



TRABALHO I

Pesquisa e Ordenação de Dados

João Luis Almeida Santos

SUMÁRIO

- 1. Introdução**
Instruções e objetivos do trabalho
- 2. Desafios e Evolução do Código**
Metodologia, descrição da lógica, sumário de decisões
- 3. Implementação e Solução Principal**
Metodologia, descrição da lógica, sumário de decisões
- 4. Análise de Complexidade**
Análise do Algoritmo e definição matemática do tempo de execução
- 5. Conclusão**
Conclusão final a partir dos resultados e processo de execução do trabalho
- 6. Referências e Código-Fonte**
Referências utilizadas no relatório final e código-fonte. (Também em anexo.)

O projeto implementado também pode ser acessado pelo repositório do GITHUB, [através deste link](#).

Este relatório apresenta o desenvolvimento e a análise de uma solução voltada para a ordenação externa de números inteiros positivos, conforme proposto no Trabalho T1 da disciplina. O objetivo central é aplicar técnicas de manipulação de arquivos e otimização de uso de memória principal, simulando um cenário em que os dados excedem a capacidade de armazenamento interno e precisam ser processados em blocos. Apesar dos comentários estarem em português, a prática para nomear as funções e variáveis foram baseadas no inglês. O documento explora a lógica aplicada e discute as decisões de design algorítmico, identifica desafios enfrentados durante o desenvolvimento e realiza uma análise da complexidade do código. O código clona o arquivo original para um arquivo que podemos ler e alterar sem alterar o arquivo principal, INPUT. Após, há a ordenação dos dados, e salvamento dos resultados no arquivo SAIDA.TXT. Este relatório estuda as decisões críticas tomadas durante a realização do trabalho, as estratégias de manipulação de arquivos temporários, e análise da complexidade.

Aqui estão as bibliotecas utilizadas e as variáveis globais usadas no código:

```
1
2 // Configurado para o meu computador, alterar se for usar no Windows
3 #ifdef _WIN32
4 #define INPUT "C:\\Users\\neoka\\CLionProjects\\Trabalho---POD\\dados.txt"
5 #endif
6 #ifdef linux
7 #define INPUT "dados.txt"
8 #endif
9 #include <ctype.h>
10 #include <stdio.h>
11 #include <stdbool.h>
12 #include <stdlib.h>
13 #include <string.h>
14
15 #define OUTPUT "saida.txt"
```

I. DESAFIOS E EVOLUÇÃO DO CÓDIGO

Na primeira versão do código, desenvolvi uma solução baseada no carregamento completo dos dados em memória principal usando um array, seguido da aplicação do *Insertion Sort*. Essa abordagem, disponível no GitHub¹, partia da premissa equivocada de que a ordenação poderia ser feita somente lendo os arquivos e colocando os números em um array. Na minha segunda tentativa, fiz um programa que realizava o processo de merge incorretamente. Para a versão final foi necessário mais pesquisa para entender como fazer apropriadamente o programa funcionar.

```
int main () {
    // Estou colocando 15 números
    int data_amount;
    printf("Insira a quantidade de elementos:\n");
    if (scanf("%d", &data_amount) != 1 || data_amount <= 0){
        printf("[ERRO] Entrada inválida\n");
        return 0;
    };
    int* list = malloc(data_amount * sizeof(int));
    if (list == NULL)
    {
        printf("[ERRO] Falha ao alocar memória.\n");
        return 1;
    }
    makeNumList(data_amount, list);

    printf("LISTA ORIGINAL:\n");
    printIntList(list, data_amount);
    insertionSort(list, data_amount);
    // Insere a lista em um arquivo, numeros separados via ";"
    intListOrderToFile(list, data_amount);
    free(list);

    return 0;
}
```

¹<https://github.com/lysz/Trabalho—POD/blob/main/antigo.c>

```
}
```

A função **externalSort(int chunk_size)** segue estes passos: contamos quantos números existem no arquivo, criamos arquivos temporários para cada bloco de tamanho **chunk_size**, usamos **qsort** para ordenar cada bloco, realizamos o merge dos blocos ordenados e, por fim, removemos os arquivos temporários.

Fizemos melhorias como contagem exata de elementos, criação dinâmica de arquivos temporários, ordenação interna eficiente com **qsort** e merge em tempo $O(n \log m)$, além de cuidar automaticamente da memória.

Na contagem dos números, a função `count_numbers_in_file` lê todo o arquivo, ignora caracteres não numéricos e conta quantos inteiros existem. Usamos **fscanf** para isso.

