# UCE Indexation & recherche

Vincent Labatut

## TP 2 | Correction

# 1 Construction

### Exercice 1

```java
private int filterTokens(List<Token> tokens)
{  int result = 0;
   // on passe chaque paire de tokens consécutifs en revue
   Iterator<Token> it = tokens.iterator();
   Token t1 = null;
   while(it.hasNext())
   {  Token t2 = it.next();
      if(t1==null)
         result = 1;
      else
      {  // si deux tokens consécutifs sont identiques,
         // on supprime le second
         if(t1.equals(t2))
            it.remove();
         // on compte les termes
         String type1 = t1.getType();
         String type2 = t2.getType();
         if(!type1.equals(type2))
            result++;
      }
      t1 = t2;
   }
   return result;
}
```

### Exercice 2

```java
private int buildPostings(List<Token> tokens, AbstractIndex index)
{  int result = 0;
   int i = 0;
   IndexEntry entry = null;
   // on traite chaque token séparément
   for(Token token: tokens)
   {  String type = token.getType();
      // si besoin, on crée une nouvelle entrée
      if(entry==null || !entry.getTerm().equals(type))
      {  entry = new IndexEntry(type);
         index.addEntry(entry,i);
         i++;
      }

      // dans tous les cas, on met àjour la liste de postings
      int docId = token.getDocId();
      Posting posting = new Posting(docId);
      entry.addPosting(posting);
      result++;
```

```
20      }
21
22      return result;
23    }
```

### Exercice 3

```
1  public AbstractIndex buildIndex(List<Token> tokens, LexiconType lexiconType)
2  { int indexSize;
3    AbstractIndex result = null;
4
5    { System.out.println(" Sorting tokens...");
6      Collections.sort(tokens);
7      System.out.println(" "+tokens.size()+" tokens sorted\n");
8    }
9
10   { System.out.println(" Filtering tokens...");
11     indexSize = filterTokens(tokens);
12     System.out.println(" "+tokens.size()+" tokens remaining, corresponding
13                                     to "+indexSize+" terms\n");
14   }
15
16   { System.out.println(" Building posting lists...");
17     switch(lexiconType)
18     { case ARRAY:
19         result = new ArrayIndex(indexSize);
20         break;
21       case HASH:
22         result = new HashIndex(indexSize);
23         break;
24       case TREE:
25         result = new TreeIndex();
26         break;
27     }
28     int postingNumber = buildPostings(tokens,result);
29     System.out.println(" "+postingNumber+" postings listed, lexicon="
30                                     +lexiconType);
31   }
32
33   return result;
34 }
```

## 2 Indexation

### Exercice 4

```
1  public static AbstractIndex indexCorpus(TokenListType tokenListType,
2              LexiconType lexiconType) throws UnsupportedEncodingException
3  { AbstractIndex result;
4    Tokenizer tokenizer;
5    Normalizer normalizer;
6    List<Token> tokens = null;
7    switch(tokenListType)
8    { case ARRAY:
9        tokens = new ArrayList<Token>();
10       break;
11     case LINKED:
12       tokens = new LinkedList<Token>();
13       break;
14   }
```

```
15     int docNbr;
16
17     // tokénization
18     { System.out.println("Tokenizing corpus...");
19       tokenizer = new Tokenizer();
20       docNbr = tokenizer.tokenizeCorpus(tokens);
21       System.out.println(tokens.size()+" tokens were found\n");
22     }
23
24     // normalisation
25     { System.out.println("Normalizing tokens...");
26       normalizer = new Normalizer();
27       normalizer.normalizeTokens(tokens);
28       System.out.println(tokens.size()+" tokens remaining after
29                                            normalization\n");
30     }
31
32     // construction
33     { System.out.println("Building index...");
34       Builder builder = new Builder();
35       result = builder.buildIndex(tokens, lexiconType);
36       System.out.println("There are "+result.getSize()+" entries in the
37                             index, token list="+tokenListType+"\n");
38     }
39
40     // affichage
41     System.out.println("Content of the final index :");
42     result.print();
43
44     // finalisation de l'index
45     result.tokenizer = tokenizer;
46     result.normalizer = normalizer;
47     result.docNbr = docNbr;
48     return result;
49   }
```

## Exercice 5

```
1   public static void main(String[] args) throws Exception
2   { // configutation de l'outil
3     Configuration.setCorpusName("wp");
4
5     // test du processus d'indexation
6     testIndexation();
7   }
8
9   private static void testIndexation() throws IOException
10  { // construction de l'index
11    AbstractIndex index = AbstractIndex.indexCorpus(TokenListType.LINKED,
12                                         LexiconType.ARRAY);
13  }
```

