

Penerapan Regex dan String Matching untuk Filter Chat pada E-Commerce

Thea Josephine Halim - 13522012
Program Studi Teknik Informatika
Sekolah Teknik Elektro dan Informatika
Institut Teknologi Bandung, Jalan Ganesha 10 Bandung
E-mail (gmail): 13522012@mahasiswa.itb.ac.id

Abstrak—Dalam era digital yang terus berkembang, e-commerce telah menjadi sektor penting yang mengalami pertumbuhan signifikan. Aplikasi belanja online memungkinkan pengguna untuk melakukan transaksi dengan mudah dan cepat. Namun, dengan meningkatnya penggunaan platform ini, tantangan seperti pencegahan transaksi di luar platform yang dapat merugikan penyedia layanan juga muncul. Untuk menjaga keamanan dan integritas transaksi, banyak platform e-commerce mengimplementasikan filter chat yang mendeteksi dan memblokir penyebaran informasi pribadi seperti nomor telepon, alamat, dan rincian rekening bank, hingga penggunaan kata-kata vulgar.

Kata Kunci—Filter Chat, E-Commerce, Regex, Pencocokan String

Abstract—In the rapidly evolving digital era, e-commerce has become a crucial sector experiencing significant growth. Online shopping applications enable users to conduct transactions easily and quickly. However, with the increasing use of these platforms, challenges such as preventing off-platform transactions that can harm service providers have also emerged. To maintain transaction security and integrity, many e-commerce platforms implement chat filters that detect and block the sharing of private information such as phone numbers, addresses, and bank account details, as well as the use of vulgar words.

Keywords—Chat Filter, E-Commerce, Regex, String Matching

I. PENDAHULUAN

Di dunia maya, kita bisa bertemu banyak orang dengan berbagai latar belakang dan sifat, menjadikan salah satu daya tarik dari kehidupan online. Kita bisa berkomunikasi dengan orang lain tanpa batas dari balik layar. Komunikasi dengan orang asing secara anonim dalam chat (pesan teks) sudah menjadi biasa, dan orang-orang bebas menulis apa saja. Sayangnya, kebebasan ini bisa saja disalahgunakan.

Dalam era digital yang terus berkembang, e-commerce atau aplikasi bisnis online telah menjadi salah satu sektor yang mengalami pertumbuhan pesat karena memungkinkan pengguna untuk melakukan transaksi dengan mudah dan cepat tanpa harus repot-repot mengunjungi toko, dapat dilakukan kapan saja, dan di mana saja. Akan tetapi, seiring dengan meningkatnya penggunaan platform ini, timbul pula berbagai

tantangan, salah satunya adalah upaya untuk mencegah transaksi di luar aplikasi yang dapat merugikan penyedia layanan maupun membahayakan konsumen. Untuk menjaga keamanan dan integritas transaksi, banyak platform e-commerce mengimplementasikan filter chat yang mendeteksi dan memblokir penyebaran informasi pribadi seperti nomor telepon, alamat, dan informasi rekening bank. Pentingnya chat filtering ini guna melindungi pengguna dari potensi penipuan dan memastikan semua transaksi tetap berada dalam sistem aplikasi yang mencegah kerugian pihak e-commerce sendiri. Selain itu, filter chat pada e-commerce juga memastikan agar chat antara penjual dan pembeli tetap sopan dengan filter chat yang dapat mendeteksi kata-kata vulgar ataupun mengandung pornografi.



Gambar 1.1 Ilustrasi Kasus Pelanggaran Chat
Sumber: [Seller Shopee Edu](#)

Salah satu primadona e-commerce akhir-akhir ini adalah Shopee, aplikasi jual beli online asal Singapura yang dimiliki perusahaan Sea Limited. Pada makalah ini akan dibahas mengenai penggunaan RegEx (Regular Expression) dan algoritma *pattern matching* pada simulasi *chat filtering* sederhana pada pesan teks antara penjual dan pembeli di e-commerce Shopee. Berdasarkan petunjuk kebijakan pelayanan pembeli pada laman pusat edukasi penjual Shopee, pesan yang dianggap melanggar meliputi:

- 1) Mengirim *chat* yang tidak pantas, kasar, atau mengandung unsur SARA (suku, agama, ras, dan antar golongan).
- 2) Mengirim *chat* pornografi atau pelecehan seksual.
- 3) Mengirim *chat* yang bersifat mengancam.
- 4) Meminta Pembeli untuk membatalkan pesanan.

- 5) Mengarahkan Pembeli untuk bertransaksi di luar aplikasi/situs Shopee.
- 6) Melakukan *spam chat*.

Bahan analisis pada makalah ini akan dibatasi pada implementasi *chat filtering* sederhana berupa deteksi kata-kata kasar ataupun mengandung pornografi, permintaan untuk membatalkan pesanan, dan niat mengarahkan pembeli pada transaksi di luar aplikasi dengan adanya *sharing* informasi pribadi (email, sosial media, nomor rekening, nomor telepon, dan website).

II. LANDASAN TEORI

A. Regular Expression (Regex)

Regular Expression atau biasa disebut dengan regex adalah semacam pola yang digunakan untuk mencari dan memanipulasi teks berdasarkan aturan tertentu. Regex memungkinkan identifikasi dan pencocokan pola teks yang kompleks dan beragam, seperti alamat email, nomor telepon, atau pola teks khusus lainnya. Regex terdiri dari karakter normal dan karakter khusus (meta-characters) yang bersama-sama membentuk pola pencarian. Regex mampu melakukan pencocokan pola teks dengan efisiensi yang tinggi untuk menentukan apakah suatu teks diterima dengan aturan yang ada.

Beberapa contoh pendefinisian pola regex adalah sebagai berikut:

1) Karakter

Semua karakter, mulai dari angka 0-9, huruf a-z termasuk varian huruf besar dan huruf kecil, serta substring karakter lainnya.

2) Metakarakter dan Karakter Spesial

Karakter khusus yang memiliki arti khusus di regex terdiri dari `.`, `+`, `*`, `?`, `^`, `$`, `(`, `)`, `[`, `]`, `{`, `}`, `|`, `\`.

