página de disco. Seu tamanho é menor do que o tamanho da página de disco. Neste caso, a página de disco deve ser preenchida com caractere '@' até completar o seu tamanho.

Descrição do Registro de Dados. Devem ser considerados campos de tamanho fixo e registros de tamanho fixo. Os registros de dados devem conter os seguintes campos:

- *chaveBusca*: armazena a chave de busca do índice - tamanho: *string* de 120 bytes.
- *byteOffset*: armazena o *byte offset* do registro que corresponde à chave de busca associada – tamanho: inteiro de 8 bytes.

Representação Gráfica do Registro de Dados. O tamanho de cada registro de dados é de 128 bytes, representado da seguinte forma:

50 bytes			8 bytes
<i>chaveBusca</i>			<i>byteOffset</i>
0	1	127

Observações Importantes.

- Cada registro de dados deve seguir estritamente a ordem definida na sua representação gráfica.

- Os nomes dos atributos também devem seguir estritamente os nomes definidos na especificação dos mesmos.

Programa

Descrição Geral. Implemente um programa em C que ofereça uma interface por meio da qual o usuário possa realizar operações sobre um índice simples ou linear e que possa inserir, remover e atualizar dados de um arquivo de dados que está indexado por esse índice simples ou linear.

Importante. A definição da sintaxe de cada comando bem como sua saída devem seguir estritamente as especificações definidas em cada funcionalidade. Para especificar a sintaxe de execução, considere que o programa seja chamado de “programaTrab4”. Essas orientações devem ser seguidas uma vez que a correção do funcionamento do programa se dará de forma automática.

Descrição Específica. O programa deve oferecer as seguintes funcionalidades:

[10] Crie um arquivo de índice secundário fortemente ligado para um arquivo de dados de entrada já existente. O campo a ser indexado é *nomeServidor*. A implementação do índice não deve ser na forma de lista invertida, ou seja, os valores repetidos de *nomeServidor* devem aparecer várias vezes no arquivo de dados do índice. Registros logicamente removidos presentes no arquivo de dados de entrada não devem ter suas chaves de busca correspondentes no arquivo de índice. Registros com nomes de servidores nulos também não devem ter suas chaves de busca correspondentes no arquivo de índice. No arquivo de índice, as chaves de busca devem seguir a ordenação alfabética, desde que as chaves de busca são do tipo *string*. Adicionalmente, quando duas ou mais chaves de busca possuírem o mesmo valor para o campo *nomeServidor*, os valores dos *byte offsets* correspondentes devem ser ordenados de forma crescente. A ordenação deve ser numérica, desde que os *byte offsets* são do tipo inteiro. A implementação dessa funcionalidade deve ser feita de forma eficiente. Antes de terminar a execução da funcionalidade, deve ser utilizada a função `binarioNaTela1` ou `binarioNaTela2`, ambas disponibilizadas na página do projeto da disciplina, para mostrar a saída do arquivo de índice secundário.

Entrada do programa para a funcionalidade [10]:

```
10 arquivoEntrada.bin arquivoIndiceNomeServidor.bin
```

onde:

- `arquivoEntrada.bin` é um arquivo binário de entrada que segue as mesmas especificações do primeiro trabalho prático, e que contém dados desordenados e registros logicamente removidos.
- `arquivoIndiceNomeServidor.bin` é o arquivo binário de índice secundário fortemente ligado que indexa o campo *nomeServidor*. Esse arquivo deve seguir as especificações definidas neste trabalho prático.

Saída caso o programa seja executado com sucesso:

```
Listar o arquivo binário arquivoIndice.bin.
```

Mensagem de saída caso algum erro seja encontrado:

```
Falha no processamento do arquivo.
```

Exemplo de execução:

```
./programaTrab4  
10 arquivo.bin arquivoIndiceNomeServidor.bin  
usar a função binarioNaTela1 ou binarioNaTela2 antes de terminar a  
execução da funcionalidade, para mostrar a saída do arquivo  
arquivoIndiceNomeServidor.bin, o qual indexa o campo nomeServidor.
```

