

Universidade Federal de São Paulo – UNIFESP  
Campus: Parque Tecnológico – São José dos Campos



Instituto de Ciência e Tecnologia – ICT

Bacharelado em Ciência e Tecnologia – BCT

**Algoritmos e Estruturas de Dados II**

## **Ordenação Externa**

Prof. Reginaldo Massanobu Kuroshu

**João Pedro da Silva Zampoli**

**RA: 168880**

**Luiza de Souza Ferreira**

**RA: 170453**

**Viviane Flor Park**

**RA: 169259**

São José dos Campos, 2024

## SUMÁRIO

<b>1 INTRODUÇÃO.....</b>	<b>4</b>
<b>2 OBJETIVOS.....</b>	<b>5</b>
2.1 Objetivo geral.....	5
2.2 Objetivos específicos.....	5
<b>3 MÉTODOS.....</b>	<b>6</b>
<b>4 RESULTADOS E DISCUSSÕES.....</b>	<b>9</b>
4.1 Tempos de Execução.....	9
4.2 Leitura e Escrita de Registros.....	11
<b>5 CONCLUSÃO.....</b>	<b>13</b>

## **RESUMO**

Neste relatório da UC “Algoritmos e Estruturas de Dados II”, o objeto de estudos é um código de ordenação externa em Python agindo sobre 3 entradas de tamanhos diferentes, tendo como objetivo comparar o desempenho do algoritmo aplicado sobre elas.

## 1 INTRODUÇÃO

Em um mundo onde o volume de dados gerados e processados está em constante crescimento, a eficiência na manipulação e ordenação desses dados tornou-se crucial. Quando o volume de dados excede a capacidade da memória principal (RAM), técnicas de ordenação especializadas são necessárias para garantir que a manipulação e análise dos dados sejam realizadas de maneira eficaz. Neste contexto, a ordenação externa surge como uma solução essencial.

A ordenação externa refere-se a métodos e algoritmos projetados para ordenar grandes conjuntos de dados que não cabem inteiramente na memória principal. Esses algoritmos são especialmente relevantes em sistemas de gerenciamento de bancos de dados e em aplicações que lidam com grandes volumes de dados, como processamento de dados em grandes servidores e sistemas de arquivos distribuídos.

O External Merge Sort é um dos algoritmos mais eficientes para ordenação externa, desenvolvido para superar as limitações das abordagens tradicionais de ordenação em memória. O algoritmo divide o conjunto de dados em blocos menores que podem ser processados na memória, ordena esses blocos individualmente e, em seguida, realiza uma fusão dos blocos ordenados para obter o conjunto de dados final ordenado. Esse processo é altamente eficiente, principalmente devido à sua capacidade de minimizar o número de operações de entrada e saída, que são geralmente o obstáculo em sistemas de ordenação externa.

Este relatório explora a aplicação dos princípios fundamentais da ordenação, em especial, a do External Merge Sort. Abordaremos a estrutura e o processo do algoritmo, e faremos uma análise do tempo de execução da leitura e escrita de registros para entradas de tamanhos diferentes. Ao final, a análise proporcionará uma compreensão clara das aplicações e benefícios desta técnica na ordenação de grandes volumes de dados.

## **2 OBJETIVOS**

### **2.1 Objetivo geral**

Avaliar o desempenho de um algoritmo de ordenação externa para 3 entradas distintas.

### **2.2 Objetivos específicos**

- Calcular a quantidade de vezes que cada registro é lido e escrito no disco para cada arquivo de entrada;
- Comparar o tempo de execução do código usado para cada entrada;
- Descrever detalhes dos experimentos realizados e discutir sobre resultados observados.

### 3 MÉTODOS

Primeiramente, foi escolhida a linguagem de programação de tipagem dinâmica Python para o código de ordenação externa. A escolha foi direcionada pela impressão que os integrantes do grupo têm de que o tratamento de arquivos nessa linguagem é mais simples.

O código segue a lógica de um External Merge Sort. Para isso, foram seguidos os passos 1 a 5, descritos abaixo:

- 1) Lemos 5000 registros do arquivo de entrada e ordenamos esse bloco utilizando HeapSort;
- 2) Escrevemos registros ordenados em um arquivo intermediário;
- 3) Repetimos 1) e 2) até que todos os registros do arquivo de entrada tenham sido lidos e gravados;
- 4) Intercalamos os arquivos intermediários dois a dois com HeapSort, guardando os registros ordenados em um novo arquivo intermediário;
- 5) Repetimos 4) enquanto o passo gerar mais de um arquivo intermediário.