- Tanda titik `.` : cocok dengan karakter apapun. Contoh regex `.uku` akan cocok dengan kuku, Duku, buku, dan tuku.
- Tanda `|` : menyatakan OR
- Tanda `^` : menyatakan negasi

3) Kelas karakter

Construct	Deskripsi
<code>[abc]</code>	a, b, atau c (simple class)
<code>[^abc]</code>	Semua karakter selain a,b,c (negasi)
<code>[a-zA-Z]</code>	a sampai z atau A sampai Z, inclusive (range)
<code>[a-d[m-p]]</code>	a sampai d atau m sampai p (gabungan)
<code>[a-z&&[def]]</code>	d, e atau f (irisan)

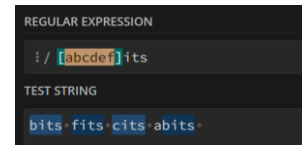
<code>[a-z&&[^bc]]</code>	a sampai z, kecuali b dan c (subtraksi)
<code>[a-z&&[^m-p]]</code>	a sampai z, dan bukan m sampai p (subtraksi)

Tabel 2.3.1 Konstruksi Regex Character Class

Sumber: [Modul Praktikum Regex](#)

a) Simple class

Karakter dalam kurung siku, akan dicari yang cocok dengan substring salah satu karakter dari a,b,c,d,e,f dengan substring lanjutan "its".

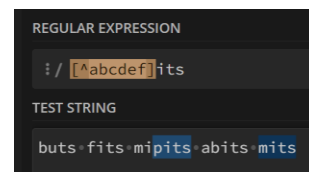


Gambar 2.1.3.1 Contoh Simple Class

Sumber: Dokumentasi Penulis

b) Negasi

Karakter `^` menandakan negasi, substring yang diikuti "its" tidak boleh mengandung substring "a,b,c,d,e,f" pada bagian awalnya.

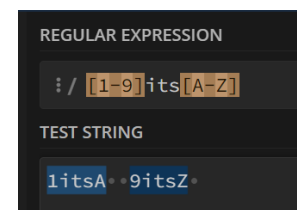


Gambar 2.1.3.2 Contoh Penggunaan Negasi

Sumber: Dokumentasi Penulis

c) Range

Karakter `-` menandakan range karakter pada suatu rentang. Substring "its" harus diawali dengan salah satu karakter angka dari rentang 1-9 dan diakhiri karakter A-Z.

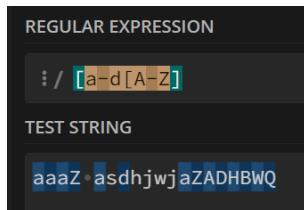


Gambar 2.1.3.3 Contoh Penggunaan Range

Sumber: Dokumentasi Penulis

d) Union

Gabungan dari range, Gambar 2.3.4 menunjukkan ilustrasi union berupa gabungan dari huruf a-d dengan akhiran karakter A-Z



Gambar 2.1.3.4 Contoh Penggunaan Union
Sumber: Dokumentasi Penulis

4) Predefined Character Class

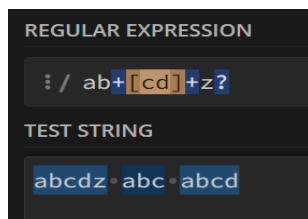
Construct	Deskripsi
.	Semua karakter
\d	Digit [0-9]
\D	Non digit [^0-9]
\s	Whitespace character [\t\n\r\b\f]
\S	Non whitespace character [^\s]
\w	Word character [a-zA-Z_0-9]
\W	Non word character [^\w]

Tabel 2.1.4.1 Predefined Character Class
Sumber: Modul Praktikum Regex

5) Quantifier dan Repetition Operator

Construct	Deskripsi
X?	X muncul satu atau tidak sama sekali
X*	X muncul nol atau banyak
X+	X muncul satu atau banyak
x{n}	X muncul tepat n kali
x{n,}	X muncul setidaknya n kali
x{n,m}	X muncul antara n sampai m kali

Tabel 2.1.5.1 Regex Quantifier dan Operator
Sumber: Modul Praktikum Regex



Gambar 2.1.5.1 Contoh Penggunaan Quantifier dan Operator

Sumber: Dokumentasi Penulis

6) Boundary Matchers

Pencocokan posisi seperti awal kalimat, awal kata, akhir kata.

Construct	Deskripsi
^	Awal baris
\$	Akhir baris
\b	Batas kata
\B	Batas bukan kata
\G	Akhir match sebelumnya
\Z	Akhir dari input tapi untuk final terminator jika ada
\z	Akhir dari input

Tabel 2.1.6.1 Regex Boundary Matchers
Sumber: Modul Praktikum Regex



Gambar 2.1.6.1 Contoh Penggunaan Regex Positioning

Sumber: Dokumentasi Penulis

Dengan penerapan pola karakter regex dan operator, kita bisa membuat sebuah aturan yang akan menentukan apakah suatu string diterima atau tidak.

B. Konsep Singkat String

Sebelum membahas algoritma pengecekan string, perlu dipahami konsep dasar dari sebuah string. Pada string S dengan panjang n kita akan memiliki:

$$S = x_0 x_1 \dots x_{m-1}$$

Maka kita memiliki prefix (substring S[0...k]) dan suffix (substring S[k...m-1]), dengan k indeks di antara 0 dan m-1.

C. Pattern/String Matching

String matching atau pattern matching, adalah proses mencari kemunculan suatu pola (pattern) tertentu dalam sebuah teks (text). Tujuan utama dari string matching adalah untuk menemukan semua posisi dalam teks di mana pola muncul. String matching adalah fundamental dalam berbagai aplikasi komputasi seperti pencarian teks, analisis DNA, analisis citra sidik jari, pemrosesan bahasa alami, dan pencarian data di basis data. Proses ini bisa dilakukan menggunakan berbagai metode, baik yang sederhana seperti brute-force hingga yang lebih kompleks dan efisien seperti algoritma Knuth-Morris-Pratt (KMP) dan Boyer-Moore (BM).

Dalam string matching, akan terdefinisi sebuah teks (T) berupa string panjang yang akan menjadi lokasi pencarian kita dengan panjang n karakter, serta sebuah pattern (P) berupa string dengan panjang m karakter ($m \ll n$) dan akan dicari dalam teks (T) yang ada. Pada string matching kita ditugaskan

untuk menentukan apakah pattern ditemukan dalam teks, dan mencari lokasi di mana ditemukannya pertama kali.

