

**PENGUKURAN KEMIRIPAN DOKUMEN BERBAHASA  
INDONESIA MENGGUNAKAN METODE RABIN-KARP DAN  
LEVENSHTEIN DISTANCE**

*Diajukan Sebagai Syarat untuk Menyelesaikan*

*Pendidikan Program Strata-1*

*di Jurusan Teknik Informatika*



Oleh :

DARA D KARNINDO

09121002053

**JURUSAN TEKNIK INFORMATIKA  
FAKULTAS ILMU KOMPUTER UNIVERSITAS SRIWIJAYA  
2019**

## **LEMBAR PENGESAHAN TUGAS AKHIR**

### **PENGUKURAN KEMIRIPAN DOKUMEN BERBAHASA INDONESIA MENGGUNAKAN ALGORITMA RABIN-KARP DAN LEVENSHTEIN DISTANCE**

Oleh :

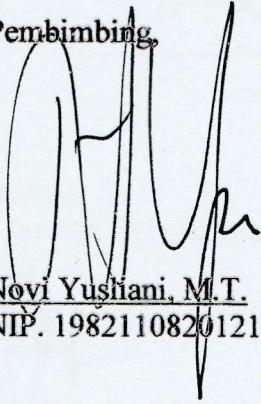
Dara D Karnindo  
09121002053

Palembang, Juli 2019

Mengetahui,  
Ketua Jurusan Teknik Informatika



Rifkie Primartha, M.T.  
NIP. 197706012009121004

Pembimbing,  
  
Novi Yusliani, M.T.  
NIP. 198211082012122001

## TANDA LULUS UJIAN SIDANG TUGAS AKHIR

Pada hari Selasa tanggal 30 Juli 2019 telah dilaksanakan ujian sidang komprehensif oleh Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya.

Nama : Dara D Karnindo

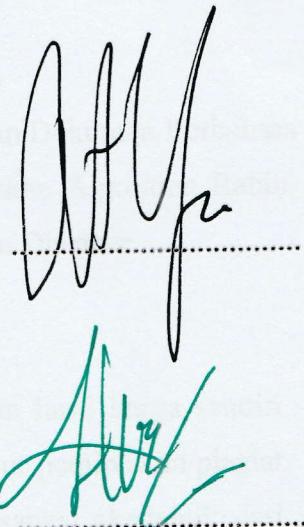
NIM : 09121002038

Judul : Pengukuran Kemiripan Dokumen Bernahasa Indonesia Menggunakan Algoritm Rabin-Karp dan Levenshtein Distance

1. Ketua Penguji

Nowi Yusliani, M.T.

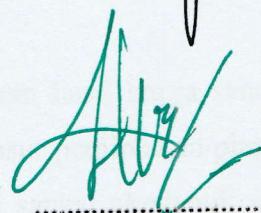
NIP. 198211082012122001



2. Penguji I

Abri Svahrini Utami, M.Kom., Ph.D.

NIP. 197812222006042003



3. Penguji II



Hafidz Novianti, M.T.

NIP. 197911012014042002

Mengetahui,  
Ketua Jurusan Teknik Informatika



Rifkie Primartha, M.T.  
NIP. 197706012009121004

## HALAMAN PERNYATAAN

Yang bertanda tangan di bawah ini :

Nama : Dara D Karnindo  
NIM : 09121002053  
Program Studi : Teknik Informatika  
Judul Skripsi : Pengukuran Kemiripan Dokumen Berbahasa Indonesia Menggunakan Algoritma Rabin-Karp dan Levenshtein Distance  
Hasil Pengecekan Software *Turnitin* : 19 %

Menyatakan bahwa Laporan Projek saya merupakan hasil karya sendiri dan bukan hasil penjiplakan/plagiat. Apabila ditemukan unsur penjiplakan/plagiat dalam laporan projek ini, maka saya bersedia menerima sanksi akademik dari Universitas Sriwijaya sesuai dengan ketentuan yang berlaku.

Demikian, pernyataan ini saya buat dengan sebenarnya dan tidak ada paksaan oleh siapapun.