Para ordenar cada bloco, carregamos até **chunk_size** números na memória e aplicamos:

```
qsort(bloco, tamanho, sizeof(int), compare);
```

No merge, repetimos a busca pelo menor valor entre todos os blocos e escrevemos esse valor no arquivo de saída, até processar todos os números. Inicialmente fiz a escolha de utilizar o bubble sort, mas descobri depois que o `qsort` (implementação do Quicksort) era bem mais prático e agilizava a escrita do código, e é uma função que já existe por padrão no C. Bastou-se criar a função de `compare` para utilizar junto.

O programa começa contando quantos números existem no arquivo. Isso é crucial para determinar quantos blocos serão criados. A função **count_numbers_in_file** ignora caracteres inválidos e conta apenas dígitos, garantindo que espaços ou símbolos acidentais não quebrem o programa. Para arquivos muito grandes, essa contagem é feita sem carregar tudo na memória. O código lê o arquivo pedaço por pedaço, como mostra este trecho de código:

```

// Function for reading data
void makeNumList(const int size, int* num_list) {
    FILE *f = fopen(INPUT, "r");

    // A função vai adicionando os caracteres dos números a uma string
    // até encontrar um caracter vazio, ponto e vírgula, ou o fim do arquivo

    if (f == NULL) {
        printf("[ERRO] Não foi possível abrir o arquivo, saindo.\n");
        return;
    }

    int i = 0;
    int j = 0;
    char number[16];
    char c;

    while ((c = fgetc(f)) != EOF && j < size) {
        if (c == ' ' || c == '\n') {
            if (i > 0) {
                number[i] = '\0';
                num_list[j++] = strtol(number, NULL, 10);
                i = 0;
            }
            continue;
        }

        if (c == ';') {
            if (i > 0) {
                number[i] = '\0';
                num_list[j++] = strtol(number, NULL, 10);
                i = 0;
            }
        } else if (i < 15) {
            number[i++] = c;
        }

        if (j >= size) {
            break;
        }
    }
}

```

```

// Trata o último número, se necessário
if (i > 0 && j < size) {
    number[i] = '\0';
    num_list[j++] = strtol(number, NULL, 10);
}

fclose(f);
}

```

Um desafio foi lidar com o último bloco, que pode ter menos números que os outros. O código ajusta automaticamente o tamanho desse bloco final usando cálculo modular:

```

int tamanho = (i == total_blocos - 1) ? total % chunk_size : chunk_size;

```

A etapa final usa uma estratégia de "menor elemento primeiro". Todos os blocos ordenados são abertos simultaneamente. O programa sempre pega o menor número disponível entre os primeiros elementos de cada bloco e escreve no arquivo final.

```

1 void externalSort(int chunk_size) {
2     printf("\n=== INICIANDO ORDENACAO EXTERNA ===\n");
3     printf("Tamanho dos blocos: %d elementos\n", chunk_size);
4
5
6     // 1. Verificar quantidade total de números
7     printf("\n[PASSO 1] Contando numeros no arquivo...\n");
8     int total = count_numbers_in_file(INPUT);
9     if(total == -1) {
10        printf("Falha contando quantidade de números no arquivo.\n");
11        return;
12    }
13    if(total <= 0) {
14        printf("Erro ao ler arquivo de entrada\n");
15        return;
16    }
17    printf("Total de numeros encontrados: %d\n", total);
18
19    // 2. Calcular quantos arquivos temporários serão criados

```

```

20     int num_arquivos = (total + chunk_size - 1) / chunk_size;
21     printf("\n[PASSO 2] Criando %d arquivos temporarios\n", num_arquivos);
22     char nome_arquivo[256];
23
24     // 3. Ler, ordenar e salvar em arquivos temporários
25     printf("\n[PASSO 3] Processando blocos:\n");
26     FILE *entrada = fopen(INPUT, "r");
27     for(int i = 0; i < num_arquivos; i++) {
28         // Calcular tamanho do bloco
29         int tamanho;
30         if (i == num_arquivos - 1) {
31             tamanho = total % chunk_size;
32         } else {
33             tamanho = chunk_size;
34         }
35         printf("\n-> Bloco %d/%d (%d elementos)\n", i+1, num_arquivos, tamanho);
36
37         // Ler números
38         printf("    Lendo do arquivo original...\n");
39         int *bloco = malloc(tamanho * sizeof(int));
40         for(int j = 0; j < tamanho; j++) {
41             fscanf(entrada, "%d", &bloco[j]);
42             while(fgetc(entrada) == ';')
43                 continue;
44         }
45
46         // Ordenar e salvar
47         printf("    Ordenando...\n");
48         qsort(bloco, tamanho, sizeof(int), compare);
49
50         sprintf(nome_arquivo, PROCESS_N, i);
51         printf("    Salvando no arquivo temporario: %s\n", nome_arquivo);
52         FILE *temp = fopen(nome_arquivo, "w");
53         for(int j = 0; j < tamanho; j++) fprintf(temp, "%d;", bloco[j]);
54         fclose(temp);
55         free(bloco);
56     }
57     fclose(entrada);
58
59     // 4. Fazer merge dos arquivos temporários
60     printf("\n[PASSO 4] Iniciando merge dos arquivos:\n");
61     // Ponteiro duplo é um array de arquivos FILE
62     // Isso facilita aqui pra baixo