Algoritma string matching dapat diklasifikasikan ke dalam dua kategori utama: algoritma eksak dan algoritma aproksimasi. Algoritma eksak mencari kecocokan sempurna antara pola dan teks (*exact match*), sedangkan algoritma aproksimasi mengizinkan beberapa perbedaan atau kesalahan antara dua string. Kesalahan ini bisa berupa penyisipan, penghapusan, atau substitusi karakter. Salah satu contoh dari algoritma pencocokan aproksimasi adalah Levenshtein Distance. Algoritma ini sangat berguna dalam aplikasi di mana kesalahan atau variasi dalam data sangat umum, seperti dalam pengenalan suara, analisis DNA, pemrosesan bahasa alami, dan pencarian teks yang tidak terstruktur. Dalam algoritma pencocokan aproksimasi akan digunakan dua metrik pengukur:

- 1) Jarak Edit: Jarak minimum antara dua string yang diukur berdasarkan operasi penyisipan, penghapusan, atau substitusi.
- 2) Persentase Kesalahan (Toleransi Kesalahan): Proporsi kesalahan yang diizinkan.

Untuk pencocokan *exact*, algoritma klasik yang sering digunakan adalah algoritma brute-force (kurang efektif), KMP, dan BM.

1) Brute Force (BF)

Brute Force adalah algoritma pencocokan string *exact* paling sederhana yang akan melakukan pengecekan pada setiap huruf pada teks dengan setiap huruf pada pattern. Sifatnya yang mengecek setiap posisi dalam teks menjadikannya memiliki kompleksitas waktu $O(n*m)$, terjadi saat karakter pertama hingga dua dari terakhir sama terus-menerus dengan teks, sedangkan kasus terbaiknya $O(n)$ terjadi ketika karakter pertama pattern selalu berbeda dengan teks. Cara kerjanya sebagai berikut:

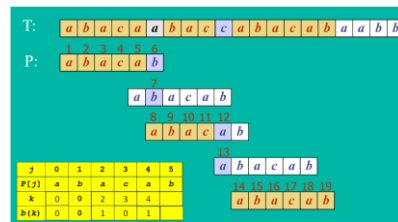
```
Teks: NOBODY NOTICED HIM
Pattern: NOT

NOBODY NOTICED HIM
1 NOT
2 NOT
3 NOT
4 NOT
5 NOT
6 NOT
7 NOT
8 NOT
```

Gambar 2.3.1.1 Pencocokan Brute Force
Sumber: [Website Homepage Rinaldi Munir](#)

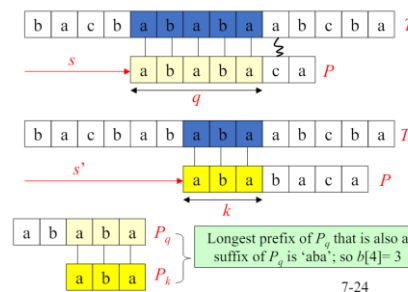
- a) Mulai dari karakter pertama dalam teks.
- b) Bandingkan karakter pertama pola dengan karakter teks saat ini. Jika karakter cocok, lanjutkan ke karakter berikutnya dari pola dan teks. Jika karakter tidak cocok, geser pola satu posisi ke kanan dalam teks dan ulangi proses pencocokan karakter pertama pola.

- c) Jika seluruh pola berhasil dicocokkan dengan bagian dari teks, maka kecocokan ditemukan.
- 2) Knuth Morris Pratt (KMP)



Gambar 2.3.2.1 Pencocokan Knuth Morris Pratt
Sumber: [Website Homepage Rinaldi Munir](#)

Knuth Morris Pratt adalah algoritma pencocokan string *exact* yang menggunakan tabel *longest prefix suffix* (LPS) atau *Border Function* untuk menentukan indeks posisi pengecekan pada pattern berikutnya serta menghindari pemeriksaan ulang pattern. Pengecekannya mirip dengan Brute Force, mengecek pattern dari kiri ke kanan, tetapi dengan batasan tertentu dengan pembuatan LPS. LPS adalah jumlah prefix dan suffix sama terpanjang yang bisa didapatkan pada *missmatch pattern* pada indeks tertentu. Ilustrasi proses pembentukan tabel LPS dengan i indeks *missmatch* pada teks dan j indeks *missmatch* pada pattern, serta k adalah j-1 adalah sebagai berikut:



Gambar 2.3.2.1 Cara Pembentukan Tabel LPS
Sumber: [Website Homepage Rinaldi Munir](#)

Berikut adalah salah satu contoh hasil penentuan border function beserta detail suffix dan prefixnya:

```
KMP E X A M P L E
Hitung border function Suffix P[1..k].
j=0 k=- b[0]=0 Suffix: - Prefix: E
j=1 k=0 b[1]=0 Suffix: - Prefix: E, EA
j=2 k=1 b[2]=0 Suffix: A, AA Prefix: E, EA, EAA
j=3 k=2 b[3]=0 Suffix: A, AA, AAA Prefix: E, EA, EAA, EAAA
j=4 k=3 b[4]=0 Suffix: A, AA, AAA, AAAA Prefix: E, EA, EAA, EAAA, EAAAA
j=5 k=4 b[5]=0 Suffix: A, AA, AAA, AAAA, AAAAA Prefix: E, EA, EAA, EAAA, EAAAA, EAAAAA
j=6 k=5 b[6]=1 Suffix: E, EA, EAA, EAAA, EAAAA, EAAAAE Prefix: E, EA, EAA, EAAA, EAAAA, EAAAAE
```

Gambar 2.3.2.2 Pembentukan Tabel LPS
Sumber: [Website Homepage Rinaldi Munir](#)

Setelah dibuat tabel LPS, pencocokan karakter akan dilakukan sebagai berikut:

- a) Mulai dari karakter pertama dalam teks dan pola.
- b) Bandingkan karakter dari pola dengan teks. Jika karakter cocok, lanjutkan ke karakter berikutnya. Jika karakter tidak cocok dan bukan karakter pertama dari pola, gunakan nilai LPS untuk menentukan posisi

III. PEMBAHASAN

A. Implementasi Kode

Kode-kode di bawah ini dapat diakses melalui link github [berikut](#). Berikut adalah implementasi kode algoritma dasar Brute Force, KMP, dan BM yang menggunakan bahasa Java. Ketiga kelas akan melakukan pengecekan apakah pattern ditemukan pada teks, jika ya akan direturn indeks pertama kali ditemukannya, jika tidak akan mengembalikan angka -1.