Palembang, Juli 2019



Dara D Karnindo  
NIM. 09121002053

**Motto :**

“Keep Fighting! Because if you fight, many doors open ”

“To deal with big problem: Faith, Fight, Family, Friends and Fun”

Kupersembahkan karya tulis ini  
Kepada :

- ❖ Kedua orang tua tercinta
- ❖ Kakak dan Adik Saya tercinta
- ❖ Mute dan Parjim Squad
- ❖ Almamaterku
- ❖ IF Reg 2012
- ❖ Teman-teman dan Sahabat

# PENGUKURAN KEMIRIPAN DOKUMEN BERBAHASA INDONESIA MENGGUNAKAN RABIN-KARP DAN LEVENSHTEIN DISTANCE

Oleh:

Dara D Karnindo

09121002053

## ABSTRAK

Pengukuran kemiripan menjadi sangat penting untuk menghindari perlaku plagiat. Terdapat beberapa metode yang digunakan untuk mendeteksi plagiarisme. Beberapa diantaranya adalah algoritma *Rabin-Karp* dan *Levenshtein Distance*. Pada penelitian ini, penulis mencoba menggabungkan algoritma *Rabin-Karp* dan *Levenshthein Distance* untuk pengukuran kemiripan dokumen berbahasa Indonesia. Algoritma *Rabin-Karp* merupakan algoritma pencarian pola pada teks menggunakan teknik hashing. Pencocokan string akan dilakukan dengan membandingkan nilai hash antara kedua dokumen menggunakan *Levenshtein Distance*. Pengujian dilakukan dengan membandingkan dokumen uji dan dokumen asli. Kombinasi nilai n-gram, base dan modulo juga diterapkan pada pengujian ini. Hasil pengujian menunjukan, penggabungan kedua metode menghasilkan persentase kemiripan yang cukup baik. Penerapan berbagai kombinasi nilai n-gram base dan modulo menghasilkan persentase kemiripan yang berbeda-beda.

**Kata kunci:** mengukur kemiripan, pencocokan string, bahasa Indonesia, *rabin-karp*, *levenshtein distance*,

Palembang, Juli 2019

Mengetahui,  
Ketua Jurusan Teknik Informatika



Pembimbing,  
Novi Yusliani, M.T.  
NIP. 198211082012122001

# **SIMILARITY MEASURE FOR INDONESIAN DOCUMENTS USING RABIN KARP ALGORITHM AND LEVENSHTEIN DISTANCE**

**Oleh:**

**Dara D Karnindo**

**09121002053**

## **ABSTRACT**

Similarity measure is very important to avoid plagiarism behavior. There are several methods and approaches used to detect plagiarism. Some of them are Rabin-Karp algorithm and Levenshtein Distance. In this study, the author tried to combine the Rabin-Karp algorithm and Levenshthein Distance to measure the similarity of Indonesian documents. Rabin-Karp algorithm is a search algorithm that searches pattern in text using hashing techniques. String matching will be done by comparing the hash value between the two documents using Levenshtein Distance. The trial software is done by comparing the modified document from the original document. A combination of n-gram values, base and modulo is also applied to this test. From the test result, it is known that combining the two methods produces a pretty good percentage of similarity. The application of various combinations of base, n-gram and modulo also produces different percentages of similarity.

**Keyword :** similarity measure, string matching, Indonesian, rabin-karp, levenshtein distance

Palembang, Juli 2019

Mengetahui,  
Ketua Jurusan Teknik Informatika



Pembimbing  
Novi Yushiani, M.T.  
NIP. 198211082012122001

## KATA PENGANTAR

Alhamdulillah, puji dan syukur kehadirat Allah SWT atas segala nikmat, rahmat dan karunia-Nya, sehingga penulis dapat menyelesaikan Tugas Akhir yang berjudul “Pengukuran Kemiripan Dokumen Berbahasa Indonesia Menggunakan Algoritma Rabin–Karp dan Levenshtein Distance”. Tugas akhir ini disusun untuk memenuhi salah satu persyaratan kelulusan tingkat sarjana pada Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya.

