

Laporan Pemerolehan Informasi UTS 2

Pencarian menggunakan Inverted Index

Dosen Pengampu : JB. Budi Darmawan S.T., M.Sc.



Disusun oleh

Nama : Dyline Melynea Fernandez

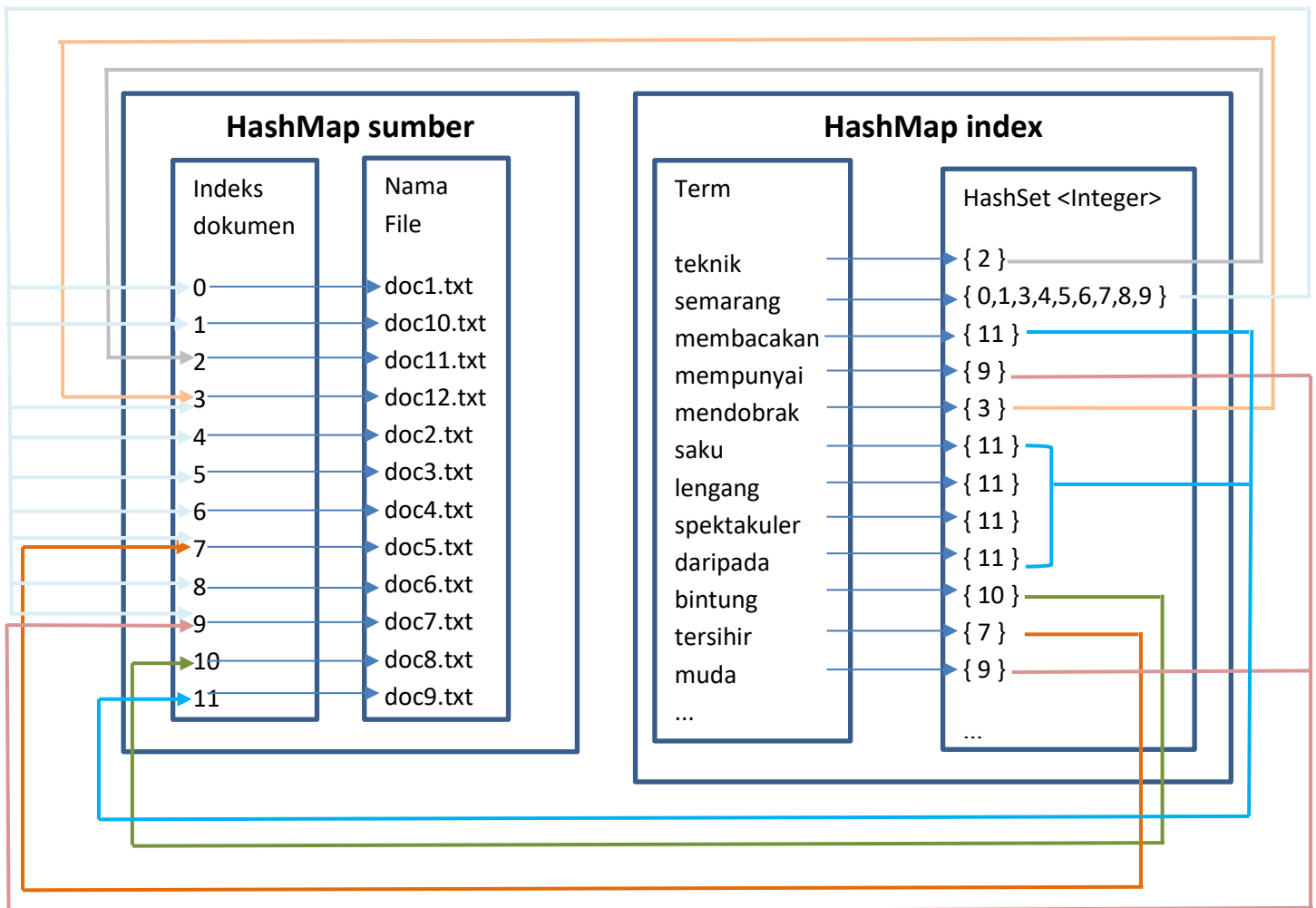
NIM : 185314125

PROGRAM STUDI INFORMATIKA
FAKULTAS SAINS & TEKNOLOGI
UNIVERSITAS SANATA DHARMA

A. Rancangan struktur data

Untuk **index** dan **sumber** menggunakan HashMap yang prinsipnya seperti kamus/dictionary. Namun value dari HashMap **index** dibuat struktur datanya yaitu HashSet yang merupakan himpunan. Sehingga value dari key HashMap **index** bisa lebih dari 1 elemen. **index** merupakan pemetaan antara term sebagai key dan id dokumen/indeks dokumen sebagai valuenya, **sumber** merupakan pemetaan antara id dokumen/indeks dokumen sebagai key dan nama file dokumen sebagai valuenya.

Ilustrasi:



Konektor warna-warni hanyalah gambaran di method search, setelah didapatkan himpunan integer, maka memanggil HashMap **sumber** untuk mendapatkan nama file dokumen.

Gambaran output HashMap **sumber** dan **index** di program

```
0=doc1.txt, 1=doc10.txt, 2=doc11.txt, 3=doc12.txt, 4=doc2.txt, 5=doc3.txt, 6=doc4.txt, 7=doc5.txt, 8=doc6.txt, 9=doc7.txt, 10=doc8.txt, 11=doc9.txt
[teknik=[2], semarang=[0, 1, 3, 4, 5, 6, 7, 8, 9], membacakan=[11], mempunyai=[9], mendobrak=[3], saku=[11], lengang=[11], spektakuler=[11],
```

B. Algoritma searching more than 1 keyword