Fungsi BFAlgo akan melakukan iterasi for loop i mulai dari indeks pertama teks hingga n-m, dengan pada tiap-tiap iterasinya dilakukan pengecekan pattern. Sebuah variable *counter counstsame* digunakan untuk menghitung jumlah kecocokan selama pengecekan pattern, apabila nilainya sudah sama dengan panjang pattern (m), maka pencarian akan dihentikan.

```
1 public class BFAlgo {
2     public static int BFAlgo(String text, String pattern){
3         int n = text.length();
4         int m = pattern.length();
5         int idxfound = -1;
6         for (int i = 0; i<n-m; i++){
7             int countsame = 0;
8             for (int j = 0; j<m; j++){
9                 if (pattern.charAt(j)!=text.charAt(i+j)){
10                    break;
11                }
12                countsame++;
13            }
14            if(countsame==m){
15                System.out.println("Found it!");
16                idxfound = i;
17                break;
18            }
19        }
20        return idxfound;
21    }
22 }
23 }
```

Gambar 3.1.1 Kode Brute Force
Sumber: Dokumentasi Penulis

Metode KMPAlgo pertama-tama akan memanggil fungsi *lpsValue* untuk mendapatkan array of integer yang berisi border function. Selanjutnya akan dilakukan iterasi for loop i mulai dari indeks 0 hingga n dan melakukan pengecekan pattern setiap kali pergeseran indeks i, dengan setiap *start index* pada pattern ditentukan oleh tabel LPS. Apabila pada suatu saat ditemukan nilai indeks iterasi j sama dengan m-1, atau selesai mengecek semua karakter pada pattern, maka pattern sudah ditemukan dan return indeks i.

```
1 public class KMPAlgo {
2     // function to build LPS / border func table in KMP
3     // will return an array of border func
4     public static int[] lpsValue(String pattern) {
5         System.out.println("Let's build the border function!\n");
6         int[] borderfunc = new int[pattern.length()];
7         int prefixbestlength = 0;
8         int idx = 1; // the lps array index starts from 1 since borderfunc[0] is always 0
9         borderfunc[0] = 0; // when j = 0 we always get border func 0
10
11         while (idx < pattern.length()) {
12             if (pattern.charAt(idx) == pattern.charAt(prefixbestlength)) {
13                 prefixbestlength++;
14                 borderfunc[idx] = prefixbestlength;
15                 idx++;
16             } else {
17                 if (prefixbestlength > 0) { // based on matching prefix
18                     prefixbestlength = borderfunc[prefixbestlength - 1];
19                 } else { // no match
20                     borderfunc[idx] = 0;
21                     idx++;
22                 }
23             }
24         }
25         return borderfunc;
26     }
27
28     // main KMP algo that will return the idx found on text
29     public static int KMPAlgo(String pattern, String text) {
30         int n = text.length();
31         int m = pattern.length();
32         int[] borderfunc = lpsValue(pattern);
33
34         // initialize starting mismatch index
35         int idx_i = 0;
36         int idx_j = 0;
37
38         while (idx_i < n) {
39             if (pattern.charAt(idx_j) == text.charAt(idx_i)) {
40                 if (idx_j == m - 1) {
41                     System.out.println("Found it!");
42                     return idx_i - m + 1;
43                 }
44                 idx_i++;
45                 idx_j++;
46             } else {
47                 if (idx_j > 0) {
48                     idx_j = borderfunc[idx_j - 1];
49                 } else {
50                     idx_i++;
51                 }
52             }
53         }
54         System.out.println("We can't find the pattern");
55         return -1;
56     }
57 }
```

Gambar 3.1.2 Kode Knuth Morris Pratt
Sumber: Dokumentasi Penulis

Metode BMAIgo pertama-tama akan memanggil fungsi *lastoccur* untuk mendapatkan array of integer yang berisi last occurrence untuk setiap karakter yang ada pada ASCII. Selanjutnya akan dilakukan iterasi for loop i mulai dari indeks 0 hingga n-1 dan melakukan pengecekan pattern, apabila terjadi *missmatch* maka akan digunakan nilai *lastocurrence* untuk menentukan indeks i pencarian selanjutnya. Apabila pada suatu saat ditemukan nilai indeks iterasi j adalah 0, atau selesai mengecek semua karakter pada pattern, maka pattern sudah ditemukan dan return indeks i.

```

1 public class BMAlgoritma {
2     public static int[] lastoccur(char[] pattern){
3         System.out.println("Let's build the last occurrence list!\n");
4         int lastoccur[] = new int[128];
5         for (int i = 0; i<128; i++){
6             lastoccur[i] = -1; //every char doesnt exist at first
7         }
8         for (int i = 0; i<pattern.length(); i++){
9             lastoccur[pattern.charAt(i)] = i;
10        }
11        return lastoccur;
12    }
13
14    // main algo of BM
15    public static int BMAlgo(string pattern, string text){
16        int n = text.length();
17        int m = pattern.length();
18        int lastoccur[] = lastoccur(pattern);
19        int idx_i = m-1;
20        int idx_j = m-1;
21
22        while (idx_i<=n-1) {
23            if(pattern.charAt(idx_i)==text.charAt(idx_i)){
24                if(idx_j==0){
25                    System.out.println("Found it!");
26                    return idx_i;
27                }
28            } else{
29                idx_i--;
30                idx_j--;
31            }
32        } else{
33            int lo = lastoccur[text.charAt(idx_i)];
34            idx_i = idx_i-m+Math.min(idx_j, 1+lo);
35            idx_j = m-1;
36        }
37    }
38    System.out.println("We can't find the pattern");
39    return -1;
40 }
41 }
42 }