Penulis menyadari bahwa dalam penulisan Tugas Akhir ini tidak terlepas dari dukungan, doa, bantuan, pengarahan maupun bimbingan dari berbagai pihak. Untuk itu penulis mengucapkan terima kasih setulus-tulusnya kepada :

1. Kedua orang tua saya dan kedua saudara saya yang selalu memberikan semangat dalam penulisan Tugas Akhir ini;
2. Pemerintah dan Universitas Sriwijaya yang telah memberikan saya kesempatan dan berbagai fasilitas dalam perkuliahan sehingga saya dapat menyelesaikan Tugas Akhir ini;
3. Bapak Jaidan, S.Pd., M.T. selaku Dekan Fakultas Ilmu Komputer Universitas Sriwijaya;
4. Bapak Rifkie Primartha, M.T. selaku Pembimbing Akademik dan Ketua Jurusan Teknik Informatika.
5. Ibu Novi Yusliani, M.T. selaku Pembimbing Tugas Akhir atas bimbingan dan masukannya penulis dapat menyelesaikan Tugas Akhir ini dengan baik;

6. Ibu Alvi Syahrini Utami, M.Kom., Ph.D. selaku dosen penguji yang telah memberikan koreksi dan masukan untuk Tugas Akhir ini;
7. Ibu Hardini Novianti, M.T. selaku dosen penguji yang telah memberikan koreksi dan masukan untuk Tugas Akhir ini;
8. Segenap staf pengajar di Fakultas Ilmu Komputer Universitas Sriwijaya yang telah mengajar, membimbing dan memberikan pemahaman tentang ilmu computer;
9. Segenap karyawan Fakultas Ilmu Komputer Universitas Sriwijaya, terutama Kak Ricy Firnando selaku Admin Jurusan Teknik Informatika atas bantuannya selama ini;
10. Mute ,Wenty Octaviani, Putri Septria, Lisa Desta Sari dan Riza Gamal Fuad yang telah memberi semangat, membantu dan menemani saya dalam penulisan Tugas Akhir ini;
11. Teman – teman Parjim Squad: Dwi Erviana, Elbananda Permana, Putri Septria, Rahmi Fadhillah Busyra dan Wenty Octaviani yang selalu menunggu, mendukung dan mendoakan saya selama penulisan Tugas Akhir ini
12. Semua anggota Teknik Informatika Regular 2012 yang telah memberi semangat kepada saya dalam menyelesaikan Tugas Akhir ini;
13. Semua teman-teman dan sahabat yang telah memberikan semangat kepada saya dalam penulisan Tugas Akhir ini.

Akhir kata, penulis menyadari bahwa laporan tugas akhir ini masih jauh dari kesempurnaan karena keterbatasan ilmu yang dimiliki penulis. Oleh karena itu, penulis mengharapkan kritik dan saran yang membangun untuk kesempurnaan tugas akhir ini. Semoga tugas akhir yang sederhana ini dapat memberikan manfaat bagi yang membutuhkan.

Palembang, Juli 2019

Penulis

## DAFTAR ISI

Halaman

HALAMAN JUDUL .....	i
LEMBAR PENGESAHAN .....	ii
DAFTAR ISI .....	iii
DAFTAR GAMBAR .....	v
DAFTAR TABEL .....	vi

### **BAB I PENDAHULUAN**

1.1 Pendahuluan.....	I-1
1.2 Latar Belakang .....	I-1
1.3 Perumusan Masalah .....	I-4
1.4 Tujuan Penelitian .....	I-4
1.5 Manfaat Penelitian .....	I-4
1.6 Batasan Masalah .....	I-5
1.7 Sistematika Penulisan.....	I-5
1.8 Kesimpulan.....	I-6

### **BAB II TINJAUAN PUSTAKA**

2.1 Pendahuluan.....	II-1
2.2 Plagiarisme .....	II-1
2.3 <i>Preprocessing</i> .....	II-4
2.4. <i>Shingle</i> .....	II-11
2.5 Algoritma <i>Rabin-Karp</i> .....	II-12
2.6 <i>Levenshtein Distance</i> .....	II-13
2.7 Pengukuran Kemiripan.....	II-15