[11] Permita a recuperação dos dados de todos os registros que satisfaçam um critério de busca determinado pelo usuário sobre o campo *nomeServidor*, usando o índice secundário fortemente ligado criado na funcionalidade [10]. Note que somente o campo *nomeServidor* deve ser utilizado como forma de busca, desde que o índice criado na funcionalidade [10] indexa chaves de busca desse campo. Essa funcionalidade pode retornar 0 registros (quando nenhum satisfaz ao critério de busca), 1 registro (quando não existem valores repetidos do nome do servidor) ou vários registros (quando existem valores repetidos do nome do servidor). Antes de executar essa funcionalidade, realize a operação de carregamento do índice secundário do disco para a memória primária. Os dados solicitados devem ser mostrados na saída padrão da seguinte forma. Para cada registro, mostre os metadados, seguidos por dois pontos (:), seguidos de um espaço em branco, seguidos pelos valores de seus campos, de forma que cada campo apareça em uma linha diferente. Os metadados devem ser lidos do registro de cabeçalho (campos *desCampo1*, *desCampo2*, *desCampo3*, *desCampo4*, *desCampo5*). Caso o valor seja nulo, escreva 'valor nao declarado'. Depois do registro, deve-se pular uma linha em branco. Depois de mostrar todos os registros, deve ser mostrado na saída padrão o número de páginas de disco acessadas, especificadas da seguinte forma: (i) número de acessos a disco na operação de carregamento do índice secundário; e (ii) número de acessos a disco para acessar o arquivo de dados.

Sintaxe do comando para a funcionalidade [11]:

```
l1 arquivo.bin arquivoIndiceNomeServidor.bin nomeServidor valor
```

Saída caso o programa seja executado com sucesso:

Para cada registro, devem ser exibidos seus campos da seguinte forma. Cada campo deve ser exibido em uma linha diferente. Primeiro, deve ser colocado o valor do metadado para aquele campo, e depois o seu valor. Campos com valores nulos devem ser representados por 'valor nao declarado'. Depois de cada registro, pule uma linha em branco. Especifique também o número de páginas de disco para carregar o arquivo de índice e o número de páginas de disco para acessar o arquivo de dados.

Mensagem de saída caso não seja encontrado o registro que contém o valor do campo ou o campo pertence a um registro que esteja removido:

Registro inexistente.

Mensagem de saída caso algum erro seja encontrado:

Falha no processamento do arquivo.

Exemplo de execução (considere que os metadados lidos do arquivo .csv sejam: numero de identificacao do servidor, salario do servidor, telefone celular do servidor, nome do servidor, cargo do servidor):

```
./programaTrab4  
l1 arquivo.bin nomeServidor CRISTIANO ANDRE DA SILVA  
numero de identificacao do servidor: 6202963  
salario do servidor: 4097.91  
telefone celular do servidor: (60)99435-1564  
nome do servidor: CRISTIANO ANDRE DA SILVA  
cargo do servidor: AGENTE ADMINISTRATIVO  
--- pular uma linha em branco ----  
Número de páginas de disco para carregar o arquivo de índice: 5  
Número de páginas de disco para acessar o arquivo de dados: 1
```

[12] Estenda a funcionalidade [4] descrita no segundo trabalho prático de forma que, depois de cada remoção lógica no arquivo de dados, a chave de busca referente ao registro logicamente removido seja removida do índice secundário fortemente ligado criado na funcionalidade [10]. A especificação detalhada da funcionalidade [12] é a mesma da funcionalidade [4]. Com relação à manipulação do índice secundário, antes de executar a funcionalidade [12], realize a operação de carregamento do índice do disco para a memória primária. Enquanto forem realizadas as n remoções, o índice deve ser atualizado somente em memória primária. Ao finalizar as n remoções, realize a operação de reescrita do índice da memória primária para o disco. Nessa operação de reescrita, use a opção wb para reescrever no arquivo.

Entrada do programa para a funcionalidade [12]:

```
12 arquivo.bin arquivoIndiceNomeServidor.bin n
nomeCampo1 valorCampo1
nomeCampo2 valorCampo2
...
nomeCampon valorCampon
```

onde:

- arquivo.bin é um arquivo binário de entrada que segue as mesmas especificações definidas no primeiro trabalho prático. As remoções a serem realizadas nessa funcionalidade devem ser feitas nesse arquivo.

