

PEMBUATAN APLIKASI PENGUKURAN TINGKAT KEMIRIPAN DOKUMEN BERBASIS WEB MENGGUNAKAN ALGORITMA WINNOWER

*Nur Fadillah Ulfa¹
Metty Mustikasari²*

^{1,2}Fakultas Ilmu Komputer dan Teknologi Informasi Universitas Gunadarma,

¹nurfadillahulfa@gmail.com

²metty@staff.gunadarma.ac.id

Abstrak

Plagiarisme atau penjiplakan adalah suatu tindakan pengambilan karangan atau pendapat orang lain dan menjadikannya seolah karangan tersebut sebagai pendapat sendiri. Untuk mengantisipasinya, dibutuhkan suatu cara yang dapat menganalisis teknik plagiat yang dilakukan. Algoritma Winnowing merupakan salah satu algoritma pada metode document fingerprinting. Metode ini akurat dalam mengidentifikasi penyalinan teks termasuk bagian kecil yang mirip dalam sekumpulan dokumen melalui fingerprint yang dihasilkan. Melalui pencocokan fingerprint akan diperoleh nilai similarity antar dokumen. Semakin kecil tingkat persentase kesamaan dokumen teks yang diuji, maka dokumen tersebut tidak termasuk plagiat, tetapi jika hasil dari pengujian pada dua dokumen semakin besar, maka disimpulkan bahwa dokumen tersebut menyerupai tindakan plagiat.

Kata kunci : Algoritma Winnowing, Dokumen, Similaritas, Plagiarism.

DESIGNING WEB-BASED APPLICATION TO MEASURE THE LEVELS OF SIMILARITY OF THE DOCUMENTS USING WINNOWER ALGORITHM

Abstract

Plagiarism is an act of taking essay or opinion of others and make it as the wreath is as his own opinion. In anticipation of this, we need a way to analyze the technique of plagiarism committed. Winnowing algorithm is one of the methods of document fingerprinting algorithms. The result shows that this method is accurate to identify the plagiarism of the text, even small parts that are similar among documents through the generated fingerprint. Through fingerprint, this would obtain the value of similarity between documents. The smaller the similarity percentage level in the document, the less the value of plagiarism found in the document. However, if the results of the tests showed bigger similarity, then it could be concluded that the documents are resembles. This is indicating an act of plagiarism

Keywords : Winnowing Algorithms, Documents, Similarity, Plagiarism

PENDAHULUAN

Pemanfaatan teknologi digital telah menjadi kebutuhan dalam era modern saat ini. Komponen yang ada di dalam dunia digital

salah satunya adalah dokumen teks.

Dokumen dalam bentuk digital memudahkan dalam hal penyimpanan, efisien, mudah dicari, bahkan mudah dalam hal penjiplakan.

Penjiplakan atau plagiarisme berarti mencontoh atau meniru atau mencuri tulisan dan karya orang lain yang kemudian diaduk sebagai karangannya sendiri dengan atau tanpa persetujuan penulisnya. Plagiarisme berasal dari bahasa Latin *plagiaris* yang berarti merampok, atau membajak. Plagiarisme merupakan tindakan pencurian atau kebohongan intelektual. Plagiarisme adalah tindakan penyalahgunaan, pencurian/perampasan, penyerbitan, pernyataan, atau menyatakan sebagai milik sendiri sebuah pikiran, ide, tulisan, atau ciptaan yang sebenarnya milik orang lain [1]. Penjiplakan dokumen digital bukanlah hal yang susah, cukup dengan menggunakan teknik *copy-paste-modify* pada sebagian isi dokumen dan bahkankeseluruhan isi dokumen sudah bisa dikatakan bahwa dokumen tersebut merupakan hasil duplikasi dari dokumen lain.

