

***CLUSTERING DOKUMEN TEKS BERDASARKAN  
FINGERPRINT BIWORD WINNOWING DENGAN  
MENGUNAKAN METODE K-MEANS***

**TUGAS AKHIR**

Diajukan Sebagai Salah Satu  
Syarat Untuk Memperoleh Gelar Sarjana Teknik  
Pada Jurusan Teknik Informatika

Oleh

**YULISKA**  
**10951007997**



**FAKULTAS SAINS DAN TEKNOLOGI  
UNIVERSITAS ISLAM NEGERI SULTAN SYARIF KASIM RIAU  
PEKANBARU  
2014**

**LEMBAR PENGESAHAN**  
**CLUSTERING DOKUMEN TEKS BERDASARKAN**  
**FINGERPRINT BIWORD WINNOWING DENGAN**  
**MENGGUNAKAN METODE K-MEANS**

**TUGAS AKHIR**

Oleh:

**YULISKA**  
**10951007997**

Telah dipertahankan di depan sidang dewan penguji  
Sebagai salah satu syarat untuk memperoleh gelar sarjana Teknik Informatika  
Fakultas Sains dan Teknologi Universitas Islam Negeri Sultan Syarif Kasim Riau  
di Pekanbaru, pada tanggal 26 Juni 2014

Pekanbaru, 26 Juni 2014

Mengesahkan

Ketua Jurusan



Dekan

**Yenita Morena, M.Si**  
**NIP. 19601125 198503 2 002**



**Elin Haerani, S.T, M.Kom**  
**NIP. 19810521 200710 2 003**

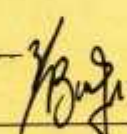
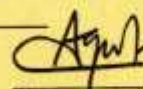
**DEWAN PENGUJI :**

**Ketua : Surya Agustian, S.T, M.Kom**

**Sekretaris : Surya Agustian, S.T, M.Kom**

**Anggota I : Jasril, S.Si, M.Sc**

**Anggota II : Elvia Budianita, S.T, M.Cs**



# ***CLUSTERING DOKUMEN TEKS BERDASARKAN FINGERPRINT BIWORD WINNOWING DENGAN MENGGUNAKAN METODE K-MEANS***

**YULISKA  
10951007997**

Tanggal Sidang : 26 Juni 2014

Periode Wisuda :

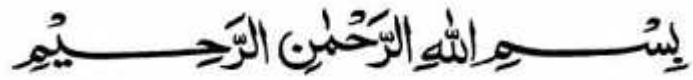
Jurusan Teknik Informatika  
Fakultas Sains Dan Teknologi  
Univesitas Islam Negeri Sultan Syarif Kasim Riau

## **ABSTRAK**

Kebutuhan akan informasi yang semakin meningkat menyebabkan informasi yang berbentuk dokumen teks mengalami penumpukan, sehingga proses pencarian menjadi sulit. Untuk itu, dilakukan suatu teknik *clustering* agar dokumen-dokumen yang memiliki tingkat kesamaan tinggi dapat berada pada satu kluster/kelompok yang sama, sehingga proses pencarian menjadi lebih mudah. Salah satu metode yang sering digunakan dalam kasus *clustering* adalah *K-Means*. Penelitian ini mencoba melakukan *clustering* terhadap dokumen teks dengan menggunakan metode *K-Means* berdasarkan *fingerprint biword winnowing* sebagai *feature* atau ciri dokumen teks, dimana *biword winnowing* adalah algoritma yang biasa dipakai untuk mendeteksi kesamaan isi suatu dokumen teks. Proses *clustering* diawali dengan mengekstrak *fingerprint* dari masing-masing dokumen yang ada pada koleksi dokumen dengan menggunakan algoritma *biword winnowing*. Kemudian dilakukan proses pembentukan dimensi dokumen yang merepresentasikan frekuensi kemunculan *fingerprint* dari masing-masing dokumen tersebut, dan dilanjutkan dengan proses *clustering* dokumen menggunakan metode *K-Means*. Selanjutnya, untuk mengetahui kualitas hasil pengelompokan dilakukan pengujian dengan menggunakan parameter *precision*. Dari keseluruhan hasil pengujian, kualitas hasil pengelompokan terhadap koleksi dokumen uji oleh aplikasi *clustering* dokumen teks ini dinilai baik dengan nilai *precision* 60%-88.33%.

