

Algoritmos e Estruturas de Dados II

Aula 28 – Processamento Cossequencial e Ordenação Externa

Prof. Luciano A. Digiampietri
digiampietri@usp.br
@digiampietri

2024

Processamento Cossequencial

- Processamento coordenado e sequencial de uma ou mais listas (no intuito de formar uma lista única)
- Exemplo: obtenção da intersecção ou união de listas ordenadas

Input	Array 1	10	15	22	80	
	Array 2	5	10	11	22	70

Output:

Union : 5 10 11 15 22 70 80 90
Intersection : 10 22

Processamento Cossequencial

- Etapas:
 - Inicialização (abrir arquivos, inicializar variáveis)
 - Sincronização (como avançar em cada lista)
 - Condições de fim de lista (o que fazer)
 - Reconhecimento de erros (houve duplicações? Elementos fora de ordem?)

Processamento Cossequencial

- Base para a Ordenação Externa
- Maximizar a manipulação em memória para minimizar o número de acessos ao disco

Ordenação Externa

- Objetivo: ordenar um arquivo muito grande, que não cabe inteiro na memória
- O que fazer: ordenar pedaços desse arquivo (em memória), e depois combinar os pedaços

Ordenação Externa

Projeto de Algoritmos – Cap.4 Ordenação – Seção 4.2

Ordenação Externa

- A ordenação externa consiste em ordenar arquivos de tamanho maior que a memória interna disponível.
- Os métodos de ordenação externa são muito diferentes dos de ordenação interna.
- Na ordenação externa os algoritmos devem diminuir o número de acesso as unidades de memória externa.
- Nas memórias externas, os dados ficam em um arquivo seqüencial.
- Apenas um registro pode ser acessado em um dado momento. Essa é uma restrição forte se comparada com as possibilidades de acesso em um vetor.
- Logo, os métodos de ordenação interna são inadequados para ordenação externa.
- Técnicas de ordenação diferentes devem ser utilizadas.

(ZIVIANI, 2010)

Ordenação Externa

Ordenação Externa

Fatores que determinam as diferenças das técnicas de ordenação externa:

1. Custo para acessar um item é algumas ordens de grandeza maior.
2. O custo principal na ordenação externa é relacionado a transferência de dados entre a memória interna e externa.
3. Existem restrições severas de acesso aos dados.
4. O desenvolvimento de métodos de ordenação externa é muito dependente do estado atual da tecnologia.
5. A variedade de tipos de unidades de memória externa torna os métodos dependentes de vários parâmetros.
6. Assim, apenas métodos gerais serão apresentados.

(ZIVIANI, 2010)

Ordenação Externa

Ordenação Externa

- Os algoritmos para ordenação externa devem reduzir o número de passadas sobre o arquivo.
- Uma boa medida de complexidade de um algoritmo de ordenação por intercalação é o número de vezes que um item é lido ou escrito na memória auxiliar.
- Os bons métodos de ordenação geralmente envolvem no total menos do que dez passadas sobre o arquivo.

(ZIVIANI, 2010)

Ordenação Externa - Fitas

- Por décadas as fitas magnéticas eram o dispositivo comum para memória secundária e, portanto, utilizadas para a ordenação externa
- Até hoje, mesmo que discos sejam usados para a ordenação, pensar em ordenação usando fitas é uma boa abstração para o entendimento dos algoritmos de ordenação externa:
 - Há um conjunto de fitas de entrada com dados a serem ordenados (intercalados)
 - Há um conjunto de fitas de saída que recebem o resultado da intercalação

Pense nessas fitas como sendo cada uma delas um disco com acesso sequencial (se não tiver o número de discos necessários, trechos sequenciais desse disco)

Ordenação Externa

Principais abordagens gerais de ordenação externa

- Intercalação Balanceada
- Intercalação usando Seleção por Substituição
- Intercalação Polifásica

Ordenação Externa

Intercalação Balanceada de Vários Caminhos

- Considere um arquivo armazenado em uma fita de entrada:

INTERCALACA OBALANCEADA

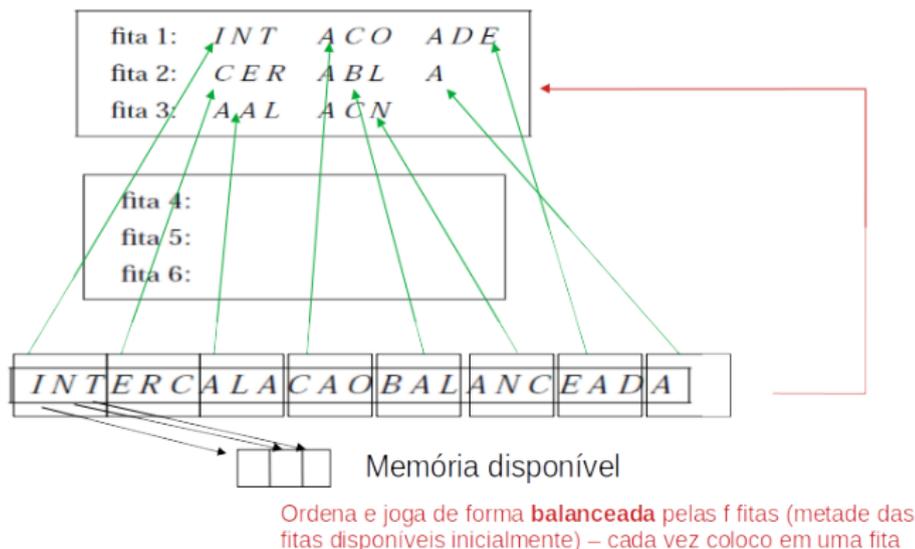
- Objetivo:
 - Ordenar os 22 registros e colocá-los em uma fita de saída.
- Os registros são lidos um após o outro.
- Considere uma memória interna com capacidade para para três registros.
- Considere que esteja disponível seis unidades de fita magnética.

(ZIVIANI, 2010)

Ordenação Externa

Intercalação Balanceada de Vários Caminhos

- Fase de criação dos segmentos ordenados (**corridas**).



Ordenação Externa

Intercalação Balanceada de Vários Caminhos

- Fase de intercalação - Primeira passada:
 1. O primeiro registro de cada fita é lido.
 2. Retire o registro contendo a menor chave.
 3. Armazene-o em uma fita de saída.
 4. Leia um novo registro da fita de onde o registro retirado é proveniente.
 5. Ao ler o terceiro registro de uma corrida sua fita fica inativa.
 6. A fita é reativada quando o terceiro registro das outras fitas forem lidos.
 7. Neste instante 1 corrida de nove registros ordenados foi formada na fita de saída.
 8. Repita o processo para as corridas restantes.

(ZIVIANI, 2010)

Ordenação Externa

Registros: INTERCALACA0BALANCEADA

Fita 1:

Fita 2:

Fita 3:

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: **INTERCALACA**OBALANCEADA

Fita 1: INT

Fita 2:

Fita 3:

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT

Fita 2: CER

Fita 3:

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALA**CA**OBALANCEADA

Fita 1: INT

Fita 2: CER

Fita 3: AAL

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALA**CA**BALANCEADA

Fita 1: INT ACO

Fita 2: CER

Fita 3: AAL

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO

Fita 2: CER ABL

Fita 3: AAL

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACA**OBALANCEADA**

Fita 1: INT ACO

Fita 2: CER ABL

Fita 3: AAL ACN

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL

Fita 3: AAL ACN

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4:

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: A

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AA

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AAC

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACE

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: ACEI

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACEIL

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACEILN

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: ACEILNR

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACEILNRT

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACEILNRT

Fita 5:

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACEILNRT

Fita 5: AAABCCLNO

Fita 6:

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: INT ACO ADE

Fita 2: CER ABL A

Fita 3: AAL ACN

Fita 4: AACEILNRT

Fita 5: AAABCCLNO

Fita 6: AADE

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1:

Fita 2:

Fita 3:

Fita 4: AACEILNRT

Fita 5: AAABCCLNO

Fita 6: AADE

Ordenação Externa

Registros: INTERCALACAOBALANCEADA

Fita 1: AAAAAAABCCCDEEILLNORT

Fita 2:

Fita 3:

Fita 4: AACEILNRT

Fita 5: AAABCCLNO

Fita 6: AADE

Ordenação Externa

Intercalação Balanceada de Vários Caminhos

- Quantas passadas são necessárias para ordenar um arquivo de tamanho arbitrário?
 - Seja n , o número de registros do arquivo.
 - Suponha que cabem m registros na memória interna.
 - A primeira etapa produz n/m corridas ordenadas.
 - Seja $P(n)$ o número de passadas para a fase de intercalação.
 - Seja f o número de fitas utilizadas em cada passada.
 - Assim:

$$P(n) = \left\lceil \log_f \left\lceil \frac{n}{m} \right\rceil \right\rceil$$

No exemplo acima, $n=22$, $m=3$ e $f=3$ temos:

$$P(n) = \left\lceil \log_3 \left\lceil \frac{22}{3} \right\rceil \right\rceil = 2.$$

(ZIVIANI, 2010)

Ordenação Externa

Intercalação Balanceada de Vários Caminhos

- No exemplo foram utilizadas $2f$ fitas para uma intercalação-de- f -caminhos.
- É possível usar apenas $f + 1$ fitas:
 - Encaminhe todos os blocos para uma única fita.
 - Redistribua as corridas entre as fitas de onde elas foram lidas.
 - O custo envolvido é uma passada a mais em cada intercalação.
- No caso do exemplo de 22 registros, apenas quatro fitas seriam suficientes:
 - A intercalação das corridas a partir das fitas 1, 2 e 3 seria toda dirigida para a fita 4.
 - Ao final, a segunda e a terceira corridas ordenadas de nove registros seriam transferidas de volta para as fitas 1 e 2, e a fita 3 usada como fita de saída

(ZIVIANI, 2010)

Ordenação Externa em Disco

- Na ordenação externa em fita, tenho $f+1$ fitas distintas \rightarrow cada uma sendo lida sequencialmente
- Na ordenação externa em disco:
 - poderia semelhantemente utilizar $f+1$ discos, cada um deles sendo lido sequencialmente
 - ou se não tiver vários discos (pelo menos não tantos quantos eu desejaria para um dado f), “simular essas f fitas” em f cilindros distintos
 - O problema é que para ler de cilindros distintos tenho que fazer um novo seek, então melhor já trazer e processar pelo menos todo o bloco (e não apenas um registro)

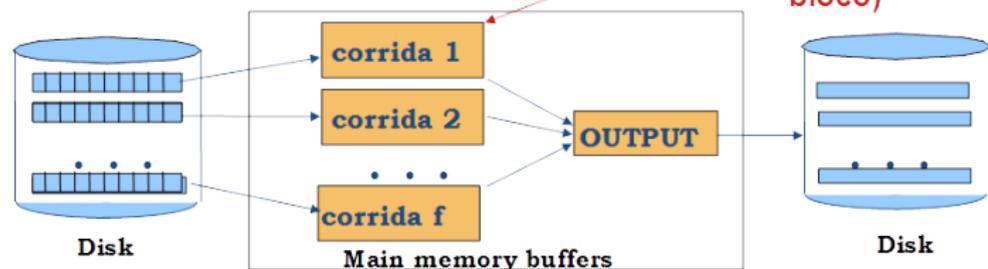
Ordenação Externa em Disco

Tamanhos considerados:

- **N**: tamanho do arquivo original **em número de blocos**
- **M**: Tamanho da memória interna disponível para a ordenação **em número de blocos**

Ordenação Externa em Disco

- Intercala f corridas de cada vez



Na verdade, apenas um trecho da corrida i (no mínimo um bloco)