As Figuras 3.1 e 3.2 mostram o código implementado.

```

1  import os
2
3  Leitura = 0
4  Escrita = 0
5  QtdNumeros = 5000
6  ParteDeExecucao = 1
7  TamanhoTotal = len(open('input2.txt').readlines())
8
9  def ContarLeitura():
10     global Leitura
11     Leitura += 1
12
13  def ContarEscrita():
14     global Escrita
15     Escrita += 1
16
17  def Heapify(Vet, Tamanho, Indice, min_heap=False):
18     Extremo = Indice
19     Esquerda = 2 * Indice + 1
20     Direita = 2 * Indice + 2
21
22     if Esquerda < Tamanho and ((not min_heap and Vet[Indice] < Vet[Esquerda]) or (min_heap and Vet[Indice] > Vet[Esquerda]]):
23         Extremo = Esquerda
24
25     if Direita < Tamanho and ((not min_heap and Vet[Extremo] < Vet[Direita]) or (min_heap and Vet[Extremo] > Vet[Direita]]):
26         Extremo = Direita
27
28     if Extremo != Indice:
29         Vet[Indice], Vet[Extremo] = Vet[Extremo], Vet[Indice]
30         Heapify(Vet, Tamanho, Extremo, min_heap)
31
32  def HeapSort(Vet):
33     Tamanho = len(Vet)
34
35     for i in range(Tamanho // 2 - 1, -1, -1):
36         Heapify(Vet, Tamanho, i, min_heap=False)
37
38     for i in range(Tamanho - 1, 0, -1):
39         Vet[i], Vet[0] = Vet[0], Vet[i]
40         Heapify(Vet, i, 0, min_heap=False)
41
42  def CriarArquivosIntermediarios():
43     Pedacos = []
44     with open('input3.txt', 'r') as f:
45         Pedaco = []
46         for Linha in f:
47             ContarLeitura()
48             Pedaco.append(int(Linha))
49             if len(Pedaco) == QtdNumeros:
50                 HeapSort(Pedaco)
51                 Pedacos.append(Pedaco[:])
52                 Pedaco = []
53         if Pedaco:
54             HeapSort(Pedaco)
55             Pedacos.append(Pedaco[:])
56     return Pedacos
57
58  def EscreverArquivosIntermediarios(Pedacos):
59     if not os.path.exists("Pedacos"):
60         os.makedirs("Pedacos")
61     for i in range(len(Pedacos)):
62         with open(f'Pedacos/Parte0{ParteDeExecucao}Pedaco_{i+1}.txt', 'w') as f:
63             for Num in Pedacos[i]:
64                 ContarEscrita()
65                 f.write(f'{Num}\n')
66

```

Figura 3.1 - Primeira parte do código implementado para a ordenação externa.

```

67 def CriarPartes():
68     global ParteDeExecucao
69     global QtdNumeros
70     global TamanhoTotal
71
72     Pedacos = CriarArquivosIntermediarios()
73     EscreverArquivosIntermediarios(Pedacos)
74     while len(Pedacos) > 1:
75         ParteDeExecucao += 1
76         QtdNumeros *= 2
77         Pedacos = CriarArquivosIntermediarios()
78         EscreverArquivosIntermediarios(Pedacos)
79     return Pedacos
80
81 def CriarOutput():
82     f = open('output3.txt', 'w')
83     g = open(f'Pedacos/Parte0{ParteDeExecucao}Pedaco_01.txt', 'r')
84
85     for Linha in g:
86         f.write(Linha)
87     f.close()
88     g.close()
89
90 Pedacos = CriarPartes()
91 CriarOutput()
92
93 print("Tamanho Total: ", TamanhoTotal)
94
95 print(f'Leituras: {Leitura}')
96 print(f'Escritas: {Escrita}')
97
98 LeituraRegistro = Leitura / TamanhoTotal
99 EscritaRegistro = Escrita / TamanhoTotal
100
101 print(f'Quantidade de Leituras dos Registros: {ParteDeExecucao}')
102 print(f'Quantidade de Escritas dos Registros: {ParteDeExecucao}')
103
104 def ApagarPasta():
105     try:
106         if not os.path.exists("Pedacos"):
107             raise ValueError("Pasta não existe")
108
109         for item in os.listdir("Pedacos"):
110             CaminhoItem = os.path.join("Pedacos", item)
111
112             if os.path.isfile(CaminhoItem):
113                 os.remove(CaminhoItem)
114                 # print(f"Arquivo {item} apagado")
115             else:
116                 raise ValueError("Item não é um arquivo")
117
118     except Exception as e:
119         print(f"Erro ao tentar remover: {e}")
120
121     try:
122         os.rmdir("Pedacos")
123         #print("Pasta Pedacos removida")
124
125     except Exception as e:
126         print(f"Erro ao tentar remover: {e}")
127
128 ApagarPasta()