1. Pisahkan setiap term/kata keywords (patokan pisah dengan whitespace)
2. Buat objek HashSet untuk menampung index term (tipe Integer)
3. Cari irisan index dari semua keywords (gunakan fungsi retainAll dari kelas HashSet)
4. Apabila ukuran dari objek hashset yang menampung index term adalah 0, maka cetak “tidak ditemukan”
5. Lakukan perulangan sebanyak jumlah elemen objek HashSet
6. Cetak mapping values dari key yang bersesuaian iterasi foreach objek HashSet

Penjelasan:

Keyword dipisah terlebih dahulu untuk setiap katanya, kemudian membuat objek HashSet untuk menampung indeks term. Untuk awal pertama membuat objek, isi parameter konstruktor HashSet diisi dengan **index.get(keywords[0].toLowerCase())** untuk inisialisasi. Setelahnya isi dari objek HashSet tersebut (**res**), akan mengalami perulangan sebanyak banyak elemen keyword yang sudah dipisahkan per kata. Pada perulangan tersebut, akan ditetapkan value **res** yang beririsan dengan hasil indeks term.

Kemudian selanjutnya memfilter kondisi jika ukuran objek **res** adalah 0, yang artinya tidak ada index term tersebut, maka cetak “tidak ditemukan”.

Setelah sudah melakukan kondisi di atas, di luar kondisi tersebut dilakukan foreach loop sebanyak elemen **res**. Dan mencetak atribut **sumber** tipe HashMap yang menampung indeks term sebagai key dengan nama filenya sebagai value. Hasil atribut itu adalah dokumen yang mewadahi keyword yang dicari.

