

UNIVERSIDADE DO MINHO

MIEI - GRUPO 49

LABORATÓRIOS DE INFORMÁTICA 3

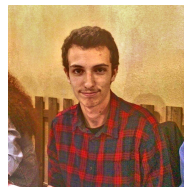
---

## Projeto Java LI3

---



(a) Sérgio Alves A78296



(b) Hugo Brandão A78582



(c) Tiago Alves A78218

# Conteúdo

<b>1</b>	<b>Introdução</b>	<b>2</b>
<b>2</b>	<b>Desenho da Solução</b>	<b>3</b>
<b>3</b>	<b>Classes</b>	<b>4</b>
3.1	Article . . . . .	4
3.2	Revision . . . . .	4
3.3	Contributor . . . . .	5
3.4	Structure . . . . .	5
3.5	QueryEngineImpl . . . . .	5
3.6	ArticleCharComparator . . . . .	6
3.7	ArticleWordComparator . . . . .	6
3.8	ContributorComparator . . . . .	7
3.9	RevisionComparator . . . . .	7
3.10	Exceptions . . . . .	8
3.10.1	NoArticlesException . . . . .	8
3.10.2	NoContributorsException . . . . .	8
3.10.3	NoRevisionsException . . . . .	8
<b>4</b>	<b>Conclusão</b>	<b>9</b>

# 1 Introdução

Este relatório explica como foi feita a implementação do mesmo projeto realizado na primeira parte desta disciplina, desta feita na linguagem de programação Java.

O projeto consistia numa implementação que permita analisar os artigos presentes em backups da Wikipedia, e retirar desses artigos todos os dados necessários para responder a uma lista de queries dada pelos professores. Contudo, tornou-se mais simples de realizar já com a parte C elaborada.

Para isto, tínhamos de criar estruturas de dados para armazenar toda a informação, e depois de carregada, estaríamos então em condições para executar as queries.

Dito isto, neste relatório são explicadas as nossas decisões tomadas para a elaboração deste projeto de Laboratórios de Informática,

## 2 Desenho da Solução

Para esta segunda parte do trabalho, e para responder corretamente às queries pretendidas, optámos por uma implementação simples mas eficiente. Dividimos o trabalho por módulos, ou seja, separamos o código por pequenos pedaços reutilizáveis de tal forma que seja fácil tanto reaproveitar estes pedaços de código, como fazer a manutenção sem que tenha impacto direto nos outros. Deste modo, a programação torna-se mais inteligente e menos suscetível a erros.

Na nossa implementação temos 5 classes fundamentais, sendo elas a *Parser*, que realiza o parse de todos os ficheiros, fornecidos pela equipa docente, para a nossa estrutura de dados, a *Article*, que guarda toda a informação necessária de um artigo, a *Revision*, com os dados de cada revisão, a *Contributor*, relativamente a cada contribuidor, e a *Structure*. A classe *Structure* define as estruturas onde vão ser guardadas todas as informações que queremos sobre os artigos. Nesta, utilizámos dois *TreeMap*, sendo um para os artigos e outro para os contribuidores. Optámos por esta decisão, devido às condições que um *TreeMap* nos oferece. Temos uma "chave", que é o id de cada artigo, associada ao respetivo artigo. Quando for necessário procurar por um deles, por exemplo, na query *contributor\_name*, basta aceder à estrutura e fazer o "get(contributor\_id)", em que *contributor\_id* corresponde ao id do contribuidor pretendido, e temos imediatamente a resposta desejada, com o retornar do nome desse contribuidor. O mesmo se aplica à escolha da estrutura para os artigos. Além disto, utilizámos um inteiro que nos dá o total de artigos nos *snapshots*.

Em relação ao *parser* escolhemos *stax*, pois é de simples implementação, com código legível, e com alguma rapidez na sua execução. Quanto à resolução das queries, foram feitas praticamente todas elas com recurso a streams, o que torna a solução rápida e eficaz.

## 3 Classes

Na seguinte secção vamos explicar as classes que decidimos criar para a resolução do projeto, incluindo diagramas e explicação das estruturas escolhidas para cada classe.