```

Figura 3.2 - Segunda parte do código implementado para a ordenação externa.

O código foi então executado para três entradas distintas, fornecidas junto ao enunciado do exercício. Para cada uma delas, o código foi executado dez vezes em uma terminal macOS, tendo os tempos de execução anotados para análise posterior. O código retorna, também, a quantidade de vezes que os registros foram lidos e escritos no disco.



## 4 RESULTADOS E DISCUSSÕES

### 4.1 Tempos de Execução

A Tabela 4.1.1 mostra os tempos de execução real, de usuário e de sistema registrados.

Tabela 4.1.1 - Tempos de execução do código para cada entrada e para cada execução, em segundos.

	input 1			input 2			input 3		
	real	usuário	sistema	real	usuário	sistema	real	usuário	sistema
1	0.391	0.304	0.028	9.189	7.784	0.338	150.445	129.856	4.765
2	0.355	0.296	0.023	9.027	7.711	0.302	150.849	130.296	4.690
3	0.351	0.292	0.024	9.665	7.867	0.315	150.614	130.122	4.771
4	0.350	0.291	0.024	8.914	7.690	0.302	152.292	130.300	4.789
5	0.361	0.294	0.025	8.934	7.708	0.304	149.901	129.641	4.667
6	0.358	0.298	0.024	9.111	7.771	0.306	151.324	129.827	4.671
7	0.358	0.294	0.025	8.865	7.699	0.294	151.033	130.613	4.815
8	0.980	0.328	0.029	8.986	7.762	0.301	153.103	130.732	4.828
9	0.394	0.304	0.028	8.936	7.698	0.297	157.701	134.210	4.829
10	0.481	0.307	0.030	8.747	7.617	0.287	150.927	129.407	4.652

A partir dos dados de tempo de execução real foram construídos os gráficos das Figuras 4.1.1, 4.1.2 e 4.1.3 por meio do RStudio.



Figura 4.1.1 - Gráfico de linha para o tempo real de execução do input 1.



Figura 4.1.2 - Gráfico de linha para o tempo real de execução do input 2.

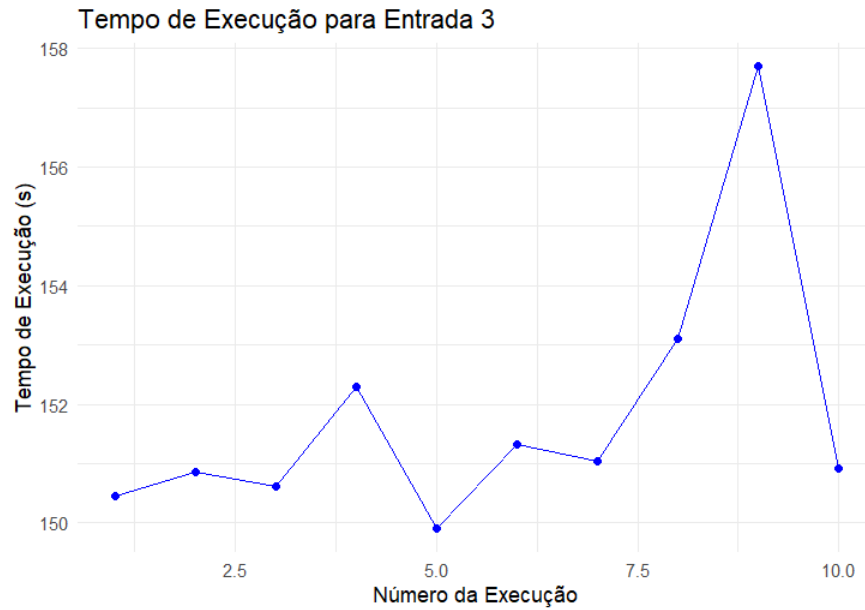


Figura 4.1.3 - Gráfico de linha para o tempo real de execução do input 3.