2.8	Penelitian Terkait .....	II-15
-----	--------------------------	-------

### **BAB III METODOLOGI PENELITIAN**

3.1	Pendahuluan .....	III-1
3.2	Pengumpulan Data .....	III-1
3.2.1	Jenis dan Sumber Data .....	III-1
3.2.2	Teknik Pengumpulan Data .....	III-1
3.3	Tahapan Penelitian .....	III-2
3.3.1	Analisa Dokumen Masukan .....	III-2
3.3.2	Menetapkan Kerangka Kerja/Framework .....	III-5
a.	<i>Preprocessing</i> .....	III-6
b.	Algoritma <i>Rabin-Karp</i> .....	III-6
c.	<i>Levenshtein Distance</i> .....	III-6
d.	Perhitungan Kemiripan.....	III-7
3.3.3	Menetapkan Kriteria Pengujian .....	III-8
3.3.4	Menentukan Format Data Pengujian .....	III-8
3.3.5	Menentukan Alat yang digunakan dalam Pelaksanaan Penelitian .....	III-10
3.3.6	Melakukan Pengujian dan Analisis Hasil Pengujian..	III-10
3.4	Metode Pengembangan Perangkat Lunak .....	III-11
3.4.1	Fase Insepsi .....	III-11
3.4.2	Fase Elaborasi.....	III-12
3.4.3	Fase Konstruksi.....	III-12
3.4.4	Fase Transisi.....	III-13
3.5	Manajemen Proyek Penelitian .....	III-13

### **BAB IV PENGEMBANGAN PERANGKAT LUNAK**

4.1	Pendahuluan .....	IV-1
4.2	Fase Insepsi .....	IV-1
4.2.1	Pemodelan Bisnis .....	IV-2

4.2.2	Kebutuhan Sistem .....	IV-4
a	Fitur Prapengolahan.....	IV-4
b	Fitur Hashing.....	IV-4
c	Fitur Pencocokan .....	IV-4
d	Fitur Perhitungan Kemiripan.....	IV-4
4.3	Analisis dan Desain .....	IV-5
4.3.1	Analisis Arsitektur Perangkat Lunak .....	IV-5
4.3.2	Analisis Data .....	IV-5
4.3.3	Analisis Dokumen Masukan .....	IV-5
4.3.4	Analisis Pra-pengolahan .....	IV-6
4.3.5	Hashing .....	IV-7
4.3.6	Pencocokan String .....	IV-12
4.3.7	Perhitungan Kemiripan .....	IV-12
4.3.8	Desain Perangkat Lunak .....	IV-13
1	<i>Use Case</i> .....	IV-13
a.	Definisi Aktor.....	IV-14
b.	Definisi <i>Use Case</i> .....	IV-14
2.	Skenario <i>Use Case</i> .....	IV-14
4.4	Fase Elaborasi .....	IV-17
4.4.1	Pemodelan Bisnis .....	IV-17
4.4.2	Perancangan Data .....	IV-18
4.4.3	Perancangan Antarmuka .....	IV-18
4.4.4	Kebutuhan Sistem .....	IV-19
4.4.5	Diagram Aktivitas .....	IV-19
4.4.6	Diagram Alur .....	IV-22
4.5	Fase Konstruksi .....	IV-24
4.5.1	Diagram Kelas .....	IV-24
4.5.2	Implementasi .....	IV-25
4.5.3	Implementasi Kelas .....	IV-25
4.5.4	Implementasi Antarmuka .....	IV-25

4.6 Fase Transisi.....	IV-27
4.6.1 Kebutuhan Sistem .....	IV-27
4.6.2 Rencana Pengujian .....	IV-27

## **BAB V HASIL DAN ANALISIS PENELITIAN**

5.1 Pendahuluan .....	V-1
5.2 Data Hasil Pengujian.....	V-1
5.2.1 Konfigurasi Percobaan I.....	V-2
5.2.2 Konfigurasi Percobaan II .....	V-6
5.2.3 Hasil Uji N-gram Terhadap Persentase Kemiripan.....	V-7
5.2.4 Hasil Uji Base Terhadap Persentase Kemiripan .....	V-8
5.2.3 Hasil Uji Modulo Terhadap Persentase Kemiripan .....	V-9

## **BAB VI KESIMPULAN DAN SARAN**

6.1 Kesimpulan .....	VI-1
6.2 Saran.....	VI-1

DAFTAR PUSTAKA.....

LAMPIRAN .....