- arquivoIndiceNomeServidor.bin é o arquivo binário referente ao índice secundário fortemente ligado que segue as especificações definidas nesse trabalho prático.

- n é o número de remoções a serem realizadas, as quais seguem as mesmas especificações definidas na funcionalidade [4] do segundo trabalho prático. Ou seja, para cada remoção, deve ser informado o nome do campo a ser considerado e seu critério de busca representado pelo valor do campo. Cada uma das n remoções deve ser especificada em uma linha diferente. Deve ser deixado um espaço em branco entre o nome do campo e o valor do campo. Os valores dos campos do tipo *string* devem ser especificados entre aspas duplas (").

Saída caso o programa seja executado com sucesso:

Listar o arquivo de índice secundário arquivoIndiceNomeServidor.bin.

Mensagem de saída caso algum erro seja encontrado:

Falha no processamento do arquivo.

Exemplo de execução:

```
./programaTrab4
12 arquivo.bin arquivoIndiceNomeServidor.bin 2
idServidor 1914208
cargoServidor "TECNICO EM SAUDE PUBLICA"
usar a função binarioNaTela1 ou binarioNaTela2 antes de terminar a
execução da funcionalidade, para mostrar a saída do arquivo
arquivoIndiceNomeServidor.bin, o qual foi atualizado frente às
remoções realizadas no arquivo de dados.
```


[13] Estenda a funcionalidade [5] descrita no segundo trabalho prático de forma que, depois de cada inserção de registro adicional no arquivo de dados, a chave de busca referente ao registro inserido seja inserida no índice secundário fortemente ligado criado na funcionalidade [10]. A especificação detalhada da funcionalidade [13] é a mesma da funcionalidade [5]. Com relação à manipulação do índice secundário, antes de executar a funcionalidade [13], realize a operação de carregamento do índice do disco para a memória primária. Enquanto forem realizadas as n inserções, o índice deve ser atualizado somente em memória primária. Ao finalizar as n inserções, realize a operação de reescrita do índice da memória primária para o disco. Nessa operação de reescrita, use a opção wb para reescrever no arquivo.

Entrada do programa para a funcionalidade [13]:

```
13 arquivo.bin arquivoIndiceNomeServidor.bin n
valorIdServidor1 valorSalarioServidor1 valorTelefoneServidor1
valorNomeServidor1 valorCargoServidor1
valorIdServidor2 valorSalarioServidor2 valorTelefoneServidor2
valorNomeServidor2 valorCargoServidor2
...
valorIdServidorn valorSalarioServidorn valorTelefoneServidorn
valorNomeServidorn valorCargoServidorn
```

onde:

- arquivo.bin é um arquivo binário de entrada que segue as mesmas especificações definidas no primeiro trabalho prático. As inserções a serem realizadas nessa funcionalidade devem ser feitas nesse arquivo.
- arquivoIndiceNomeServidor.bin é o arquivo binário referente ao índice secundário fortemente ligado que segue as especificações definidas nesse trabalho prático.
- n é o número de inserções a serem realizadas, as quais seguem as mesmas especificações definidas na funcionalidade [5] do segundo trabalho prático. Ou seja, para cada inserção, deve ser informado os valores a serem inseridos no arquivo de dados, para os campos especificados na mesma ordem que a definida no primeiro trabalho prático, a saber: idServidor, salarioServidor, telefoneServidor, nomeServidor, cargoServidor. Valores nulos devem ser identificados, na entrada da funcionalidade, por NULO. Cada uma das n inserções deve ser especificada em uma linha diferente. Deve ser deixado um espaço em branco entre os valores dos campos. Os valores dos campos do tipo *string* devem ser especificados entre aspas duplas (").

Saída caso o programa seja executado com sucesso:

Listar o arquivo de índice secundário arquivoIndiceNomeServidor.bin.

Mensagem de saída caso algum erro seja encontrado:

Falha no processamento do arquivo.