```

Gambar 3.1.3 Kode Boyer Moore
Sumber: Dokumentasi Penulis

B. Filter Kata-Kata Kasar

Dalam penanganan filter kata-kata kasar ini akan digunakan sebuah kamus yang memuat kata-kata kasar yang kerap digunakan masyarakat Indonesia. Kata-kata pada kamus tersebut akan menjadi pattern, dan kita akan melakukan pencarian dengan KMP dan BM pada teks. Algoritma yang terdefinisi ada pada parameter algo.

```

1 import java.io.IOException;
2 import java.io.File;
3 import java.io.File;
4 import java.io.File;
5 import java.util.List;
6 public class FilterBadwords {
7     public static void filterbadwords(string text, int algo) throws IOException{
8         System.out.println("CHECKING BADWORDS...");
9         // Membaca dictionary dari file badwords
10        List<String> dictionary = Files.readAllLines(Paths.get("data/badwords.txt"));
11        if (algo == 0){
12            // Use KMP
13            for (String word : dictionary) {
14                if (BMAlgoritma.KMPAlgo(word, text)!=-1) {
15                    System.out.println("This word is not allowed: " + word);
16                    break;
17                }
18            }
19        } else if (algo == 1){
20            // Use BM
21            for (String word : dictionary) {
22                if (BMAlgoritma.BMAlgo(word, text)!=-1) {
23                    System.out.println("This word is not allowed: " + word);
24                    break;
25                }
26            }
27        } else{
28            // Use BF
29            for (String word : dictionary) {
30                if (BMAlgoritma.BFAlgo(word, text)!=-1) {
31                    System.out.println(">>> This word is not allowed: " + word);
32                    break;
33                }
34            }
35        }
36    }
37 }
38 }
39 }
40 }

```

Gambar 3.2.1 Kode Filter Kata-Kata Tidak Senonoh
Sumber: Dokumentasi Penulis

C. Filter Pengajuan Pembatalan

Filter pengajuan pembatalan dilakukan dengan cara yang sama seperti filter kata-kata kasar. Pertama kali akan dilakukan pembacaan file kamus terdefinisi, kemudian akan dilakukan iterasi pada teks untuk mencari kemunculan pattern dari kamus berdasarkan algoritma yang terpilih. Algoritma yang digunakan akan terdefinisi pada parameter algo.

```

1 import java.io.IOException;
2 import java.io.File;
3 import java.io.File;
4 import java.io.File;
5 import java.util.List;
6 public class Cancel {
7     public static void filtercancel(string text, int algo) throws IOException{
8         // Membaca dictionary dari file cancel
9         System.out.println("CHECKING CANCEL...");
10        List<String> dictionary = Files.readAllLines(Paths.get("data/cancel.txt"));
11        if (algo == 0){
12            // Use KMP
13            for (String word : dictionary) {
14                if (BMAlgoritma.KMPAlgo(word, text)!=-1) {
15                    System.out.println(">>> This word is not allowed: " + word);
16                    break;
17                }
18            }
19        } else if (algo == 1){
20            // Use BM
21            for (String word : dictionary) {
22                if (BMAlgoritma.BMAlgo(word, text)!=-1) {
23                    System.out.println(">>> This word is not allowed: " + word);
24                    break;
25                }
26            }
27        } else{
28            // Use BF
29            for (String word : dictionary) {
30                if (BMAlgoritma.BFAlgo(word, text)!=-1) {
31                    System.out.println(">>> This word is not allowed: " + word);
32                    break;
33                }
34            }
35        }
36    }
37 }
38 }
39 }

```

Gambar 3.1.3 Kode Filter Cancel
Sumber: Dokumentasi Penulis

D. Translate

Sudah tidak asing adanya orang-orang yang menggunakan angka sebagai pengganti huruf, seperti huruf “e” yang digantikan oleh angka 3, huruf “s” digantikan oleh angka 5, dan seterusnya. Maka dari itu, sebelum dilakukan pemrosesan *pattern match*, akan dilakukan terjemahan ke bentuk aslinya.

```

1 import java.util.HashMap;
2 import java.util.HashMap;
3 import java.util.HashMap;
4 public class Translate {
5
6     public static String translate(String input) {
7         numPattern = "[1-9]";
8         numPattern = numPattern.compile(numPattern);
9         matcher = pattern.matcher(input);
10        StringBuffer result = new StringBuffer();
11
12        while (matcher.find()) {
13            replacement;
14            switch (matcher.group()) {
15                case "1":
16                    replacement = "l"; break;
17                case "4":
18                    replacement = "a"; break;
19                case "6":
20                    replacement = "g"; break;
21                case "5":
22                    replacement = "s"; break;
23                case "3":
24                    replacement = "e"; break;
25                case "0":
26                    replacement = "o"; break;
27                default:
28                    replacement = matcher.group(); break;
29            }
30            matcher.appendReplacement(result, replacement);
31            matcher.appendTail(result);
32        }
33        String translated = result.toString().toLowerCase();
34        return translated.toString();
35    }
36 }
37 }
38 }
39 }
40 }