Praktek penjiplakan dokumen ini seringkali diterapkan oleh akademisi baik tingkat sekolah maupun perguruan tinggi. Tindakan plagiat yang dilakukan oleh siswa maupun mahasiswa ini sangattidak mencerminkan sikap kreatif dan terpelajar sebagai kaum intelektual. Demi menyelesaikan tugas-tugasnya dengan cepat, seseorang dapat melakukan teknik *copy-paste-modify* tanpa perlu mempelajari dan mengeksplorasi materi terlebih dahulu. Kadang kala tindakan penjiplakan ini dimodifikasi dengan mengganti kata-kata yang mengandung sinonim, dengan maksud agar terlihat berbeda dari pekerjaan teman. Hal semacam ini dapat menimbulkan masalah terhadap evaluasi hasil belajarsiswa/mahasiswa.

Kesamaan dokumen bukan hanya ditinjau dari isi kata yang digunakan sebagai penyusunan kalimat atau ahsama, akan tetapi juga dikanmiripkan apabila isi dokumen memiliki makna yang sama. Penelitian pengukuran kesamaan dokumen Bahasa Indonesia yang ada hanya mengukur kesamaan kata atau pun kalimat, belum mempertimbangkan struktur kalimat, jumlah kalimat, posisi kalimat dan makna kata untuk membandingkan kalimat [2].

Proses pendeteksian penjiplakan ini menggunakan algoritma *Winnowing* yang mana *output*-nya berupa sekumpulan nilai *hash* yang didapatkan melalui metode *k-gram*. Sedangkan konsep *synonym recognition* ini dimaksudkan untuk dapat mengenali kata-kata yang mengandung sinonim sebagai tindak penjiplakan.

METODE PENELITIAN

Algoritma Winnowing

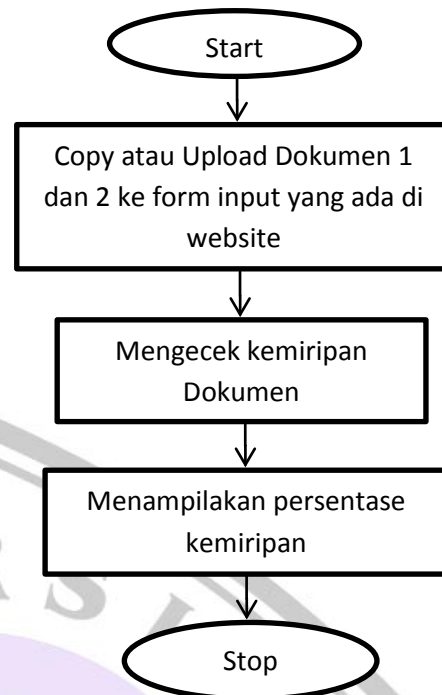
Banyak cara atau metode yang dapat digunakan untuk mendeteksi penjiplakan dalam file text. Namun ada 3 kebutuhan mendasar yang harus dipenuhi oleh algoritma deteksi penjiplakan [3]. Kebutuhan pertama adalah *Whitespace Insensitivity* yang berarti dalam melakukan pencocokan terhadap file teks seharusnya tidak terpengaruh oleh spasi, jenis huruf, tanda baca dan sebagainya. Kedua, *noise Suppression* yang berarti menghindari penemuan kecocokan dengan panjang kata yang terlalu kecil atau kurang relevan, misal: „the“. Panjang kata yang ditengarai merupakan penjiplakan harus cukup untuk membuktikan bahwa kata-kata

tersebut telah dijiplak dan bukan merupakan kata yang umum digunakan. Ketiga, *position Independence* yang berarti penemuan kecocokan/kesamaan harus tidak bergantung pada posisi kata-kata. Meskipun posisinya tidak sama, kecocokan harus dapat ditemukan.

Winnowing adalah algoritma yang digunakan untuk melakukan proses document fingerprinting. Proses ini ditujukan agar dapat mengidentifikasi penjiplakan, termasuk bagian-bagian kecil yang mirip dalam dokumen yang berjumlah banyak. Input dari proses document fingerprinting adalah file teks. Kemudian outputnya akan berupa sekumpulan nilai hash yang disebut fingerprint. Fingerprint inilah yang akan dijadikan dasar perbandingan antara file-file teks yang telah dimasukkan [4].