### 3.1 Article

Para cada artigo é necessário guardar toda a sua informação, mais concretamente, o seu id, o número de revisões que cada artigo tem, e uma estrutura para guardar todas as revisões. Para esta última, utilizámos um *TreeSet* proporcionando assim uma pesquisa rápida.

```
1 public class Article{
2     private int nrevisions;           /* Numero de revisoes do
        artigo */
3     private long id;                 /* 0 id do artigo */
4     private TreeSet<Revision> revList; /* Lista de revisoes */
5     ...
6 }
```

### 3.2 Revision

Esta classe guarda tudo o que é necessário sabermos sobre uma revisão de um dado artigo. Aqui guardamos o seu id, o id do contribuidor dessa revisão, o título, data, número de palavras e número de caracteres.

É uma classe simples, legível mas de grande importância para a execução correta deste projeto.

```
1 public class Revision{
2     private long id;                 /* id of the revision */
3     private long idcontributor;      /* id of the contributor that
        wrote the */
4     private String title;            /* title of the revision */
5     private String date;             /* date of the creation of the
        revisions */
6     private int nwords;              /* number of words of the revision's
        article */
7     private int nchars;              /* number of characters of the revision
        's article */
8 }
```

### 3.3 Contributor

Aqui, estão todos os dados de um contribuidor, sendo eles o seu id, o número de contribuições que realizou em todos os artigos que nos foram fornecidos, e o seu nome.

Implementámos um método *compareNContributions* que é útil para o comparador de contribuidores, chamado na query *top\_10\_contributors*.

```
1 public class Contributor{
2     private long id;                /* id do contribuidor */
3     private int ncontributions;    /* Numero de contribuicoes */
4     private String name;          /* Nome do contribuidor */
5 }
```

### 3.4 Structure

Esta é a classe principal do projeto. Aqui é guardada a informação de todos os artigos, revisões e contribuidores. Tal como foi dito acima, na secção *Desenho da solução*, usámos um *TreeMap* para os artigos, e um outro para os contribuidores.

O método fundamental aqui é o *insert*, que faz exatamente o que o nome diz, isto é, insere toda a informação nas estruturas, artigos para o *articles* e contribuidores para o *contributors*. Desta forma, bastava apenas, nas queries, aceder a esta classe e implementar algoritmos para solucioná-las.

```
1 public class Structure{
2     private int allArticles;        /* Numero total de
3         artigos */
4     private Map<Long,Article> articles; /* Lista de artigos
5         */
6     private Map<Long,Contributor> contributors; /* Lista de
7         contribuidores */
8 }
```

### 3.5 QueryEngineImpl

Nesta classe está toda a resolução das queries pretendidas. Optámos por, na maior parte, usar streams. Esta escolha deve-se ao facto de, desta forma, a execução tornar-se mais rápida e eficiente. *All\_Revisions*, *Top\_10\_Contributors*, *Top\_20\_Largest\_Articles*, *Top\_N\_Articles\_With\_More\_Words*, e *Titles\_With\_Prefix* foram implementadas com streams. As restantes, não sendo feitas desta forma, foram igualmente elaboradas de modo que a sua execução fosse também eficiente.

```
1 public long all_revisions() {
2     try{
3         return this.struct.getArticles().values().stream()
4             .mapToLong(Article::getNRevisions)
5             .sum();
6     }catch(NoArticlesException e){System.out.println("Sem artigos");}
7     return 0;
8 }
```

### 3.6 ArticleCharComparator

Esta class serve principalmente para ordenar um artigo segundo o número de caracteres da revisão mais recente. É necessário para a query *top\_twenty\_Largest\_Articles*, como se pode ver no excerto de código abaixo. Ordenamos a stream de artigos através deste comparador, obtendo uma stream ordenada como pretendíamos.

```
1 public ArrayList<Long> top_20_largest_articles() {
2     ArrayList<Long> top_twenty_ids = new ArrayList<Long>(20);
3     try{
4         this.struct.getArticles().values()
5             .stream()
6             .sorted(new ArticleCharComparator())
7             .limit(20)
8             .forEach(a -> top_twenty_ids.add((a.
9                 getId())));
10    }catch (NoArticlesException e){System.out.println("Sem artigos");}
11    return top_twenty_ids;
12 }
```