## DAFTAR TABEL

	Halaman
II-1 Aturan Pemenggalan Prefiks "me" .....	II-7
II-2 Aturan Pemenggalan Prefiks "pe" .....	II-8
II-3 Aturan Pemenggalan Prefiks "be" .....	II-9
II-4 Aturan Pemenggalan Prefiks "te" .....	II-9
II-5 Modifikasi tabel 3 .....	II-10
II-6 Modifikasi tabel 4 .....	II-10
II-7 Proses <i>Preprocessin</i> .....	II-11
III-1 Corpus Dokumen Plagiat .....	III-2
III-2 Contoh Corpus Dokumen Plagiat .....	III-3
III-3 Inisialisasi awal Algoritma <i>Levenshtein</i> .....	III-7
III-4 Rancangan Hasil Pendekripsi .....	III-8
III-5 Rancangan Pengujian Nilai Parameter .....	III-8
III-6 Rancangan Hasil Uji Tiap Parameter .....	III-8
IV-1. Kebutuhan Fungsional Perangkat Lunak .....	IV-4
IV-2. Kebutuhan Non-Fungsional Perangkat Lunak .....	IV-4
IV-3. Contoh Proses Shingle .....	IV-7
IV-4. Hasil Shingle dengan n=3 .....	IV-8
IV-5. Nilai Hash Pada Kedua String1 dan String2 .....	IV-11

IV-6. Hasil Perhitungan Levenshtein Distance .....	IV-12
IV-7. Definisi Aktor .....	IV-14
IV-8. Definisi Use Case .....	IV-14
IV-9. Skenario Use Case Mengukur Kemiripan Dokumen .....	IV-15
IV-10 Skenario Use Case Melakukan Pra-pengolahan.....	IV-16
IV-11 Hasil Perhitungan Levenshtein Distance .....	IV-12
IV-12 Daftar Kelas .....	IV-25
IV-13 Rencana Pengujian Use Case Mengukur Kemiripan Dokumen .....	IV-27
IV-14 Rencana Pengujian Use Case Melakukan Pra-pengolahan .....	IV-28
IV-15 Pengujian Use Case Mengukur Kemiripan Dokumen .....	IV-29
IV-10 Pengujian Use Case Melakukan Pra-pengolahan .....	IV-30
V-1. Rincian Dokumen Pembanding .....	V-1
V-2. Rincian Dokumen Uji .....	V-2
V-3. Hasil Percobaan I .....	V-3
V-4. Pengujian Nilai Parameter Terhadap Hasil Persentase Kemiripan ....	V-6
V-5. Hasil Uji N-gram .....	V-7
V-6. Hasil Uji Base .....	V-8
V-7. Hasil Uji Base 2 .....	V-8
V-8. Hasil Uji Modulo .....	V-9

## **DAFTAR GAMBAR**

	Halaman
II-1 Tahapan <i>External Plagiarism Detection</i> .....	II-3
III-1 Diagram Tahapan Perangkat Lunak .....	III-5
III-2 Penjadwalan Penelitian dalam Bentuk WBS .....	III-14
III-3 Gantt Chart Penjadwalan untuk Tahap Menentukan Ruang Lingkup dan Unit Penelitian.....	III-17
III-4 Gantt Chart Penjadwalan untuk Tahap Menentukan Dasar Teori yang Berkaitan dengan Penelitian.....	III-17
III-5 Gantt Chart Penjadwalan untuk Tahap Menentukan Kriteria Pengujian.III-18	
III-6 Gantt Chart Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Melakukan Penelitian Fase Insepsi .....	III-18
III-7 Gantt Chart Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Melakukan Penelitian Fase Elaborasi .....	III-19
III-8 Gantt Chart Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Melakukan Penelitian Fase Konstruksi.....	III-19
III-9 Gantt Chart Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Melakukan Penelitian Fase Transisi .....	III-20
III-10 Gantt Chart Penjadwalan untuk Tahap Melakukan Pengujian Penelitian.....	III-20
III-11 Gantt Chart Penjadwalan untuk Tahap Analisa Hasil Pengujian dan Pembuatan Kesimpulan.....	III-21
IV-1. Diagram Proses Pra-pengolahan.....	IV-6