## Exercice 6

```
1   public static List<String> getFileNamesFromPostings(List<Posting> postings)
2   { List<String> result = new LinkedList<String>();
3     File folder = new File(getCorpusFolder());
4     String fn[] = folder.list();
5     Arrays.sort(fn);
6
```

```
7      for(Posting posting: postings)
8      {  int docId = posting.getDocId();
9         String name = fn[docId];
10        result.add(name);
11     }
12
13     return result;
14   }
```

## 3 Durée de traitement

### Exercice 8

Modification de la méthode `buildIndex` :

```
1   public AbstractIndex buildIndex(List<Token> tokens, LexiconType lexiconType)
2   {  int indexSize;
3      AbstractIndex result = null;
4
5      {  System.out.println(" Sorting tokens...");
6         long start = System.currentTimeMillis();
7         Collections.sort(tokens);
8         long end = System.currentTimeMillis();
9         System.out.println(" "+tokens.size()+" tokens sorted,
10                                         duration="+(end-start)+" ms\n");
11     }
12
13     {  System.out.println(" Filtering tokens...");
14        long start = System.currentTimeMillis();
15        indexSize = filterTokens(tokens);
16        long end = System.currentTimeMillis();
17        System.out.println(" "+tokens.size()+"tokens remaining, corresponding
18                    to "+indexSize+" terms, duration="+(end-start)+" ms\n");
19     }
20
21     {  System.out.println(" Building posting lists...");
22        long start = System.currentTimeMillis();
23        switch(lexiconType)
24        {  case ARRAY:
25              result = new ArrayIndex(indexSize);
26              break;
27           case HASH:
28              result = new HashIndex(indexSize);
29              break;
30           case TREE:
31              result = new TreeIndex();
32              break;
33        }
34        int postingNumber = buildPostings(tokens,result);
35        long end = System.currentTimeMillis();
36        System.out.println(" "+postingNumber+" postings listed,
37                    lexicon="+lexiconType+",  duration="+(end-start)+" ms");
38     }
39
40     return result;
41   }
```

Modification de la méthode `indexCorpus` :

```
1   public static AbstractIndex indexCorpus(TokenListType tokenListType,
2                LexiconType lexiconType) throws UnsupportedEncodingException
```

```
3   { AbstractIndex result;
4     Tokenizer tokenizer;
5     Normalizer normalizer;
6     List<Token> tokens = null;
7     switch(tokenListType)
8     { case ARRAY:
9         tokens = new ArrayList<Token>();
10        break;
11      case LINKED:
12        tokens = new LinkedList<Token>();
13        break;
14    }
15    int docNbr;
16    long startTotal = System.currentTimeMillis();
17
18    // tokénization
19    { System.out.println("Tokenizing corpus...");
20      long start = System.currentTimeMillis();
21      tokenizer = new Tokenizer();
22      docNbr = tokenizer.tokenizeCorpus(tokens);
23      long end = System.currentTimeMillis();
24      System.out.println(tokens.size()+" tokens were found,
25                          duration="+(end-start)+" ms\n");
26    }
27
28    // normalisation
29    { System.out.println("Normalizing tokens...");
30      long start = System.currentTimeMillis();
31      normalizer = new Normalizer();
32      normalizer.normalizeTokens(tokens);
33      long end = System.currentTimeMillis();
34      System.out.println(tokens.size()+" tokens remaining after
35                normalization, duration="+(end-start)+" ms\n");
36    }
37
38    // construction
39    { System.out.println("Building index...");
40      long start = System.currentTimeMillis();
41      Builder builder = new Builder();
42      result = builder.buildIndex(tokens);
43      long end = System.currentTimeMillis();
44      System.out.println("There are "+result.getSize()+" entries in the
45      index, token list="+tokenListType+", duration="+(end-start)+" ms\n");
46    }
47
48    long endTotal = System.currentTimeMillis();
49    System.out.println("Total duration="+(endTotal-startTotal)+" ms\n");
50
51    // affichage
52    System.out.println("Content of the final index :");
53    result.print();
54
55    // finalisation de l'index
56    result.normalizer = normalizer;
57    result.tokenizer = tokenizer;
58    result.docNbr = docNbr;
59    return result;
60  }
```

