

Penggunaan Algoritma Commentz-Walter dalam String Matching

Vicko Novianto | 13513092¹

Program Studi Teknik Informatika

Sekolah Teknik Elektro dan Informatika

Institut Teknologi Bandung, Jl. Ganesha 10 Bandung 40132, Indonesia

¹13513092@std.stei.itb.ac.id

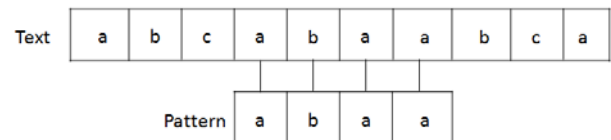
Abstract—Algoritma pencocokan string adalah algoritma untuk menemukan pola di dalam teks. Banyak algoritma pencocokan string yang terkenal, salah satunya adalah algoritma Commentz-Walter. Algoritma ini menggabungkan ide Boyer-Moore dan Aho-Corasick. Algoritma ini memiliki kompleksitas linear pada umumnya. Algoritma ini menggunakan struktur trie untuk menyimpan kata-kata kunci.

Index Terms—Commentz-Walter, Pola, Shift, Teks, Trie

I. PENDAHULUAN

Algoritma pencocokan string adalah algoritma untuk menemukan satu atau beberapa string/pola di dalam string yang lebih besar atau teks. [1] Pencocokan string mempunyai peran yang besar dalam menyelesaikan persoalan di berbagai bidang, misalnya deteksi intrusi dalam jaringan, aplikasi bioinformatik, deteksi plagiarisme, keamanan informasi, pengenalan pola, pencocokan dokumen dan penambangan data (*text mining*). Algoritma pencocokan string bisa dibagi dua, yaitu algoritma pencocokan string eksak (*exact string matching algorithm*) dan algoritma pencocokan string non-eksak (*approximate string matching algorithm*). Algoritma pencocokan string ada yang hanya bisa mencocokkan satu pola dan ada yang bisa mencocokkan banyak pola.

Misalkan Σ adalah set alfabet yang terbatas. Secara formal, baik pola maupun teks adalah vektor dari elemen Σ . Asumsikan teks adalah array $T[1..n]$ dengan panjang n dan pola adalah array $P[1..m]$ dengan panjang m dan $m \leq n$. Pola P ditemukan pada pergeseran (*shift*) s di teks T (dengan kata lain P ditemukan mulai pada posisi $s+1$ di teks T) jika $0 \leq s \leq n-m$ dan $T[s+1..s+m] = P[1..m]$. Jika P ditemukan pada s yang berada pada range $0 \leq s \leq n-m$ maka dikatakan valid. Algoritma pencocokan string adalah masalah menemukan semua s yang valid dengan pola P di dalam teks T . [1]



Gambar 1. Ilustrasi pencocokan string

II. DASAR TEORI

Algoritma pencocokan string yang menggunakan banyak pola menggunakan teks $T = t_1 t_2 \dots t_n$ dimana t_i adalah karakter ke- i dan kita ingin mencari bersamaan untuk sebuah himpunan string $P = \{p_1, p_2, \dots, p_r\}$ dengan p_i adalah string dengan panjang m_i dengan $i=1..r$ [2]. Ada banyak algoritma pencocokan string yang menggunakan banyak pola yang bervariasi dalam kecepatannya dihitung dari kompleksitasnya. Beberapa contohnya :

1. Algoritma pencocokan string Aho-Corasick
Kompleksitas dari algoritma ini adalah linear terhadap panjang pola ditambah panjang dari teks yang dicari ditambah jumlah kecocokan. Algoritma ini membangun sebuah finite state machine yang merupai trie dengan sisi tambahan antara simpul internalnya. Sisi tambahan ini memungkinkan transisi yang cepat bila terjadi ketidakcocokkan (misalkan saat mencari cat dalam trie yang tidak mengandung cat namun mengandung cart, maka akan gagal pada simpul ca) ke cabang trie yang lain yang memiliki prefiks yang sama. Ini memungkinkan kita tidak perlu melakukan runut-balik.
2. Algoritma pencocokan string Rabin Karp
Algoritma ini menggunakan hash untuk menemukan himpunan pola dalam sebuah teks. Kompleksitas terburuknya adalah $O(nm)$.
3. Algoritma pencocokan string Commentz-Walter

III. ALGORITMA COMMENTZ-WALTER

A. Struktur data

Untuk menggambarkan beberapa kata kunci dalam cara yang baik, struktur data yang digunakan adalah *trie*.