**Kata Kunci :** *Biword Winnowing, Clustering, Fingerprint, K-Means, Precision.*

## KATA PENGANTAR



Segala puji bagi Allah, Pencipta mausia dan alam semesta karena senantiasa melimpahkan rahmat dan karunianya yang tidak akan sanggup untuk penulis hitung, sehingga penulis dapat menyelesaikan laporan tugas akhir perkuliahan ini dengan baik. Shalawat dan salam rindu yang terdalam semoga senantiasa tercurah dan tersampaikan kepada Nabi Muhammad, keluarganya, para sahabatnya dan para pengikutnya yang setia hingga akhir zaman.

Tugas akhir dengan judul “*Clustering* dokumen teks berdasarkan *fingerprint biword winnowing* dengan menggunakan metode *K-Means*” ini penulis susun dengan maksud untuk memenuhi tugas akhir perkuliahan yang dibebankan kepada penulis di samping untuk menambah dan memperdalam kemampuan penulis dalam memahami materi perkuliahan.

Dalam kesempatan ini penulis juga ingin mengucapkan terima kasih kepada seluruh pihak yang telah membantu penulis dalam menyelesaikan laporan tugas akhir ini. Untuk itu penulis mengucapkan terima kasih yang tak terhingga kepada :

1. Allah SWT yang selalu melimpahkan rahmat dan nikmat-Nya kepada penulis hingga terselesaikannya Tugas Akhir ini.
2. Bapak Prof. DR. H. Munzir Hitami, MA, selaku Rektor Universitas Islam Negeri Sultan Syarif Kasim Riau.
3. Ibu Dra. Hj.Yenita Morena, M.Si, selaku Dekan Fakultas Sains dan Teknologi, Universitas Islam Negeri Sultan Syarif Kasim Riau
4. Ibu Maylisda dan Ayah Zakaria tercinta, yang telah mendo'akan dan memberikan dukungan yang sangat luar biasa kepada penulis dalam menyelesaikan tugas akhir ini. Semoga beliau selalu dalam lindungan Allah SWT serta segala ketulusan dan pengorbanan beliau diridhoi oleh Allah SWT. Amin.
5. Ibu Elin Haerani, S.T, M.Kom, selaku Ketua Jurusan Teknik Informatika.

6. Bapak Novriyanto, S.T, M.Sc, selaku Pembimbing Akademis Penulis. Terima kasih untuk segala bimbingan dan masukan yang telah diberikan kepada penulis.
7. Bapak Surya Agustian, S.T, M.Kom, selaku Pembimbing tugas akhir Penulis yang selalu sabar dan meluangkan waktu untuk membimbing, memberikan saran dan kritik dalam penyusunan tugas akhir ini.
8. Bapak Jasril, S.Si, M.Sc, selaku Penguji I tugas akhir. Terima kasih untuk segala masukan yang telah diberikan kepada penulis demi kesempurnaan laporan tugas akhir ini.
9. Ibu Elvia Budianita, S.T, M.Cs, Selaku Penguji II tugas akhir. Terima kasih untuk segala masukan yang telah diberikan kepada penulis demi kesempurnaan laporan tugas akhir ini.
10. Bapak Muhammad Affandes, S.T, M.T, selaku Koordinator Tugas akhir Jurusan Teknik Informatika.
11. Adik-adikku Rully Dwi Andika dan Febila Anzari yang telah memberikan dorongan dan semangat dalam menyelesaikan tugas akhir ini.
12. Seluruh keluarga besar Penulis, yang selalu memberikan semangat untuk menyelesaikan tugas akhir ini.
13. Seluruh dosen dan staf Fakultas Sains dan Teknologi khususnya pada Jurusan Teknik Informatika. Terima kasih atas ilmu yang telah diberikan.
14. Teman-teman seperjuangan penulis di UIN Suska Riau, teman-teman TIF 2009, Khususnya TIF C. Semoga yang telah lulus mendapat apa yang dicita-citakan dan yang belum lulus diberi kemudahan oleh Allah untuk menyusul teman-teman yang telah lulus. Amin.
15. Seluruh sahabat penulis, terima kasih untuk segala dukungannya.
16. Seluruh pihak yang belum penulis cantumkan, terima kasih atas dukungannya kepada penulis.