- No máximo, $f = M - 1$ (M = número de blocos da memória disponível, 1 bloco será utilizado para saída – resultado da intercalação)

Intercalação em f vias ($f \leq M$)

- Fase de ordenação (quanto maiores as corridas iniciais, melhor!):
 - Lê M blocos de cada vez (lota a memória) e ordena formando uma corrida
 - $\lceil N/M \rceil$ corridas
- Fase de intercalação:
 - M blocos de memória precisam ser divididos para as f corridas \rightarrow cada via conterá M/f blocos (ou seja, $1/f$ da corrida) \rightarrow cada via deverá ser lida f vezes
 - Cada passo (passada sobre todo o arquivo):
 - Leitura: necessários $f * f$ acessos ao disco (f vias f vezes cada, sendo cada vez uma leitura sequencial de um bloco)
 - CPU: Cada intercalação $O(M)$ \rightarrow $O(N)$ no total
 - Escrita: N/b seeks (b = número de blocos do buffer de resultado)
 - Número de passos (incluindo geração inicial) = $1 + P(N) =$

Intercalação em f vias ($f \leq M$)

- Fase de ordenação (quanto maiores as corridas iniciais, melhor!):
 - lê M blocos de cada vez (lota a memória) e ordena formando uma corrida
 - $\lceil N/M \rceil$ corridas
- Fase de intercalação:
 - M blocos de memória precisam ser divididos para as f corridas \rightarrow cada via conterá M/f blocos (ou seja, $1/f$ da corrida) \rightarrow cada via deverá ser lida f vezes
 - Cada passo (passada sobre todo o arquivo): \rightarrow **Apenas f seeks se cada corrida estivesse em um disco**
 - Leitura: necessários $f \cdot f$ acessos ao disco (f vias f vezes cada, sendo cada vez uma leitura sequencial de um bloco)
 - CPU: Cada intercalação $O(M)$ $\rightarrow O(N)$ no total
 - Escrita: N/b seeks (b = número de blocos do buffer de resultado)
 - Número de passos (incluindo geração inicial) = $1 + P(N) \approx 1 + \lceil (\log_2 \lceil (N/M) \rceil) \rceil$
 - Aumentar o f diminui o nr de passos mas aumenta o número de seeks de leituras
 - $f = M \rightarrow$ não há leitura sequencial da corrida se ela não estiver em um disco dedicado a ela
 - **Custo = $2 \cdot N \cdot (1 + P(N))$ seeks** (leitura e escrita de cada bloco $P(N)+1$ vezes)

Obs: necessário usar um algoritmo de ordenação interna que ordene *in loco* (sem usar vetor auxiliar)

Intercalação em f vias ($f \leq M$)

Exemplo: número de passos $P(N)$ de intercaulações em f vias para $M = f+1$ blocos

N	f=2	f=4	f=8	f=16	f=128	f=256
100	7	4	3	2	1	1
1,000	10	5	4	3	2	2
10,000	13	7	5	4	2	2
100,000	17	9	6	5	3	3
1,000,000	20	10	7	5	3	3
10,000,000	23	12	8	6	4	3
100,000,000	26	14	9	7	4	4
1,000,000,000	30	15	10	8	5	4

Referência

- Slides baseados no material da profa. Ariane Machado Lima - ACH2024
- ELMARIS, R.; NAVATHE, S. B. Fundamentals of Database Systems. 4 ed. Ed. Pearson-Addison Wesley. Cap 13.8. 4 ed. Pearson. 2004
- RAMAKRISHNAN, R.; GEHRKE, J. Database Management Systems. 3 ed. Ed. McGraw-Hill. 2002 Cap 13
- ZIVIANI, N. Projeto de Algoritmos. Com Implementações em Pascal e C. Cengage. 3ª Edição, 2010.

Algoritmos e Estruturas de Dados II

Aula 28 – Processamento Cossequencial e Ordenação Externa

Prof. Luciano A. Digiampietri
digiampietri@usp.br
@digiampietri

2024