```

Gambar 3.4.1 Kode Translate
Sumber: Dokumentasi Penulis

E. Filter Sharing Informasi Pribadi

Dalam filtering informasi pribadi, akan dilakukan empat kali pengecekan terpisah untuk nomor telepon, alamat email, username sosial media, dan tautan. Setiap pengecekan akan dipisah supaya bisa diketahui penyebab ditolaknya suatu teks. Fungsi filterPrivateInfo akan menyatukan seluruh fungsi pengecekan dan mengeluarkan output yang sesuai dengan hasil pengecekan. Pada pengecekan email dan username sosial media ditambahkan pada satu blok *if conditional* yang sama karena email pasti mengandung karakter “@”. Berikut adalah detail peraturan regex yang digunakan:

```

1 import java.util.regex.*;
2
3 public class FilterPrivateInfo {
4     public static void FilterPrivateInfo(String text) {
5         System.out.println("CHECKING PRIVATE INFO...");
6         boolean containsPhoneNumber = containsPhoneNumber(text);
7         boolean containsEmailAddress = containsEmailAddress(text);
8         boolean containsSocialMedia = containsSocialMedia(text);
9         boolean containsOutsideLink = containsOutsideLink(text);
10        boolean safe = true;
11
12        if (containsPhoneNumber) {
13            safe = false;
14        }
15        if (containsEmailAddress) {
16            safe = false;
17        }
18        else if (containsSocialMedia) {
19            safe = false;
20        }
21        if (containsOutsideLink) {
22            safe = false;
23        }
24        if (safe) {
25            System.out.println("The text does not contain prohibited information: " + text);
26        }
27    }
28
29    // Regex pattern for matching phone numbers
30    public static boolean containsPhoneNumber(String text) {
31        String phonePattern = "(\\+62\\s?\\d{1,3})([-\\s/\\d{3,4})?)?";
32        Pattern compiledPattern = Pattern.compile(phonePattern);
33        Matcher matcher = compiledPattern.matcher(text);
34        if (matcher.find()) {
35            System.out.println("The text contains a phone number: " + matcher.group());
36            return true;
37        }
38        return false;
39    }
40
41
42    // Regex pattern for matching email addresses
43    public static boolean containsEmailAddress(String text) {
44        String emailPattern = "\\b[A-Z0-9_+]+@[A-Z0-9-]+\\.([A-Z]{2,})\\b";
45        Pattern compiledPattern = Pattern.compile(emailPattern, Pattern.CASE_INSENSITIVE);
46        Matcher matcher = compiledPattern.matcher(text);
47        if (matcher.find()) {
48            System.out.println("The text contains an email address: " + matcher.group());
49            return true;
50        }
51        return false;
52    }
53
54
55    // Regex pattern for matching social media handles
56    public static boolean containsSocialMedia(String text) {
57        String socmedPattern = "@([A-Za-z0-9_]{3,25})";
58        Pattern compiledPattern = Pattern.compile(socmedPattern);
59        Matcher matcher = compiledPattern.matcher(text);
60        if (matcher.find()) {
61            System.out.println("The text contains a social media handle: " + matcher.group());
62            return true;
63        }
64        return false;
65    }
66
67
68    // Regex pattern for matching links
69    public static boolean containsOutsideLink(String text) {
70        String linkPattern = "\\b(https?|ftp)://[!-a-zA-Z0-9+&@#/?=|:;,~]*[-a-zA-Z0-9+&@#/%*~]";
71        Pattern compiledPattern = Pattern.compile(linkPattern, Pattern.CASE_INSENSITIVE);
72        Matcher matcher = compiledPattern.matcher(text);
73        if (matcher.find()) {
74            System.out.println("The text contains an outside link: " + matcher.group());
75            return true;
76        }
77        return false;
78    }
79
80
81 }

```

Gambar 3.5.1 Kode Filter Informasi Pribadi
Sumber: Dokumentasi Penulis

1) Untuk pengecekan pola nomor telepon ada pada deklarasi variable phonePattern: `(\\+62\\s?\\d{1,3})([-\\s/\\d{3,4})?)?/(\\(\\d+\\|\\s?\\d+\\)/\\d{3}(\\|`

`s\\d+)+\\|\\d+([-\\s/\\d+)+\\|\\d{8,12}`. Pola ini juga akan menangani pemisahan dengan spasi dan “-“ untuk mencegah usaha *bypassing*. Regex ini akan dibatasi untuk nomor telepon Indonesia saja.

Potongan Regex	Deskripsi
<code>\\+62\\s?\\d{1,3}</code>	Menangkap nomor telepon yang dimulai dengan kode negara Indonesia (+62), diikuti oleh 1 hingga 3 digit. <code>\\s?</code> berarti ada atau tidak ada spasi.
<code>([-\\s/\\d{3,4})?([-\\s/\\d{4,})?)?</code>	Mengizinkan format dengan tanda hubung atau spasi, misalnya -3456 atau -3456-7890.
<code>(\\(\\d+\\ \\s?\\d+\\)</code>	Menangkap nomor telepon dengan format kode area dalam kurung, misalnya (021) 123456.
<code>\\d{3}(\\ \\s\\d+)+</code>	Menangkap format nomor telepon dengan grup tiga digit diikuti oleh spasi dan digit lainnya, misalnya 123 4567 8901.
<code>\\d+([-\\s/\\d+)+</code>	Menangkap nomor telepon dengan tanda hubung atau spasi, misalnya 123-4567-8901.
<code>\\d{8,12}</code>	Menangkap nomor telepon dengan panjang minimal 8 digit karena nomor telepon Indonesia dimulai dari 8 digit dan maksimal 12 digit, misalnya 12345678 dan 467321875

Tabel 3.5.1 Penjelasan Regex Nomor Telepon

2) Untuk pengecekan pola email ada pada deklarasi variabel emailPattern: `\\b[A-Z0-9_+]+@[A-Z0-9-]+\\.([A-Z]{2,})\\b`

Potongan Regex	Deskripsi
<code>[A-Z0-9_+]+</code>	Menangkap karakter yang valid dalam bagian username email (huruf, angka, titik, underscore, persen, plus, dan tanda minus).
<code>@([A-Z0-9-]+)</code>	Menangkap karakter yang valid setelah simbol @ (huruf, angka, titik, dan tanda minus).
<code>\\.([A-Z]{2,})</code>	Menangkap domain level atas yang dimulai dengan titik diikuti oleh minimal dua huruf, misalnya .com dan .id

Tabel 3.5.2 Penjelasan Regex Alamat Email

3) Untuk pengecekan pola username sosial media ada pada deklarasi variabel socmedPattern: `@([A-Za-z0-9_]{3,25})`

Potongan Regex	Deskripsi
<code>@</code>	Menangkap simbol @ yang umumnya digunakan dalam username media sosial, misalnya @regexinstaofficial
<code>[A-Za-z0-9_]{3,25}</code>	Menangkap karakter yang valid dalam handle media sosial (huruf, angka, titik, dan underscore) dengan panjang antara 3 hingga 25 karakter.

Tabel 3.5.3 Penjelasan Regex Username Sosial Media

4) Untuk pengecekan pola link ada pada deklarasi

variabel linkPattern: `\\b(https?|ftp):\\/[-a-zA-Z0-9+&@#/%?~_!:,;]*[-a-zA-Z0-9+&@#/%~_!:]\\b`
 Pola ini juga akan menangani pemisahan dengan spasi dan karakter lain untuk mencegah usaha *bypassing*.