The search phase of algorithm B in detail:

Initial phase:

```
v ← root r      (v is the "present" node of T)
i ← wmin        (i points to the document letter
                 above the nodes of depth 1 .)
j ← 0          (j indicates the depth of the
                present node v.)
```

While $i \leq \text{length document}$ do

Scan phase:

begin

```
while there is some son v' of v labeled by  $d_{i-j}$  do
begin
```

```
  v ← v'
```

```
  j ← j + 1
```

```
  output: (w, i) for each w of out(v)
```

```
end
```

shift phase:

begin

```
  i ← i + S(v,  $d_{i-j}$ )
```

```
  j ← 0
```

end end

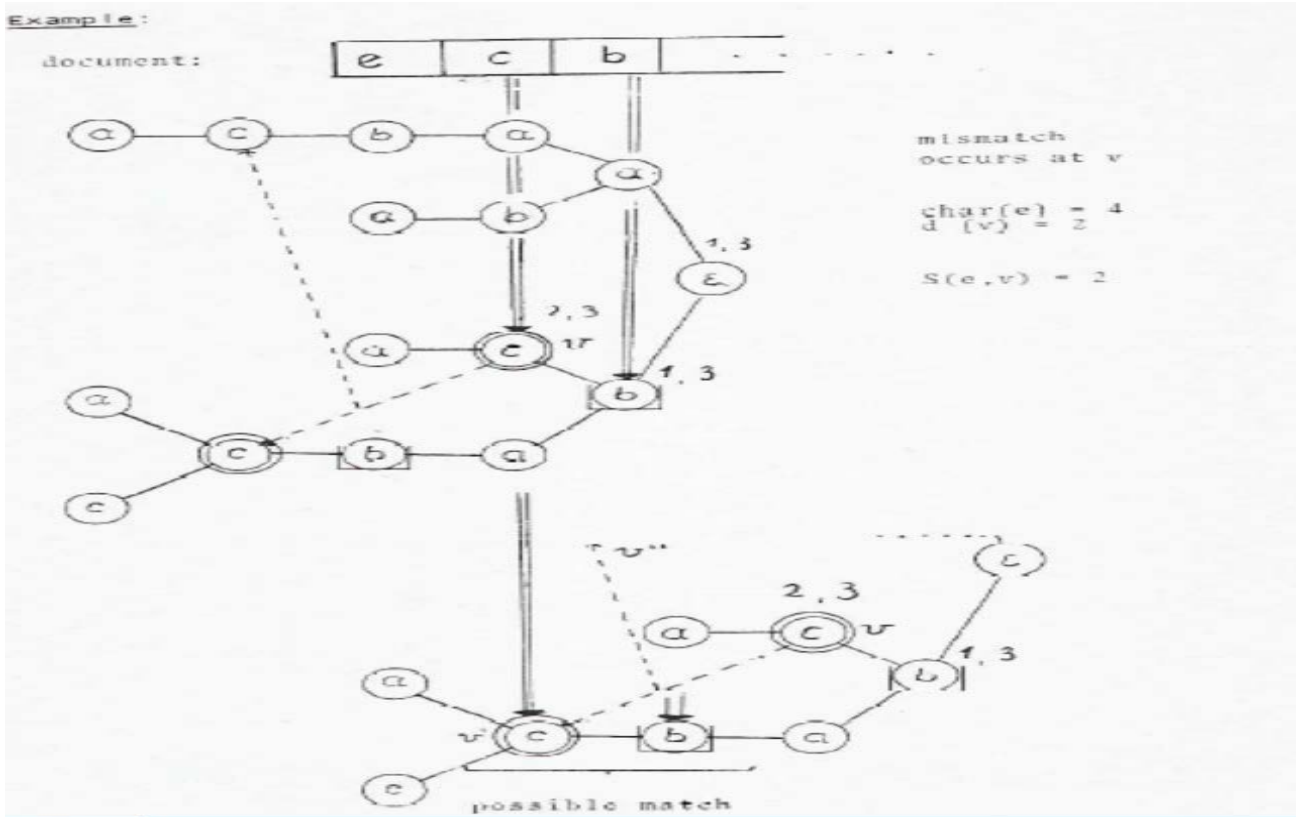
where $S(v, d_{i-j})$ is the length of the shift defined by

```
 $S(v, d_{i-j}) = \min(\max(\text{shift1}(v), \text{char}(d_{i-j}) - j - 1),$   
 $\text{shift2}(v)).$ 
```