```



```

61 FILE **arquivos = malloc(num_arquivos * sizeof(FILE*));
62 int *valores = malloc(num_arquivos * sizeof(int));
63
64 // Abrir arquivos temporários
65 printf("Abrindo arquivos para merge...\n");
66 for(int i = 0; i < num_arquivos; i++) {
67     sprintf(nome_arquivo, PROCESS_N, i);
68     arquivos[i] = fopen(nome_arquivo, "r");
69     fscanf(arquivos[i], "%d;", &valores[i]);
70 }
71
72 // Escrever resultado ordenado
73 int elementos_processados = 0;
74 FILE *saida = fopen(OUTPUT, "w");
75 printf("\n[PASSO 5] Gerando arquivo final:\n");
76 while(1) {
77     // Encontrar menor valor
78     int menor_valor = INT_MAX;
79     int indice_menor = -1;
80
81     for(int i = 0; i < num_arquivos; i++) {
82         if(valores[i] < menor_valor) {
83             menor_valor = valores[i];
84             indice_menor = i;
85         }
86     }
87
88     if(indice_menor == -1) break;
89
90     // Escrever e atualizar
91     fprintf(saida, "%d;", menor_valor);
92     elementos_processados++;
93
94     if(elementos_processados % 1000 == 0) {
95         printf("-> Progresso: %d elementos ordenados\n", elementos_processados);
96     }
97
98     if(fscanf(arquivos[indice_menor], "%d;", &valores[indice_menor]) != 1) {
99         printf("Arquivo %d esvaziado\n", indice_menor);
100         valores[indice_menor] = INT_MAX;
101     }

```

```

102     }
103
104     // 5. Limpeza final
105     printf("\n[PASSO 6] Finalizando processo:\n");
106     for(int i = 0; i < num_arquivos; i++) {
107         fclose(arquivos[i]);
108         sprintf(nome_arquivo, PROCESS_N, i);
109         remove(nome_arquivo);
110         printf("Removido arquivo temporario: %s\n", nome_arquivo);
111     }
112     free(arquivos);
113     free(valores);
114     fclose(saida);
115
116     printf("\n=== PROCESSO CONCLUIDO ===\n");
117     printf("Arquivo final gerado: %s\n", OUTPUT);
118     printf("Total de elementos ordenados: %d\n\n", elementos_processados);
119 }
120 }

```

II. ANÁLISE DE COMPLEXIDADE

- **n**: Número total de elementos
- **k**: Tamanho do bloco (chunk_size)
- **m**: Número de blocos ($m = n/k$)

```

1  int compare(const void* a, const void* b) {
2      return (*(int*)a - *(int*)b); // O(1)
3  }
4
5  void cloneFile(const char* origin, const char* destination) {
6      O(M) onde M = tamanho do arquivo
7  }
8
9  int count_numbers_in_file(const char *filename) {
10     O(M) onde M = tamanho do arquivo
11 }
12

```

```

13 void externalSort(int chunk_size) {
14     // O(M + N log C + N²/C)
15 }
16
17 int main() {
18     // O(1) + O(M + N log C + N²/C)
19     // = O(M + N log C + N²/C)
20 }
21 }
22

```

$$\mathcal{O}\left(M + N \log C + \frac{N^2}{C}\right)$$

III. CONCLUSÃO

A partir do presente trabalho, é possível concluir que o algoritmo desenvolvido demonstrou com sucesso os princípios fundamentais da ordenação externa utilizando arquivos de texto. O código, quando executado, se apresenta como demonstrado na figura 1. A fase de sort performa mais rapidamente quando executado com mais chunks. E a de merge, ao contrário.

```

Quantidade de elementos por vez:
1
Saída final disponível em [ SAIDA.TXT ]

```

Figura 1: Console na execução do código.