Potongan Regex	Deskripsi
<code>(https? ftp):\\/</code>	Menangkap protokol yang umum digunakan (http, https, ftp).
<code>[-a-zA-Z0-9+&@#/%?~_!:,;]*</code>	Menangkap karakter yang valid dalam URL (huruf, angka, dan berbagai simbol yang umum digunakan dalam URL).
<code>[-a-zA-Z0-9+&@#/%~_!]</code>	Menangkap karakter yang valid di akhir URL, misalnya path/to/resources

Tabel 3.5.4 Penjelasan Regex Tautan

F. Perbandingan dan Analisis Algoritma

```

Welcome to Chat Filter!
Here we are going to determine if your chat violates our rule:
- no vulgar words
- intention to ask customer to cancel an order
- sharing private information that could lead to dangerous out of app transaction

Enter a text message to check:
Tulisan di deskripsi tuh dibaca! Batalin nih pesanan! Sudah jelas-jelas ditulis kalau minimal beli itu 10 pcs kok malah Cuma beli 1 biji. Ya kita rugi dong! Rada Tolol emang ini

Choose your algorithm (KMP/BM/BF):
>> KMP

CHECKING BADWORDS...
Found it!
This word is not allowed: tolol

CHECKING CANCEL...
Found it!
>>> This word is not allowed: batal

CHECKING PRIVATE INFO...
The text does not contain prohibited information: Tulisan di deskripsi tuh dibaca! Batalin nih pesanan! Sudah jelas-jelas ditulis kalau minimal beli itu 10 pcs kok malah Cuma beli 1 biji. Ya kita rugi dong! Rada Tolol emang ini
KMP Execution Time: 28 ms
    
```

Gambar 3.6.1 Contoh Hasil Kompilasi Program
 Sumber: Dokumentasi Penulis

```

Welcome to Chat Filter!
Here we are going to determine if your chat violates our rule:
- no vulgar words
- intention to ask customer to cancel an order
- sharing private information that could lead to dangerous out of app transaction

Enter a text message to check:
Kalau mau harga grosir bs cek di link ini ya kak, https://www.tokopedia.com/xhaetaymisibri. Buat selalu update info bisa jg nih kak cek Ig kita di @xhaetayshop

Choose your algorithm (KMP/BM/BF):
>> BF

CHECKING BADWORDS...

CHECKING CANCEL...

CHECKING PRIVATE INFO...
The text contains a social media handle: @xhaetayshop
The text contains an outside link: https://www.tokopedia.com/xhaetaymisibri
BF Execution Time: 23 ms

Do you want to try again? (Y/N)
    
```

Gambar 3.6.2 Contoh Hasil Kompilasi Program 2
 Sumber: Dokumentasi Penulis

Pada program filter chat ini akan dilakukan pengujian pada algoritma KMP dan BM. Pengujian akan dilakukan dengan inputan yang sama.

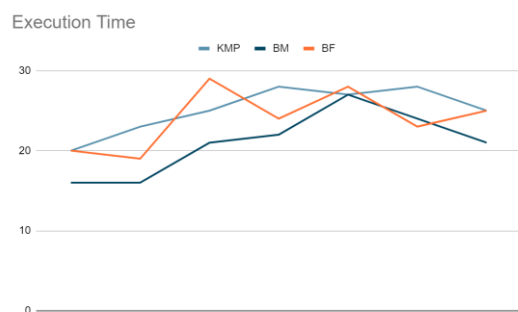
Test Case	Teks Input
1.	Hi welcome to my humble shop! For far more cheap prices you can contact me thru WAs 084 3984 3287 or at email@mail.com thank you!
2.	Kalau emang gasuka yaudah batalin aja ngapa, sulit banget dih t0l0l.
3.	Kalau mau harga grosir bs cek di link ini ya kak, di https://www.tokopedia.com/xhaetaymisibri. Buat selalu update

	info bisa jg nih kak cek Ig kita di @xhaetayshop. Di sana banyak diskon juga kok buat pembelian per pcs jadi ajak teman kakak ya! Ajak 1 teman bisa dapat promo buy 2 get 1 di pembelian berikutnya dengan menyertakan foto bukti struk pembelian. Yang pesanan ini boleh dibatalin dulu aja gapapa kok.
4.	Maaf Kak ini kok barangnya sepertinya ada yang rusak ya. Waktu ak buka kok ada robek di bagian lengan bajunya, nih ak ada video unboxing buat buktinya. Aku ajukan claim ya kak? Ongkir yang tanggung saya saja tidak apa -- Heh cok jangan ngarang ya, barang kita ini udh dicek sbmlm ngirim!
5.	Pagi kak, Oya kak maaf sekali boleh tolong cancel dl orderannya tidak ya? Stok kami lagi habis, boleh nnt chat aja ke (+62) 857 999 000 biar dikasih harga spesial buat kakak, pesan lewat situ saja
6.	Tulisan di deskripsi tuh dibaca! Batalin nih pesanan! Sudah jelas-jelas ditulis kalau minimal beli itu 10 pcs kok malah Cuma beli 1 biji. Ya kita rugi dong! Rada Tolol emang ini
7.	Kakak boleh nih cek toko kita yang lain lagi di https://www.tokopedia.com/getmorebeauty kami baru launching toko di sana tuh jadi banyak promonya!