Selanjutnya, penulis menyadari bahwa dalam penulisan laporan tugas akhir ini masih terdapat kekurangan. Oleh karena itu, penulis mengharapkan kritik dan saran yang membangun demi terciptanya laporan yang lebih baik di masa yang akan datang melalui email penulis [yuliska\\_mz@yahoo.com](mailto:yuliska_mz@yahoo.com) atau [hatake.bee@gmail.com](mailto:hatake.bee@gmail.com). Penulis berharap laporan ini dapat bermanfaat bagi para

pembaca, baik dari pihak UIN SUSKA RIAU maupun dari kalangan akademis serta masyarakat umum.

Pekanbaru, Juni 2014

Penulis

# DAFTAR ISI

	Halaman
HALAMAN JUDUL LAPORAN.....	i
LEMBAR PERSETUJUAN.....	ii
LEMBAR PENGESAHAN .....	iii
LEMBAR HAK ATAS KEKAYAAN INTELEKTUAL.....	iv
LEMBAR PERNYATAAN .....	v
LEMBAR PERSEMBAHAN .....	vi
ABSTRAK .....	vii
<i>ABSTRACT</i> .....	viii
KATA PENGANTAR .....	ix
DAFTAR ISI.....	xii
DAFTAR GAMBAR .....	xvi
DAFTAR TABEL.....	xviii
DAFTAR LAMPIRAN.....	xx
DAFTAR RUMUS .....	xxi
DAFTAR SIMBOL.....	xxii
BAB I PENDAHULUAN .....	I-1
1.1 Latar Belakang .....	I-1
1.2 Rumusan Masalah .....	I-2
1.3 Batasan Masalah. ....	I-2
1.4 Tujuan Penelitian .....	I-3
1.5 Sistematika Penulisan .....	I-3
BAB II LANDASAN TEORI .....	II-1
2.1 <i>Text Mining</i> .....	II-1
2.1.1 <i>Text Preprocessing</i> .....	II-1
2.1.2 <i>Text Transformation</i> .....	II-2
2.1.2.1 Algoritma <i>Winnowing</i> .....	II-2
2.1.2.2 Algoritma <i>Biword Winnowing</i> .....	II-5
2.1.2.3 <i>Stemming</i> .....	II-11
2.1.3 <i>Feature Selection</i> .....	II-12