IV-2. Diagram Use Case .....	IV-13
IV-3. Rancangan Antarmuka Perangkat Lunak .....	IV-19
IV-4. Diagram Aktivitas Mengukur Kemiripan Dokumen .....	IV-21
IV-5. Diagram Aktivitas Melakukan Pra-Pengolahan .....	IV-22
IV-6. Diagram Sekuen Mengukur Kemiripan Dokumen .....	IV-23
IV-7. Diagram Kelas .....	IV-24
IV-8. Antarmuka Perangkat Lunak Sebelum Pengukuran Kemiripan.....	IV-26
IV-9. Antarmuka Perangkat Lunak Sebelum Pengukuran Kemiripan.....	IV-26
IV-10. Diagram Kelas .....	IV-25
IV-11. Rancangan Antarmuka .....	IV-27
IV-12. Tampilan Antarmuka Perangkat Lunak Sebelum Proses .....	IV-33
IV-13. Tampilan Antarmuka Perangkat Lunak Setelah Proses .....	IV-33
V-1. Grafik Analisis Nilai Kesalahan Untuk Setiap Kategori Dokumen ....	V-4
V-2. Grafik Perbandingan Hasil Kemiripan Sistem dengan Sebenarnya ..	V-5

## **BAB I**

### **PENDAHULUAN**

#### **1.1 Pendahuluan**

Pada bab ini membahas latar belakang masalah, rumusan masalah, tujuan, dan manfaat penelitian, serta batasan masalah. Bab ini akan memberikan penjelasan umum mengenai keseluruhan penelitian.

Pendahuluan dimulai dengan penjelasan mengenai tantangan dan tujuan proses menemukan pengetahuan baru pada deteksi kemiripan teks serta penelitian yang berkaitan dengan menerapkan algoritma *rabin-karp* dan *levenshtein distance* yang menjadi latar belakang dari penelitian ini.

#### **1.2 Latar Belakang**

Dalam dunia akademik, terutama dalam proses penulisan karya ilmiah atau skripsi, ditemukan banyak terjadi peniruan atau penjiplakan isi tulisan, baik dalam bentuk sejumlah kalimat, paragraf, bab, dan bahkan penjiplakan total tanpa ada yang diubah sama sekali, dan tanpa memberikan sumber identitas sang penulis asli. Peristiwa ini disebut sebagai plagiarisme. Kata plagiarisme berakar dari dua kata Latin, *plagiarius* yang berarti penculik, dan *plagiare* yang berarti mencuri. Plagarisme dapat dikategorikan sebagai tindakan kriminal pelanggaran hak cipta karena dilakukan dengan cara menerbitkan karya orang lain dan menjadikannya seolah-olah karya milik sendiri (Joy & Luck, 1999).

Bull, et. al. (2001) pada penelitiannya mengkhawatirkan pengembangan yang pesat dalam penggunaan internet sebagai sarana belajar dan penelitian, akan meningkatkan praktik plagiat di berbagai bidang akibat dari kemudahan mengakses materi di internet. Masalah plagiarisme ini tumbuh menjadi sangat serius di kalangan akademis terutama universitas-universitas di Indonesia. Ini dibuktikan dengan banyaknya ditemukan plagiarisme pada *paper* atau proyek tugas akhir mahasiswa (Krisnawati & Schulz, 2014)

Untuk mengatasi hal tersebut perlu dilakukan pendekripsi terhadap dokumen-dokumen yang dicurigai sebagai dokumen plagiat. Deteksi plagiarisme dapat dilakukan secara manual dan otomatis. Deteksi secara manual dilakukan oleh manusia, namun memiliki kelemahan dari segi keefektifan jika dilakukan pada dokumen dalam jumlah yang sangat banyak (Clough, 2003). Oleh karena itu diperlukan suatu sistem yang dapat mendekripsi plagiarisme secara otomatis.

Beberapa metode pencocokan string seperti *Levenshtein Distance* (LD) dan *Smith-Waterman* (SW) (Su et al., 2008), *Conceptual Similarity* (CS) dan *Graph-Based* (GB) (Osman, Salim, Binwahlan, Hentably, & Ali, 2011), *Rabin-Karp* (RK) (Firdaus, 2003) ((Salmuasih & Sunyoto, 2013) telah diusulkan untuk mendekripsi plagiarisme pada dokumen. Metode LD mempunyai kelebihan akurasi yang baik pada proses pencocokan string, yaitu dengan cara membandingkan rangkaian string secara keseluruhan, namun algoritma ini lemah pada efisiensi waktu jika diterapkan pada aplikasi deteksi yang berskala besar (berbasis web) (Su et al, 2008). Metode RK mempunyai kelebihan sangat baik untuk pencarian *string* dengan pola banyak dan

ukuran yang besar (Firdaus, 2003), dan tetap baik jika digunakan dalam aplikasi berbasis web (Mutiara & Agustina, 2008). Namun, *rabin-karp* hanya menghitung berdasarkan jumlah *hash* yang memiliki nilai yang sama pada kedua dokumen yang dibandingkan. Jumlah *identical hash* tentu saja akan berpengaruh pada persentasi kemiripan kedua dokumen. Karena itu, penulis mencoba mengkombinasikan *levenshtein distance* dan *rabin-karp* pada penelitian ini.