Gambar 3. Algoritma fase pencarian

Keluaran dari fase pencarian algoritma pencocokan string Commentz-Walker adalah list dari pasangan (w, i) dengan w adalah sebuah kata dan i adalah sebuah bilangan bulat yang menggambarkan kemunculan w , misalnya (W, i) elemen dari keluaran algoritma ini dan W adalah kata kunci dari K dan $d_{i-|W|+1}, \dots, d_i = W$.

Algoritma ini menggabungkan ide dari algoritma pencocokan string Aho-Corasick dan Boyer-Moore. Kita mendasarkan trie dengan kata kunci yang dibalik. Misalkan w_{\min} melambangkan panjang minimal dari beberapa kata kunci. Algoritma ini mulai menaruh akar r pada T dibawah $d_{w_{\min}+1}$. Selanjutnya kita akan memindai dokumen dari kanan ke kiri sampai terjadi ketidakcocokkan (*mismatch*). Asumsikan kita baru saja memindai kata-kata dokumen yang cocok d_{i-m+1}, \dots, d_i dan sebuah ketidakcocokkan terjadi pada huruf d_{i-m} lalu kita geser akar dari trie ke kanan sebesar beberapa angka dari huruf S yang dihitung dari huruf-huruf dalam dokumen d_{i-m}, \dots, d_i .



Gambar 4. Ilustrasi fase pencarian

Tentunya, setiap pasangan (w, i) yang ditemukan menggambarkan beberapa kemunculan dari kata kunci W . Jadi algoritma ini menemukan setiap kemunculan dari beberapa kata kunci dalam dokumen D .

C. Fase preprocessing

Masukan dari fase pra-pemrosesan adalah himpunan kata kunci $K = \{W_1, \dots, W_n\}$. Keluarannya adalah trie T yang mengandung kata kunci yang dibalik dan fungsi out , shift_1 , shift_2 , dan char .

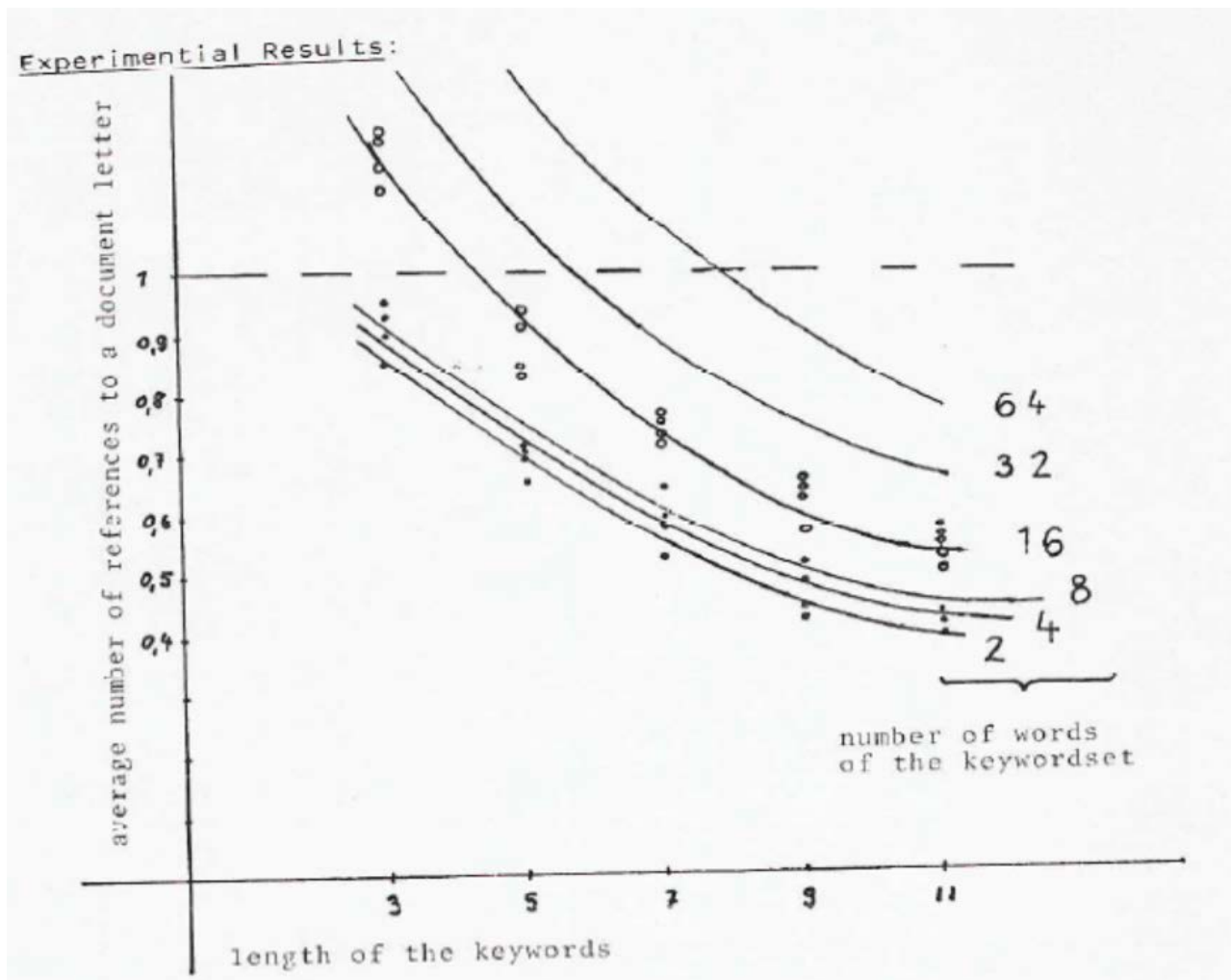
Waktu untuk fase ini adalah linear terhadap total panjang dari kata kunci yaitu W_1, \dots, W_r .