### 3.7 ArticleWordComparator

Do modo que funciona a classe anterior, o mesmo se aplica a esta, mas em vez de ordenar por número de caracteres, ordena por número de palavras. Esta classe é aplicável na query *top\_N\_Articles\_With\_More\_Words*.

```
1 public ArrayList<Long> top_N_articles_with_more_words(int n) {
2     ArrayList<Long> top_n_ids = new ArrayList<Long>(n);
3     try{
4         this.struct.getArticles().values().stream()
5             .sorted(new
6                 ArticleWordComparator())
7             .limit(n)
8             .forEach(a -> top_n_ids.add((a
9                 .getId())));
10    }catch (NoArticlesException e){System.out.println("Sem artigos");}
11    return top_n_ids;
12 }
```

```
1 public class ArticleWordComparator implements Comparator<Article> {
2     public int compare(Article a1, Article a2){
3         return Integer.valueOf(a2.getBiggestRevision()
4             .getNWords()).compareTo(Integer.valueOf(a1
5             .getBiggestRevision().getNWords()));
6     }
7 }
```

### 3.8 ContributorComparator

Para resolver a query dos 10 contribuidores com mais contribuições, criámos um comparador de contribuidores, sendo este comparador invocado para ordenar uma stream de estes mesmos. A ordenação é feita segundo o número de contribuições de cada contribuidor. Uma vez com a stream ordenada, bastou-nos fazer "limit(10)" e imprimir os resultados.

```
1 public class ContributorComparator implements Comparator<
    Contributor>{
2     public int compare(Contributor c1, Contributor c2){
3         return c1.compareNContributions(c2);
4     }
5 }
```

### 3.9 RevisionComparator

Mais uma vez, esta classe é usada para ajudar na realização de uma das queries, desta feita a *Titles\_With\_Prefix*. O objetivo é ordenar a stream de revisões por título(ordem alfabética). Obviamente, são filtradas as revisões cujo título tem como prefixo a string passada como argumento.

```
1 public class RevisionComparator implements Comparator<Revision> {
2     public int compare(Revision r1, Revision r2){
3         return r1.getDate().compareTo(r2.getDate());
4     }
5 }

1 public ArrayList<String> titles_with_prefix(String prefix) {
2     ArrayList<String> res = new ArrayList<String>();
3     try{
4         this.struct.getArticles().values()
5             .stream()
6             .map(Article::getRecentRevision)
7             .filter(r -> r.getTitle().startsWith(
                prefix))
8             .map(Revision::getTitle)
9             .sorted(new RevisionTitleComparator())
10            .forEach(title -> res.add(title));
11 }catch(NoArticlesException e){System.out.println("Sem artigos");}
12 return res;
13 }
```



## 3.10 Exceptions

Como muitas das queries necessitam de ter acesso às estruturas, e visto que essas estruturas podem estar vazias (null), temos que prevenir erros futuros. Para isto, e como se mostra a seguir, criámos exceptions para os artigos, contribuidores e revisões. Quando uma query invocar um método que lance um destas exceptions, é imprimida uma mensagem de erro.

### 3.10.1 NoArticlesException

```
1 public class NoArticlesException extends Exception implements
   Serializable{
2     public NoArticlesException(){ super(); }
3
4     public NoArticlesException(String message){super(message);}
5 }
```

### 3.10.2 NoContributorsException

```
1 public class NoContributorsException extends Exception implements
   Serializable
2 {
3     public NoContributorsException(){super();}
4
5     public NoContributorsException(String message){super(message);}
6 }
```

### 3.10.3 NoRevisionsException

```
1 public class NoRevisionsException extends Exception implements
   Serializable
2 {
3     public NoRevisionsException(){super();}
4
5     public NoRevisionsException(String message){super(message);}
6 }
```

## 4 Conclusão

Em suma, o trabalho foi realizado com sucesso apesar de alguns percalces. Tivemos de nos adaptar a um novo estilo de paradigma, ser eficazes e implementar as soluções mais ótimas para conseguir uma performance que satisfizesse os objetivos do grupo.

Achamos, no entanto, que este trabalho nos ajudou bastante a entender como a otimização e o processamento de grandes volumes de dados não é igual em todos os paradigmas da programação, tal como os seus cuidados.

Após a conclusão deste trabalho, pensamos que todos os elementos deste grupo saem com uma percepção mais clara do que é ser eficiente na hora de aceder a estruturas de dados e de processar grande quantidade de informação.

Ficamos, então, contentes com o nosso desempenho e com o resultado obtido, acreditando que todos os objetivos foram cumpridos.