*Levenshtein distance* akan digunakan untuk menggantikan perhitungan kemiripan *hash* pada *rabin-karp*. Algoritma ini bekerja dengan menghitung jarak antara *string* satu dengan lainnya sehingga derajat kemiripan kedua dokumen dapat ditentukan berdasarkan bobot jaraknya. Dengan menggunakan *levenshtein distance* hasil perhitungan jarak *hash* antara kedua dokumen diharapkan mampu mempengaruhi dan menghasilkan tingkat kemiripan yang lebih baik.

### 1.3 Perumusan Masalah

Berdasarkan latar belakang masalah yang telah dijelaskan, rumusan masalah pada penelitian ini adalah:

1. Bagaimana mengukur kemiripan dari dua dokumen menggunakan algoritma *rabin-karp* dan *levenshtein distance*?
2. Bagaimana persentase kemiripan dokumen berbahasa indonesia yang dihasilkan menggunakan algoritma *rabin-karp* dan *levenshtein distance*?
3. Bagaimana pengaruh nilai n-gram, base dan modulo terhadap persentase kemiripan dokumen yang dihasilkan?

## 1.4 Tujuan Penelitian

Tujuan penelitian ini adalah sebagai berikut :

1. Menerapkan algoritma *rabin-karp* dan *levenshtein distance* ke dalam perangkat lunak untuk pengukuran kemiripan dokumen berbahasa Indonesia.
2. Mengetahui persentase kemiripan dokumen menggunakan algoritma *rabin-karp* dan *levenshtein distance*.
3. Mengetahui pengaruh n-gram, base dan modulo terhadap persentase kemiripan yang dihasilkan .

## 1.5 Manfaat Penelitian

Manfaat penelitian ini adalah sebagai berikut :

1. Hasil penelitian ini dapat digunakan sebagai rujukan dalam penelitian yang terkait dengan deteksi plagiarisme;
2. Memahami Algoritma *Rabin Karp* dan *Levenshtein Distance* sebagai metode pada pendekripsi kemiripan dokumen berbahasa Indonesia

## 1.6 Batasan Masalah

Batasan masalah pada penelitian ini adalah sebagai berikut:

1. Data set yang digunakan adalah dokumen bahasa Indonesia.
2. Penelitian ini hanya mengukur persentase kemiripan dari dua dokumen.

3. Perangkat lunak hanya menerima dokumen berekstensi\*.txt
4. Data yang diuji hanya berupa teks.
5. Sistem tidak memperhatikan kesalahan ejaan ataupun penulisan pada dokumen.

## **1.7 Sistematika Penulisan**

Sistematika penulisan skripsi ini adalah sebagai berikut:

### **BAB I. PENDAHULUAN**

Bab I berisi uraian mengenai latar belakang, perumusan masalah, tujuan dan manfaat penelitian, batasan masalah, metodologi penelitian, model proses pengembangan perangkat lunak dan sistematika penulisan.

### **BAB II. KAJIAN LITERATUR**

Bab II berisi landasan dasar teori yang digunakan dalam penelitian ini, seperti analisis *preprocessing*, algoritma *Rabin-Karp*, *Levenshtein Distance*, dan penelitian-penelitian lain yang relevan dengan penelitian ini.

### **BAB III. METODOLOGI PENELITIAN**

Pada bab III ini akan membahas mengenai tahapan yang akan dilaksanakan pada penelitian ini. Pada tahapan penelitian ini akan dibahas dengan rinci dengan mengacu pada suatu kerangka kerja. Serta pada bagian akhir bab ini akan dijabarkan perancangan manajemen proyek perangkat lunak untuk pelaksanaan penelitian ini.