Analisis Metode

Analisis metode ini merupakan penjelasan mengenai tahap-tahap pada aplikasi pendeteksian dokumen similaritas. Tahap pertama yang dilakukan adalah copy atau upload dokumen 1 dan 2 ke form input yang ada di website, tahap kedua yaitu mengecek kemiripan dokumen, kemudian tahap ketiga adalah menampilkan persentase kemiripan dan selesai. Pada Gambar 1 berikut ditunjukkan diagram alur untuk penelitian ini.



Gambar 1. Aplikasi Pendeteksian Dokumen Similaritas

Pada algoritma winnowing ini terdapat beberapa tahap yaitu tahap pertama penghapusan karakter-karakter yang tidak relevan (white-space insensitivity), antara lain spasi atau tanda baca. Kemudian tahap kedua yaitu pembentukan rangkaian gram dengan ukuran k . Tahap ketiga adalah perhitungan nilai hash menggunakan rolling hash kedalam fingerprinting dan tahap terakhir adalah menentukan persentase kesamaan antara 2 dokumen dengan persamaan Jaccard Coefficient.

Berikut ini adalah contoh kasus dan cara manual untuk menghitung nilai similaritas menggunakan algoritma winnowing :

Contoh kasus

Teks 1 : bunga mawar merah

Teks 2 : bunga mawar putih

1. Langkah pertama :

menghilangkan tanda baca dan spasi

Teks 1 : bungamawarmerah

Teks 2 : bungamawarputih

2. Langkah kedua :

pembentukanrangkaiannilai k-gram denganukuran 5.

Teks1

:bungaungamngamagamawamawamawarawarmwarmearmerrmeramerah

Teks2

:bungaungamngamagamawamawamawarawarpwarpuarputrputiputih

3. Langkahketiga :

melakukanperhitungannilai-nilai hash dari setiap gram menggunakan rolling hash. Formula untukmenghitungnilai hash dapatdilihatpadapersamaan (1) berikut.

$$c_1 * b^{(k-1)} + c_2 * b^{(k-2)} + \dots + c_{(k-1)} * b + c_k \quad (1)$$

Menghitungnilai hash dari kata “bunga” :

$$b * 11^4 + u * 11^3 + n * 11^2 + g * 11^1 + a * 11^0 = 98 * 14641 + 117 * 1331 + 110 * 121 + 103 * 11 + 97 * 1 = 1434818 + 155727 + 13310 + 1133 + 97 = 1605085$$

Setelah mendapatkan nilai hash dari kata “bunga” maka untuk mencari nilai hash kata kedua yaitu “ungam” tidak perlu menggunakan rumus 1 lagi, karena pada kata kedua terdapat juga karakter pada kata pertama sehingga menggunakan rumus kedua untuk mencari nilai hash pada kata kedua dan seterusnya yang dapatdilihatpadapersamaan (2) berikut.

$$H_{(c_2 \dots c_{k+1})} = (H_{(c_1 \dots c_k)} - c_1 * b^{(k-1)}) * b + c_{(k+1)} \quad (2)$$

$$H_{c_2} = (1605085 - c_1 * 11^4) * 11 + c_{10} = (1605085 - 98 * 14641) * 11 +$$

$$109 = (1605085 - 1434818) * 11 + 109 = 170267 * 11 + 109 = 1873046$$

Dengan begini tidak perlu melakukan iterasi dari indeks pertama sampai terakhir untuk menghitung nilai hash untuk gram ke-2 sampai terakhir. Hal ini tentu dapat menghemat biaya komputasi saat menghitung nilai hash dari sebuah gram. Hasil perhitungan nilai-nilai hash dari setiap gram menggunakan rolling hash Teks 1:

1605085, 1873046, 1760636, 1651505, 1578399, 1740556, 1591666, 1886480, 1586325, 1827725, 1745265.

4. Langkahkeempat :membentuk window dari nilai-nilai hash denganukuran 4.

(1605085 1873046 1760636 1651505),
(1873046 1760636 1651505 1578399),
(1760636 1651505 1578399 1740556),
(1651505 1578399 1740556 1591666),
(1578399 1740556 1591666 1886480),
(1740556 1591666 1886480 1586325),
(1591666 1886480 1586325 1827725),
(1886480 1586325 1827725 1745265).