Após isso, foi feito o gráfico da Figura 4.1.4 comparando o tempo de execução real das 3 entradas por meio do RStudio.

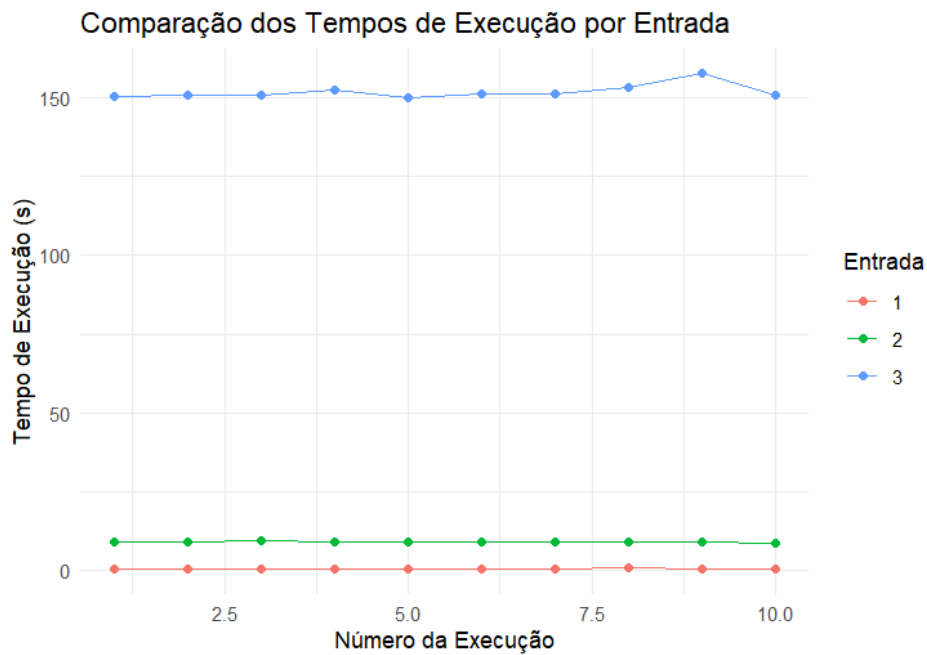


Figura 4.1.4 - Gráfico de linhas comparando tempo real de execução entre as 3 entradas.

## 4.2 Leitura e Escrita de Registros

As Figuras 4.2.1, 4.2.2 e 4.2.3 mostram as saídas do código para uma das execuções das entradas 1, 2 e 3, respectivamente. Essas saídas fornecem a

quantidade de vezes que os registros foram lidos e escritos no disco, referentes tanto a um único registro, quanto a todos os registros do arquivo de entrada.

```
Tamanho Total: 10000
Leituras: 20000
Escritas: 20000
Quantidade de Leituras dos Registros: 2
Quantidade de Escritas dos Registros: 2

real    0m0.391s
user    0m0.304s
sys     0m0.028s
```

Figura 4.2.1 - Saída da 1ª execução do input 1.

```
Tamanho Total: 100000
Leituras: 600000
Escritas: 600000
Quantidade de Leituras dos Registros: 6
Quantidade de Escritas dos Registros: 6

real    0m9.189s
user    0m7.784s
sys     0m0.338s
```

Figura 4.2.2 - Saída da 1ª execução do input 2.

```
Tamanho Total: 1000000
Leituras: 9000000
Escritas: 9000000
Quantidade de Leituras dos Registros: 9
Quantidade de Escritas dos Registros: 9

real    2m30.927s
user    2m9.407s
sys     0m4.652s
```

Figura 4.2.3 - Saída da 10ª execução do input 3.

## **5 CONCLUSÃO**

De acordo com os resultados obtidos, os tempos de execução foram pequenos de acordo com cada entrada fornecida, além disso, foi observado que o número de acesso às unidades de memória externa foi baixo comparado com o tamanho de cada entrada fornecida, portanto foi mostrado a importância da ordenação externa para arquivos de tamanho maior que a memória interna disponível.