D. Waktu eksekusi

Waktu eksekusi dibagi menjadi dua bagian, yaitu untuk melakukan fase memindai dan melakukan perhitungan $\text{shift } S(v, d_{i,j})$ jika dibutuhkan. Total waktunya adalah linear terhadap total dari banyaknya perbandingan karakter.

Sama seperti algoritma pencocokan string Boyer-Moore, jika ukuran alfabet A nya besar, hanya perlu diperiksa $|D| / w_{\min}$ huruf pada kasus terburuk.

Berdasarkan eksperimen yang dilakukan terhadap 100 judul buku Computer Science dalam bahasa Inggris dan Jerman dan subjek yang terkait. Judul-judul buku adalah dokumen. Alfabet yang digunakan adalah A-Z dan 0-9 dan blank.



Gambar 5. Hasil eksperimen

Gambar di atas menggambarkan bahwa algoritma ini memiliki kompleksitas rata-rata sublinear.

IV. KESIMPULAN

Algoritma pencocokan string Commentz-Walter adalah algoritma yang cukup baik karena menggabungkan ide Boyer-Moore ditambah Aho-Corasick. Algoritma ini cocok untuk teks yang alfabetnya banyak dan bisa mencari beberapa pola sekaligus.

V. UCAPAN TERIMA KASIH

Saya ucapkan terimakasih pada dosen IF2211 karena telah membimbing saya dalam membuat makalah ini.

REFERENSI

- [1] Thomas H Corman, Charles E. Leiserson, Ronald L. Rivest & Clifford Stein "Introduction to Algorithms-String matching", EEE Edition, 2nd Edition, Page no 906-907.
- [2] D. Huson , "multiple string matching", Comp. Sequence Analysis, Nov 17 , 2004
- [3] <http://www.hs-albsig.de/studium/wirtschaftsinformatik/Documents/commentzwalterextab.pdf>
- [4] <https://www.google.co.id/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=0CCIQFjAA&url=http%3A%2F%2Fciteseerx.ist.psu.edu%2Fviewdoc%2Fdownload%3Fdoi%3D10.1.1.402.6234%26rep%3Drep1%26type%3Dpdf&ei=nPVHVaq8JeO1mwWchYHAAw&usq=AFQjCNGV-VZg7hGkPufcC5pICWpS8kNVv-Q&sig2=ajxL840CuYHXvP08pNfNuQ&bvm=bv.92291466,d.dGY>

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, 4 Mei 2015



Vicko Novianto / 13513092