Tabel 3.6.1 Test Case dan Input

Test Case	Panjang Kata	Lama Eksekusi (ms)		
		KMP	BM	BF
1.	129	23	16	19
2.	69	20	16	20
3.	410	25	21	25
4.	287	28	24	23
5.	196	27	27	28
6.	177	28	22	25
7.	146	25	21	29
Rata-Rata	202	25.14	21	24.14

Tabel 3.6.2 Waktu Eksekusi



Gambar 3.6.3 Grafik Hasil Kompilasi Program
 Sumber: Dokumentasi Penulis

Berdasarkan hasil pengujian waktu eksekusi algoritma, didapatkan bahwa algoritma BM menjadi yang paling optimal di antara ketiganya, diikuti dengan algoritma BF di posisi kedua dan KMP di posisi ketiga. Sifat algoritma BM yang unggul dalam pencocokan string dengan banyak jenis karakter membuatnya mudah dalam bergeser indeks, berbeda dengan

algoritma KMP yang melakukan pengecekan mirip BF, tetapi dengan bantuan tabel LPS. Banyaknya macam karakter pada string akan membuat KMP lebih sering *missmatch* di indeks awal pattern, padahal KMP akan jauh lebih efisien apabila indeks *missmatch* ada pada bagian akhir pattern. *Missmatch* yang banyak terjadi di awal menyebabkan KMP bekerja seperti BF dengan bergeser indeks teks satu per satu dan mengecek pattern mulai dari indeks ke-0. Ditambah dengan perlunya *preprocessing* membangun tabel LPS membuat *overhead* pada algoritma KMP sehingga menjadikannya lebih lambat daripada BF, terlebih pada panjang string yang tidak panjang. Pada algoritma BM sebenarnya juga terdapat *overhead* karena harus *preprocessing* tabel *last occurrence*, tetapi karena keunggulannya dalam melakukan loncatan jauh membuatnya mengungguli algoritma BF pada kebanyakan kasus. Maka dari itu, BF mungkin jadi lebih tepat untuk kasus pencarian string pendek karena tidak memerlukan *preprocessing*.

IV. KESIMPULAN

Dalam pengecekan pola informasi pribadi seperti email, nomor telepon, links, dan username sosial media dibutuhkan Regex karena sifatnya yang fleksibel dan proses deteksinya tergantung dengan aturan. Berbeda dengan pengecekan kata-kata kasar dan pengajuan pengembalian, kita gunakan algoritma KMP dan BM karena sifatnya yang mencocokkan secara eksak dengan kata-kata yang sudah terdefinisi dalam sebuah kamus. Pada pencocokan string, algoritma BM jauh lebih efektif daripada KMP karena sifat teks yang memiliki banyak macam huruf. *Overhead* pada pembangunan tabel LPS pada KMP dan *last occurrence* pada BM dapat menyebabkan kinerja yang lebih lambat daripada BF. Maka dari itu, terdapat beberapa kasus di mana BF menjadi lebih cepat daripada KMP dan BM, terutama ketika panjang kata yang pendek.

UCAPAN TERIMA KASIH

Puji syukur saya panjatkan kepada Tuhan Yang Maha Esa atas penyertaan dan berkat yang telah diberikan selama ini sehingga makalah yang berjudul "Penerapan Regex dan String Matching untuk Filtering Informasi Sensitif dalam Pesan Chat" ini dapat diselesaikan dengan baik. Saya juga ingin mengucapkan terima kasih pada pihak-pihak yang telah membantu dalam penyusunan makalah ini:

1. Ibu Dr. Nur Ulfa Maulidevi, selaku dosen pengajar mata kuliah IF2211 Strategi Algoritma K02, atas bimbingan dan pengajaran yang sangat membantu saya dalam memahami materi,
2. Orang tua saya yang selalu mendukung dan membantu saya selama ini,
3. Teman-teman dan sahabat-sahabat yang memberikan saya semangat,

Penulis-penulis jurnal, artikel, dan laman website yang tidak bisa saya sebutkan satu per satu atas sumber informasi yang sangat bermanfaat sebagai bahan referensi makalah ini.

TAUTAN KODE

<https://github.com/pandaandsushi/ChatFilter>

REFERENSI

- [1] Delta. (2022, Nov 26). @username Regular Expression for social media with JavaScript. <https://stackoverflow.com/questions/74579624/username-regular-expression-for-social-media-with-javascript> (diakses pada 12 Juni 2024)
- [2] Drizki. (2023, Juli 10). Indonesian Badwords. <https://github.com/drizki/indonesian-badwords/blob/main/src/dict.json> (diakses pada 12 Juni 2024)
- [3] Makragic, A. (2017, Juni 21). Regex for Indonesian Number. <https://stackoverflow.com/questions/44670612/regex-for-indonesian-phone-number> (diakses pada 12 Juni 2024)
- [4] Moore, A. (2013, Sep 29). Regular Expression to Match URLs in Java. <https://stackoverflow.com/questions/163360/regular-expression-to-match-urls-in-java> (diakses pada 12 Juni 2024)
- [5] Munir, R. (2024) Homepage Rinaldi Munir. Sekolah Teknik Elektro dan Informatika (STEI) ITB. *String Matching dengan Regex*. Retrieved from <https://informatika.stei.itb.ac.id/~rinaldi.munir/Stmik/2022-2023/String-Matching-dengan-Regex-2019.pdf> (diakses pada 11 Juni 2024)
- [6] Munir, R. (2023) Homepage Rinaldi Munir. Sekolah Teknik Elektro dan Informatika (STEI) ITB. *Pencocokan String*. Retrieved from <https://informatika.stei.itb.ac.id/~rinaldi.munir/Stmik/2020-2021/Pencocokan-string-2021.pdf> (diakses pada 11 Juni 2024)
- [7] NTU. (2018, November). *Regular Expressions (Regex)*. Retrieved from <https://www3.ntu.edu.sg/home/ehchua/programming/howto/Regexe.html> (diakses pada 11 Juni 2024)
- [8] Rock, D. (2014, Feb 4). RegEx for Validationg Emails. <https://stackoverflow.com/questions/21608294/regex-for-validating-emails> (diakses pada 12 Juni 2024)
- [9] Seller Shopee Center (2024, April 30). Pelanggaran Chat. Retrieved from <https://seller.shopee.co.id/edu/article/11462/Pelanggaran-Chat> (diakses pada 10 Juni 2024)
- [10] Wibisono, Y. et al. (2020, April 11). "Modul Praktikum: Pengantar Regular Expression". Retrieved from <https://docs.google.com/document/d/1Is6h1A6m-Zhzw6e5erwMNUAG0D1iwL-eVmVMS2XQoc/edit> (diakses pada 10 Juni 2024)

PERNYATAAN

Dengan ini saya menyatakan bahwa makalah yang saya tulis ini adalah tulisan saya sendiri, bukan saduran, atau terjemahan dari makalah orang lain, dan bukan plagiasi.

Bandung, 12 Juni 2024



Thea Josephine Halim
(13522012)