## Exercice 9

**Remarque :** Tous les résultats présentés ici sont obtenus avec le corpus `wp`.

Utilisation d'une liste à base de tableau pour stocker les tokens :

```
Tokenizing corpus...
7612885 tokens were found, duration=10498 ms

Normalizing tokens...
7612885 tokens remaining after normalization, duration=11648 ms

Building index...
Sorting tokens...
7612885 tokens sorted, duration=5322 ms

Filtering tokens...
2177994 tokens remaining, corresponding to 142792 terms, duration=16504892 ms

Building posting lists...
2177994 postings listed, lexicon=ARRAY, duration=415 ms
There are 142792 entries in the index, token list=ARRAY, duration=16510630 ms

Total duration=16532776 ms
```

Utilisation d'une liste chaînée pour stocker les tokens :

```
Tokenizing corpus...
7612885 tokens were found, duration=21879 ms

Normalizing tokens...
7612885 tokens remaining after normalization, duration=15493 ms

Building index...
Sorting tokens...
7612885 tokens sorted, duration=6362 ms

Filtering tokens...
2177994 tokens remaining, corresponding to 142792 terms, duration=1029 ms

Building posting lists...
2177994 postings listed, lexicon=ARRAY, duration=559 ms
There are 142792 entries in the index, token list=LINKED, duration=7954 ms

Total duration=45326 ms
```

**Conclusion.** pour le cas traité ici, les listes chaînées sont beaucoup plus efficaces. La différence se fait surtout lors du filtrage des tokens. En effet, `ArrayList` est basée sur un tableau et nécessite des redimensionnements, ce qui est régulièrement le cas lors du filtrage. Or, le redimensionnement nécessite de recopier l'intégralité du tableau à une adresse différente, ce qui est coûteux en temps.

Comparons maintenant les structures de données proposées pour le lexique, en utilisant des listes chaînées pour stocker les tokens durant l'indexation (puisque c'est le type de liste qui semble être le plus adapté à ce cas-là).

Table de hachage :

```
Tokenizing corpus...
7612885 tokens were found, duration=18166 ms
```

```
Normalizing tokens...
7612885 tokens remaining after normalization, duration=20696 ms

Building index...
Sorting tokens...
7612885 tokens sorted, duration=6372 ms

Filtering tokens...
2177994 tokens remaining, corresponding to 142792 terms, duration=1012 ms

Building posting lists...
2177994 postings listed, lexicon=HASH, duration=512 ms
There are 142792 entries in the index, token list=LINKED, duration=7902 ms

Total duration=46764 ms
```

Arbre de recherche :

```
Tokenizing corpus...
7612885 tokens were found, duration=17959 ms

Normalizing tokens...
7612885 tokens remaining after normalization, duration=27539 ms

Building index...
Sorting tokens...
7612885 tokens sorted, duration=6136 ms

Filtering tokens...
2177994 tokens remaining, corresponding to 142792 terms, duration=1065 ms

Building posting lists...
2177994 postings listed, lexicon=TREE, duration=568 ms
There are 142792 entries in the index, token list=LINKED, duration=7772 ms

Total duration=53272 ms
```

Pour le tableau, l'insertion se fait en temps constant, puisque on reçoit directement la position concernée. Pour la table de hachage, en théorie l'insertion s'y fait (en moyenne) en $O(1)$, c'est-à-dire là aussi en temps constant : ceci explique la proximité des temps mesurés pour ces deux structures de données. Pour l'arbre de recherche, l'insertion est (en moyenne) en $O(\log n)$, ce qui peut expliquer le temps plus élevé que pour les deux autres structures de données.

# 4 Stockage

### Exercice 10

```java
public void write() throws IOException
{ System.out.println("Writing the index");
    long start = System.currentTimeMillis();

    String fileName = FileTools.getIndexFile();
    File file = new File(fileName);
    FileOutputStream fos = new FileOutputStream(file);
    ObjectOutputStream oos = new ObjectOutputStream(fos);

    oos.writeObject(this);

    oos.close();
```

```
13
14      long end = System.currentTimeMillis();
15      System.out.println("Index written, duration="+(end-start)+" ms\n");
16    }
```

### Exercice 11

```
1    public static AbstractIndex read() throws IOException, ClassNotFoundException
2    { System.out.println("Loading the index");
3      long start = System.currentTimeMillis();
4
5      String fileName = FileTools.getIndexFile();
6      File file = new File(fileName);
7      FileInputStream fis = new FileInputStream(file);
8      ObjectInputStream ois = new ObjectInputStream(fis);
9
10     AbstractIndex result = (AbstractIndex) ois.readObject();
11
12     ois.close();
13
14     long end = System.currentTimeMillis();
15     System.out.println("Index loaded, duration="+(end-start)+" ms\n");
16     return result;
17   }
```