C. Program (captured)

```
1 // @author dyline
2
3 import java.io.*;
4 import java.util.*;
5
6 class Index {
7
8     Map<Integer, String> sumber;
9     HashMap<String, HashSet<Integer>> index;
10
11     Index() {
12         sumber = new HashMap<Integer, String>();
13         index = new HashMap<String, HashSet<Integer>>();
14     }
15
16     public void buatIndex(String[] files) {
17         int i = 0;
18         for (String namaFile : files) {
19
20             try (BufferedReader file = new BufferedReader(new FileReader("../invertedindexirsystem\\" + namaFile))) {
21                 + namaFile))) {
22                 BufferedReader stopwords = new BufferedReader(
23                     new FileReader("../invertedindexirsystem\\stopwordsin.txt"));
24
25                 sumber.put(i, namaFile);
26                 String ln;
27                 String readStopwords;
28                 String cleanTerms = "";
29                 while ((ln = file.readLine()) != null) {
30                     readStopwords = stopwords.readLine();
31                     String[] terms = ln.split("\\W+");
32                     String[] stopword = readStopwords.split("\\W+");
33
34                     for (String term : terms) {
35                         term = term.toLowerCase();
36                         if (!term.equals(stopword)) {
```

```

36         cleanTerms = term;
37     }
38     if (!index.containsKey(cleanTerms)) {
39         index.put(cleanTerms, new HashSet<Integer>());
40     }
41     index.get(cleanTerms).add(i);
42 }
43 }
44 } catch (IOException e) {
45     System.out.println("File " + namaFile + " tidak ditemukan. skip");
46 }
47 i++;
48 }
49 index.remove("");
50 System.out.println(sumber.entrySet());
51 System.out.println(index.entrySet());
52 }
53
54 public void search(String key) {
55     String[] keywords = key.split("\\W+");
56     HashSet<Integer> res = new HashSet<Integer>(index.get(keywords[0].toLowerCase()));
57     for (String kw : keywords) {
58         res.retainAll(index.get(kw));
59     }
60
61     if (res.size() == 0) {
62         System.out.println("Tidak ditemukan");
63         return;
64     }
65     System.out.println("Ditemukan di: ");
66     for (int num : res) {
67         System.out.println("\t" + sumber.get(num));
68     }
69 }
70 }
71
72 public class InvertedIndex {
73
74     public static void main(String args[]) throws IOException {
75         Index index = new Index();
76         String path = "..\\invertedindexirsystem\\Koleksi";
77         File name = new File(path);
78         if (name.exists()) {
79             if (name.isDirectory()) {
80                 String directory[] = name.list();
81                 index.buatIndex(directory);
82                 System.out.println("kata kunci: ");
83                 BufferedReader in = new BufferedReader(new InputStreamReader(System.in));
84                 String keywords = in.readLine();
85
86                 index.search(keywords);
87             }
88         }
89     }
90 }

```

Output:

```
Output - invertedindexsystem (run) x InvertedIndex.java x MainInvertedIndexBasic.java x Gabung.java x File.java x HashSet.java x Ab..
run:
[0=doc1.txt, 1=doc10.txt, 2=doc11.txt, 3=doc12.txt, 4=doc2.txt, 5=doc3.txt, 6=doc4.txt, 7=doc5.txt, 8=doc6.txt, 9=doc7.txt, 10=doc8.txt, 11=
[teknik=[2], semarang=[0, 1, 3, 4, 5, 6, 7, 8, 9], membacakan=[11], mempunyai=[9], mendobrak={3}, saku=[11], lengang=[11], spektakuler=[11],
kata kunci:
semarang mendobrak
Ditemukan di:
    doc12.txt
BUILD SUCCESSFUL (total time: 6 seconds)

Output - invertedindexsystem (run) x InvertedIndex.java x MainInvertedIndexBasic.java x Gabung.java x File.java x HashSet.java x Ab..
run:
[0=doc1.txt, 1=doc10.txt, 2=doc11.txt, 3=doc12.txt, 4=doc2.txt, 5=doc3.txt, 6=doc4.txt, 7=doc5.txt, 8=doc6.txt, 9=doc7.txt, 10=doc8.txt, 11=
[teknik=[2], semarang=[0, 1, 3, 4, 5, 6, 7, 8, 9], membacakan=[11], mempunyai=[9], mendobrak=[3], saku=[11], lengang=[11], spektakuler=[11],
kata kunci:
semarang
Ditemukan di:
    doc1.txt
    doc10.txt
    doc12.txt
    doc2.txt
    doc3.txt
    doc4.txt
    doc5.txt
    doc6.txt
    doc7.txt
BUILD SUCCESSFUL (total time: 3 seconds)

run:
[0=doc1.txt, 1=doc10.txt, 2=doc11.txt, 3=doc12.txt, 4=doc2.txt, 5=doc3.txt, 6=doc4.txt, 7=doc5.txt, 8=doc6.txt, 9=doc7.txt, 10=doc8.txt, 11=
[teknik=[2], semarang=[0, 1, 3, 4, 5, 6, 7, 8, 9], membacakan=[11], mempunyai=[9], mendobrak=[3], saku=[11], lengang=[11], spektakuler=[11],
kata kunci:
mendobrak
Ditemukan di:
    doc12.txt
BUILD SUCCESSFUL (total time: 7 seconds)
```