## **1.8 Kesimpulan**

Pada penelitian ini metode yang digunakan untuk mengukur kemiripan teks bahasa Indonesia yaitu algoritma *rabin-karp* dan *levenshtein distance*. Pengujian akan dilakukan dengan melihat persentase hasil kemiripan teks dari sistem. Kemudian hasil pengujian akan dibandingkan dengan hasil kemiripan yang telah ditentukan secara manual.

## DAFTAR PUSTAKA

- Arifin, A. Z., Mahendra, I. P. A. K., & Ciptaningtyas, H. T. (2009). ENHANCED CONFIX STRIPPING STEMMER AND ANTS ALGORITHM FOR CLASSIFYING NEWS DOCUMENT IN INDONESIAN LANGUAGE. *The 5th International Conference on Information & Communication Technology and Systems.*
- Baruah, D., & Kakoti Mahanta, A. (2013). A New Similarity Measure with Length Factor for Plagiarism Detection. *International Journal of Computer Applications*, 72(14), 14–17.  
<https://doi.org/10.5120/12561-8671>
- Clough, P. (2003). Old and new challenges in automatic plagiarism detection. *National Plagiarism Advisory Service.*
- Firdaus, H. B. (2003). Deteksi Plagiat Dokumen Menggunakan Algoritma Rabin-Karp. *Jurnal Ilmu Komputer Dan Teknologi Informasi.*
- Gusmita, R. H., Durachman, Y., Harun, S., Firmansyah, A. F., Sukmana, H. T., & Suhaimi, A. (2014). A rule-based question answering system on relevant documents of Indonesian Quran Translation. *2014 International Conference on Cyber and IT Service Management, CITSM 2014*. <https://doi.org/10.1109/CITSM.2014.7042185>
- Haryanto, E. V. (2011). Rancang Bangun Prototype Mesin Pencari String Menggunakan Metode Fuzzy String Matching. *Konfensi Nasional Sistem Dan Infromatika*, 1(Pencarian Infromasi), 1–26.
- Joy, M., & Luck, M. (1999). Plagiarism in programming assignments. *IEEE Transactions on Education*. <https://doi.org/10.1109/13.762946>
- Krisnawati, L. D., & Schulz, K. U. (2014). *Plagiarism Detection for Indonesian Texts*. <https://doi.org/10.1145/2539150.2539213>
- Martin, B. (1994). Plagiarism: A Misplaced Emphasis, by Brian Martin. *Information Ethics*.

Retrieved from <https://www.uow.edu.au/~bmartin/pubs/94jie.html>

- Mutiara, A. B., & Agustina, S. (2008). *Anti Plagiarism Application with Algorithm Karp-Rabin at Thesis in Gunadarma University.* 3–10. Retrieved from <http://arxiv.org/abs/0811.4349>
- Nugroho, E. (2011). Perancangan Sistem Deteksi Plagiarisme Dokumen Teks Dengan Menggunakan Algoritma Rabin-Karp. In *Journal of Strategic Studies*. <https://doi.org/10.1080/01402390.2011.569130>
- Osman, A. H., Salim, N., Binwahlan, M. S., Hentably, H., & Ali, A. M. (2011). Conceptual similarity and graph-based method for plagiarism detection. *Journal of Theoretical and Applied Information Technology*.
- Pressman, R. S., & Maxim, B. R. (2015). Software Engineering : A Practitioner's Approach, Eighth Edition. In *ACM SIGSOFT Software Engineering Notes*. <https://doi.org/10.1145/1226816.1226822>
- Salmuasih, & Sunyoto, A. (2013). Implementasi Algoritma Rabin Karp untuk Pedeteksi Plagiat Dokumen Teks Menggunakan Konsep Similarity. *Seminar Nasional Aplikasi Teknologi Informasi (SNATi)*.
- Stein, B., & zu Eissen, S. M. (2006). *Near Similarity Search and Plagiarism Analysis*. [https://doi.org/10.1007/3-540-31314-1\\_52](https://doi.org/10.1007/3-540-31314-1_52)
- Su, Z., Ahn, B. R., Eom, K. Y., Kang, M. K., Kim, J. P., & Kim, M. K. (2008). Plagiarism detection using the Levenshtein distance and Smith-Waterman algorithm. *3rd International Conference on Innovative Computing Information and Control, ICICIC'08*. <https://doi.org/10.1109/ICICIC.2008.422>