Exemplo de execução:

```
./programaTrab4
13 arquivo.bin arquivoIndiceNomeServidor.bin 2
1234 6099.57 NULO NULO "ASSISTENTE EM CIENCIA E TECNOLOGIA"
2132 NULO NULO "JOAO JOSE PEDRO DA SILVA TEIXEIRA" "ANALISTA DE
SEGURO"
```

usar a função binarioNaTela1 ou binarioNaTela2 antes de terminar a execução da funcionalidade, para mostrar a saída do arquivo arquivoIndiceNomeServidor.bin, o qual foi atualizado frente às inserções realizadas no arquivo de dados.

[14] Permita a realização de estatísticas considerando a recuperação dos dados de todos os registros que satisfaçam um critério de busca determinado pelo usuário sobre o campo *nomeServidor*. Essa funcionalidade deve ser implementada da seguinte forma. Dado o nome de um servidor como entrada, primeiramente execute a funcionalidade [3] definida no primeiro trabalho prático. Em seguida, execute a funcionalidade [11] definida no presente trabalho prático. Ao final, liste a diferença no número de páginas de disco acessadas, contabilizadas da seguinte forma: número de páginas de disco acessadas na funcionalidade [3] *menos* o número de páginas de disco para acessar o arquivo de dados na funcionalidade [11]. Note que o número de páginas de disco para carregar o arquivo de índice da funcionalidade [11] foi desconsiderado no cálculo porque, em uma aplicação real, o carregamento do índice é feito apenas uma única vez no início da execução do programa, sendo o índice reescrito em disco somente no final da execução do programa. Antes de exibir o resultado da funcionalidade [3], exiba a frase “*** Realizando a busca sem o auxílio de índice”. Depois de exibir o resultado da funcionalidade [3]. Adicionalmente, antes de exibir no resultado da funcionalidade [11], exiba a frase “*** Realizando a busca com o auxílio de um índice secundário fortemente ligado”. Na sequência, exiba a diferença no número de páginas de disco acessadas.

Sintaxe do comando para a funcionalidade [14]:

```
14 arquivo.bin arquivoIndiceNomeServidor.bin nomeServidor valor
```

Saída caso o programa seja executado com sucesso:

Liste na sequência descrita a seguir:

- 1- A frase: *** Realizando a busca sem o auxílio de índice
- 2- O resultado da execução da funcionalidade [3], conforme especificado no segundo trabalho prático
- 3- A frase: *** Realizando a busca com o auxílio de um índice secundário fortemente ligado
- 4- O resultado da execução da funcionalidade [11], conforme especificado neste trabalho prático
- 5- A frase: A diferença no número de páginas de disco acessadas:
- 6- O resultado da diferença, dado pelo número de páginas de disco acessadas na funcionalidade [3] *menos* o número de páginas de disco para acessar o arquivo de dados na funcionalidade [11].

Mensagem de saída caso não seja encontrado o registro que contém o valor do campo ou o campo pertence a um registro que esteja removido:

Registro inexistente.

Mensagem de saída caso algum erro seja encontrado:

Falha no processamento do arquivo.

Exemplo de execução (considere que os metadados lidos do arquivo .csv sejam: numero de identificacao do servidor, salario do servidor, telefone celular do servidor, nome do servidor, cargo do servidor):

```
./programaTrab4
14 arquivo.bin arquivoIndiceServidor.bin nomeServidor CRISTIANO ANDRE
DA SILVA
*** Realizando a busca sem o auxílio de índice
numero de identificacao do servidor: 6202963
salario do servidor: 4097.91
telefone celular do servidor: (60)99435-1564
nome do servidor: CRISTIANO ANDRE DA SILVA
cargo do servidor: AGENTE ADMINISTRATIVO
--- pular uma linha em branco ----
Número de páginas de disco acessadas: 2
*** Realizando a busca com o auxílio de um índice secundário
fortemente ligado
numero de identificacao do servidor: 6202963
salario do servidor: 4097.91
telefone celular do servidor: (60)99435-1564
nome do servidor: CRISTIANO ANDRE DA SILVA
cargo do servidor: AGENTE ADMINISTRATIVO
--- pular uma linha em branco ----
Número de páginas de disco para carregar o arquivo de índice: x
Número de páginas de disco para acessar o arquivo de dados: 1
--- pular uma linha em branco ----
Diferença no número de páginas de disco acessadas: 1
```

Restrições

As seguintes restrições têm que ser garantidas no desenvolvimento do trabalho.