5. Langkahkelima :memilihnilai hash terkecildari setiap widow untukdijadikansebagai fingerprint, bilaterdapatnilai minimum hash yang samamakahanyaditulissatusajayaitu nilai hash yang paling kecildpada window yang pertamaditemukannilai hash tersebut.

Hasil fingerprint Teks1 : formula padapersamaan (3)
[1605085,0] [1578399,4] berikut.

$$\text{Similaritas } (d_i d_j) = \frac{|w(d_i) \cap (d_j)|}{|w(d_i) \cup (d_j)|} \quad (3)$$

hash-nya tidak perlu dituliskan lagi. Nilai 0, 4 dan 8 merupakan nilai indeks dari hash yang terbentuk oleh k-gram. Parameter nilai pada window digunakan untuk mengam- bil perwakilan nilai hash sebagai bagian fingerprint yang tepat.

Untuk perhitungan teks 2 sama seperti teks 1, berikut ini adalah nilai-nilai hash dan pembentukan window dari teks 2 :

Nilai hash Teks 2:

1605085 1873046 1760636

1651505 1578399 1740556

1591669 1886529 1586866

1833684 1810814

Nilai Window Teks 2

dengan ukuran 4 :

(1605085 1873046 1760636

1651505)

(1873046 1760636 1651505
1578399)

(1760636 1651505 1578399
1740556)

(1651505 1578399 1740556
1591669)

(1578399 1740556 1591669
1886529)

(1740556 1591669 1886529
1586866)

(1591669 1886529 1586866
1833684)

(1886529 1586866 1833684
1810814)

Hasil fingerprint teks2 :

[1605085,0] [1578399,4]

[1586866,8]

6. Langkah keenam

: pengukuran nilai similaritas.

Nilai similaritas ditentukan dengan

$$S = \frac{2}{4} \times 100\% = 50\%$$

Jadi hasil dari tingkat kesamaan teks 1 dan teks 2 adalah 50%

Use case Diagram

Padagambar *Use case* untuk *user* ini terdapat langkah-langkah yang akan dijelaskan, antara lain langkah pertama adalah user dapat meng-upload dokumen, kemudian user input k-gram dan window. Setelah itu dokumen di cek dan menampilkan hasil presentasinya a. Use case diagram Upload Files dapat dilihat pada Gambar 2 berikut.



Gambar 2. Use case diagram Upload Files

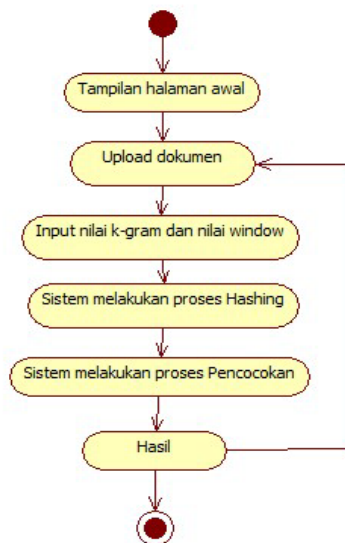
Sedangkan untuk use case diagram Text only langkah awal dimulai dari user mengcopy dokumen yang berekstensi .txt, tahap kedua input k-gram dan window lalu kita bisa melihat hasil kesamaan dokumen tersebut berupa persentase. Use case diagram Text Only dapat dilihat pada Gambar 3 berikut.



Gambar 3. Use case diagram Text Only

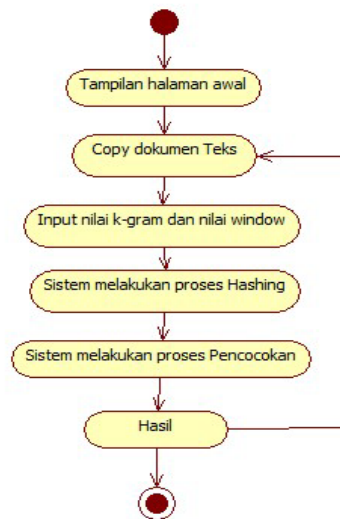
Activity Diagram

Langkah pertama dari activity diagram ini dimulai dengan mengakses halaman awal pada website. Kemudian *user* dapat menuju ke halaman selanjutnya yaitu halaman upload files, *user* dapat meng-upload 2. Lalu dokumen tersebut akan di proses. Setelah itu akan terlihat hasilnya dan akan menunjukkan berapa persentasenya. Activity Diagram Upload Files dapat dilihat pada Gambar 4 berikut.