2.1.3.1 <i>Document Representation</i> .....	II-12
2.1.3.2 Fungsi Jarak Dan Similaritas Antar Dokumen .....	II-14
2.1.4 <i>Data Mining</i> .....	II-15
2.1.4.1 <i>Clustering</i> Dokumen Dan Metodenya.....	II-15
2.1.4.2 <i>Hierarchical Agglomerative Clustering</i> .....	II-16
2.1.4.3 <i>Divisive Clustering</i> .....	II-19
2.1.4.4 <i>K-Means</i> .....	II-19
2.1.5 <i>Evaluation</i> .....	II-20
BAB III METODOLOGI PENELITIAN.....	III-1
3.1 Identifikasi Masalah.....	III-2
3.2 Pengumpulan Data.....	III-2
3.3 Analisa Aplikasi.....	III-3
3.4 Perancangan Aplikasi .....	III-4
3.5 Implementasi.....	III-5
3.5 Pengujian .....	III-6
3.6 Kesimpulan dan Saran .....	III-6
BAB IV ANALISA DAN PERANCANGAN.....	IV-1
4.1 Analisa Kebutuhan Data .....	IV-1
4.2 Penyelesaian Masalah.....	IV-2
4.2.1 <i>Text Preprocessing</i> .....	IV-2
4.2.2 <i>Text Transformation</i> .....	IV-4
4.2.2.1 Proses <i>Tokenizing</i> .....	IV-4
4.2.2.2 Proses Enkripsi MD5.....	IV-6
4.2.2.3 Proses <i>Hashing</i> .....	IV-9
4.2.2.4 Pembentukan Window.....	IV-12
4.2.2.5 Pembentukan <i>Fingerprint</i> .....	IV-15
4.2.3 <i>Feature Selection</i> .....	IV-16
4.2.3.1 Penggabungan Dan Pengurutan <i>Fingerprint</i> ...	IV-16
4.2.3.2 Pembentukan Dimensi.....	IV-17
4.2.3.3 Reduksi Dimensi .....	IV-18
4.2.4 <i>Data Mining</i> .....	IV-20



4.2.4.1 Pembentukan <i>Centroid</i> .....	IV-20
4.2.4.2 Perhitungan <i>Similarity</i> .....	IV-22
4.2.4.3 Pengelompokkan .....	IV-23
4.2.4.4 Pembentukan <i>Centroid</i> Baru .....	IV-24
4.3 Perancangan Aplikasi .....	IV-25
4.3.1 Perancangan Struktur Menu .....	IV-25
4.3.2 Perancangan <i>Pseudo Code</i> .....	IV-26
4.3.2.1 <i>Pseudo Code Text Preprocessing Dan Text Transformation</i> .....	IV-26
4.3.2.2 <i>Pseudo Code Feature Selection</i> .....	IV-27
4.3.2.3 <i>Pseudo Code Data Mining</i> .....	IV-28
4.3.3 Perancangan <i>Interface</i> .....	IV-29
4.3.3.1 Rancangan Menu Halaman Utama ( <i>Home</i> )....	IV-30
4.3.3.2 Rancangan Menu <i>Form</i> Kluster .....	IV-30
4.3.3.3 Rancangan Menu <i>Text Preprocessing</i> .....	IV-31
4.3.3.4 Rancangan Menu <i>Text Transformation</i> .....	IV-32
4.3.3.5 Rancangan Menu <i>Feature Selection</i> .....	IV-32
4.3.3.6 Rancangan Menu <i>Data Mining</i> .....	IV-33
4.3.3.6 Rancangan Menu <i>Help</i> .....	IV-34
BAB V IMPLEMENTASI DAN PENGUJIAN .....	V-1
5.1 Tahapan Implementasi.....	V-1
5.1.1 Batasan Implementasi.....	V-1
5.1.2 Lingkungan Implementasi .....	V-1
5.1.3 Implementasi <i>Interface</i> Aplikasi .....	V-2
5.2 Pengujian Aplikasi.....	V-10
5.2.1 Rencana Pengujian .....	V-10
5.2.1.1 Pengujian Hasil Pengelompokan Dengan Pembentukan <i>Centroid</i> Cara I .....	V-12
5.2.1.2 Pengujian Hasil Pengelompokan Dengan Pembentukan <i>Centroid</i> Cara I .....	V-21
5.2.4 Kesimpulan Pengujian.....	V-30
BAB VI PENUTUP .....	VI-1

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

DAFTAR PUSTAKA

LAMPIRAN

DAFTAR RIWAYAT HIDUP