[1] O arquivo de dados deve ser gravado em disco no **modo binário**. O modo texto não pode ser usado.

[2] Os dados do registro descrevem os nomes dos campos, os quais não podem ser alterados. Ademais, todos os campos devem estar presentes na implementação, e nenhum campo adicional pode ser incluído. O tamanho e a ordem de cada campo deve obrigatoriamente seguir a especificação.

[3] Deve haver a manipulação de valores nulos, conforme as instruções definidas.

[4] Não é necessário realizar o tratamento de truncamento de dados.

[5] Devem ser exibidos avisos ou mensagens de erro de acordo com a especificação de cada funcionalidade.

[6] Os dados devem ser obrigatoriamente escritos e lidos campo a campo. Ou seja, não é possível escrever e ler os dados registro a registro.

[7] O(s) aluno(s) que desenvolveu(desenvolveram) o trabalho prático deve(m) constar como comentário no início do código (i.e. NUSP e nome do aluno). Para trabalhos desenvolvidos por mais do que um aluno, não será atribuída nota ao aluno cujos dados não constarem no código fonte.

[8] Todo código fonte deve ser documentado. A **documentação interna** inclui, dentre outros, a documentação de procedimentos, de funções, de variáveis, de partes do código fonte que realizam tarefas específicas. O código fonte deve ser documentado tanto em nível de rotinas quanto em nível de variáveis e blocos funcionais.

[9] A implementação deve ser realizada usando a linguagem de programação C. As funções das bibliotecas <stdio.h> devem ser utilizadas para operações relacionadas à escrita e leitura dos arquivos. A implementação não pode ser feita em qualquer outra linguagem de programação. O programa executará no [run.codes].

Fundamentação Teórica

Conceitos e características dos diversos métodos para representar os conceitos de campo e de registro em um arquivo de dados podem ser encontrados nos *slides* de sala de aula e também nas páginas 96 a 107 do livro *File Structures (second edition)*, de Michael J. Folk e Bill Zoellick.

Material para Entregar

Arquivo compactado. Deve ser preparado um arquivo .zip contendo:

- Código fonte do programa devidamente documentado.
- Makefile para a compilação do programa.

Instruções para fazer o arquivo makefile. No [run.codes] tem uma orientação para que, no makefile, a diretiva “all” contenha apenas o comando para compilar seu programa e, na diretiva “run”, apenas o comando para executá-lo. Assim, a forma mais simples de se fazer o arquivo makefile é:

```
all:  
gcc -o programaTrab2 *.c  
run:  
./programaTrab2
```

Lembrando que *.c já engloba todos os arquivos .c presentes no seu zip."

Instruções de entrega. A entrega deve ser feita via [run.codes]:

- página: <https://run.codes/Users/login>
- código de matrícula para a **Turma A: SRAZ**
- código de matrícula para a **Turma B: 29QF**

Critério de Correção

Critério de avaliação do trabalho. Na correção do trabalho, serão ponderados os seguintes aspectos.

- Corretude da execução do programa.
- Atendimento às especificações do registro de cabeçalho e dos registros de dados.
- Atendimento às especificações da sintaxe dos comandos de cada funcionalidade e do formato de saída da execução de cada funcionalidade.
- Qualidade da documentação entregue.

Restrições adicionais sobre o critério de correção.

- A não execução de um programa devido a erros de compilação implica que a nota final da parte do trabalho será igual a zero (0).
- O não atendimento às especificações do registro de cabeçalho e dos registros de dados implica que haverá uma diminuição expressiva na nota do trabalho.
- O não atendimento às especificações de sintaxe dos comandos de cada funcionalidade e do formato de saída da execução de cada funcionalidade implica que haverá uma diminuição expressiva na nota do trabalho.
- A ausência da documentação implica que haverá uma diminuição expressiva na nota do trabalho.
- A inserção de palavras ofensivas nos arquivos e em qualquer outro material entregue implica que a nota final da parte do trabalho será igual a zero (0).
- Em caso de plágio, as notas dos trabalhos envolvidos serão zero (0).

Data de Entrega do Trabalho

Na data especificada na página da disciplina.

Bom Trabalho !