Gambar 4. Activity Diagram Upload Files

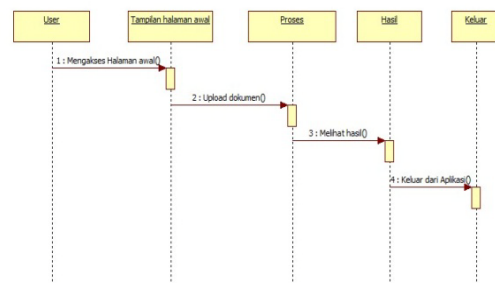
Sedangkan untuk activity diagram Text Only adalah *user* dapat meng-copy 2 dokumen. Lalu dokumen tersebut akan di proses. Setelah itu akan terlihat hasilnya dan akan menunjukkan berapa persentasenya. Activity diagram Text Only dapat dilihat pada Gambar 5 berikut.



Gambar 5. Activity Diagram Text Only

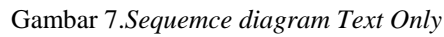
Sequence Diagram

Tahap pertama dari sequence diagram adalah *User* dapat mengakses halaman upload files untuk meng-upload dokumen dan melihat hasilnya dari *finger printing* lalu keluar dari aplikasi. Sequence diagram upload files dapat dilihat pada Gambar 6 berikut.



Gambar 6. Sequence diagram Upload Files

Sedangkan Sequence diagram untuk Text Only adalah *User* dapat mengakses halaman Text Only untuk meng-copy dokumen dan melihat hasilnya dari *finger printing* lalu keluar dari aplikasi. Sequence diagram text only dapat dilihat pada Gambar 7 berikut.



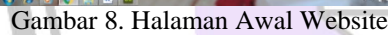
Pada tampilan di Text only ada kolom kosong yang akan di isi teks. Selanjutnya adagram dan window. Lalu ada juga tombol Cek. Halaman text only dapat dilihat pada Gambar 10 berikut.



HASIL DAN PEMBAHASAN

Tampilan Aplikasi

Tampilan awal terdiri dari sekilas tentang pengertian Algoritma Wnnowing. Halaman awal website dapat dilihat pada Gambar 8 berikut.



Tampilan selanjutnya adalah tampilan keterangan teks 1 dan teks 2 yang sudah di uji. Tampilan keterangan ini dibuat menggunakan table yaitu `table` `table-bordered`. Halamanketerangan dapat dilihat pada Gambar 11 berikut.

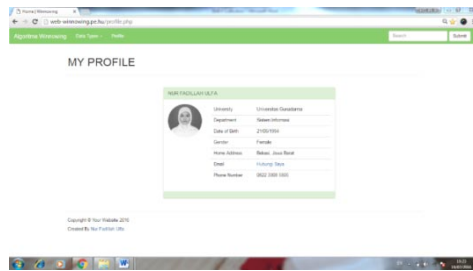
| Keterangan | Ts1 | Ts2 |
|-----------------|---|---|
| Burgulay Gum | burgulay ugung gungulay gungulay amulay amulay awulay awulay amulay amulay | burgulay ugung gungulay gungulay amulay amulay awulay awulay amulay amulay |
| Nila wani | 802055 1873246 1706305 161525 1678289 1742556 1851965 1084260 1085225 1627225 1741255 | 802055 1873246 1706305 161525 1678289 1742556 1851965 1084260 1085665 1835354 1810154 |
| Wimaw | (802055 1873246 1706305 161525 1678289) (1873246 1706305 161525 1678289 1742556) (1706305 161525 1678289 1742556 1851965) (161525 1678289 1742556 1851965 1084260) (1742556 1851965 1084260 1085225 1627225) (1085225 1627225 1627225 1627225 1742625) | (802055 1873246 1706305 161525 1678289) (1873246 1706305 161525 1678289 1742556) (1706305 161525 1678289 1742556 1851965) (161525 1678289 1742556 1851965 1084260) (1742556 1851965 1084260 1085225 1627225) (1085225 1627225 1627225 1627225 1810154) |
| Fagugugulay | {1873246, 1706305, 161525} | {1873246, 1706305, 161525} |
| Kesamar | | (1873246) = 33.3% |

Created By: Nur Fadhila Ulita

Setelah pembuatan tampilan awal website, selanjutnya membuat tampilan untuk upload dokumen dan text only. Pada tampilan kedua ini terdapat penjelasan winnowing, kemudian adatulisan file-1 dan file-2serta tombol browse Selanjutnyaada gram dan juga window. Lalu tombol Cek. Halaman upload files dapat dilihat pada Gambar 9 berikut.



Tampilan selanjutnya adalah tampilan My Profile. Tampilan halaman My profile ini terdapat penjelasan tentang data diri pembuat website. Terdiri dari tabel dan kolom-kolom berisi nama, alamat, universitas, tanggal lahir, jenis kelamin, email, dan nomor telephone. Halaman my profile dapat dilihat pada Gambar 12 berikut.



Gambar 12. Halaman *My Profile*.

Perbandingan Eksekusi Dokumen dan Waktu

Berikut ini adalah tabel hasil uji coba perbandingan ukuran dokumen dan berapa lama waktu yang digunakan.

Tabel 1 Perbandingan Eksekusi Dokumen dan Waktu

| Ukuran Dok 1 | Ukuran Dok 2 | K-gram | Window | Waktu |
|--------------|--------------|--------|--------|------------------|
| 1,18 MB | 1,12 MB | 4 | 4 | 14 detik |
| 2,48 MB | 2,30 MB | 5 | 5 | 19 detik |
| 3,51 MB | 3,08 MB | 6 | 6 | 25 detik |
| 4,61 MB | 4,31 MB | 7 | 7 | 42 detik |
| 5,90 MB | 5,74 MB | 8 | 8 | 2 menit 10 detik |

KESIMPULAN DAN SARAN

Dari hasil uji coba yang telah dilakukan yaitu pengecekan kemiripan dokumen menggunakan Algoritma Winnowing maka dapat diambil kesimpulan bahwa aplikasi website ini akan memberikan hasil berupa persentase dan keterangan bahwa kedua dokumen yang diuji termasuk plagiat atau tidak. Aplikasi ini memiliki 5 kategori persentase kemiripannya sesuai dengan teori yang ada yaitu tidak plagiat, sedikit kesamaan, plagiat tingkat sedang, mendekati plagiarisme dan plagiarisme. Waktu proses untuk

pendeteksian ini lebih lama jika memproses file yang cukup besar.

DAFTAR PUSTAKA

- [1] Ridhatillah, Ardini 2003, *Dealing with Plagiarism in the Information System Research Community: A Look at Factors That Drive Plagiarism and Ways to Address Them*, MIS Quarterly; Vol. 27, No. 4, pp. 511-532.
- [2] Kurniawati, A., Sekarwati, Kemal A., dan Wicaksana, I wayan Simri. 2012. *Arsitektur Untuk Aplikasi Deteksi Kesamaan Dokumen Bahasa Indonesia*. Konferensi Nasional Sistem Informasi 2012, STMIK - STIKOM Bali 23-25 Pebruari 2012. pp. 297-302.
- [3] Riyan Pratama, Mudafiq. 2013 *Penerapan Teknik Document Fingerprinting Pada Sistem Pendeteksi Plagiarisme eDokumentasi Terkelompok Menggunakan Algoritma Winnowing Dengan Metode K-Gram*. Jurusan Teknik Informatika Universitas Muhammadiyah Malang.
- [4] Schleimer, Saul, Wilkerson, Daniel s., Aiken, Alex. 2003. *Winnowing: Local Algorithms for Document Fingerprinting*. International conference on management of data. Proceedings of the 2003 ACM SIGMOD international conference on Management of data . pp. 76-85