

**APLIKASI MESIN PENCARI DOKUMEN *CROSS LANGUAGE*
BAHASA INGGRIS – BAHASA INDONESIA MENGGUNAKAN
*VECTOR SPACE MODEL***



SKRIPSI

**Disusun Sebagai Salah Satu Syarat
untuk Memperoleh Gelar Sarjana Komputer
pada Jurusan Ilmu Komputer/ Informatika**

Disusun oleh:

ABDUL MAGHFIR ZAKIR

24010311130054

**JURUSAN ILMU KOMPUTER/INFORMATIKA
FAKULTAS SAINS DAN MATEMATIKA
UNIVERSITAS DIPONEGORO**

2015

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Saya yang bertanda tangan di bawah ini :

Nama : Abdul Maghfir Zakir

NIM : 24010311130054

Judul : Aplikasi Mesin Pencari Dokumen *Cross Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*.

Dengan ini saya menyatakan bahwa dalam tugas akhir/ skripsi ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar kesarjanaan di suatu Perguruan Tinggi, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan di dalam daftar pustaka.



HALAMAN PENGESAHAN

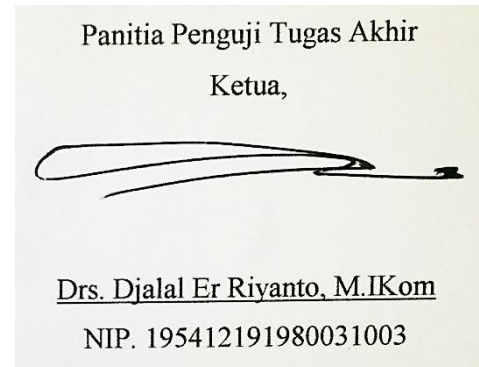
Judul : Aplikasi Mesin Pencari Dokumen *Cross Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*.

Nama : Abdul Maghfir Zakir

NIM : 24010311130054

Telah diujikan pada sidang tugas akhir tanggal 22 Desember 2015 dan dinyatakan lulus pada tanggal **14 Januari 2016**.

Semarang, Januari 2016




HALAMAN PENGESAHAN

Judul : Aplikasi Mesin Pencari Dokumen *Cross Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*.

Nama : Abdul Magfir Zakir

NIM : 24010311130054

Telah diujikan pada sidang tugas akhir tanggal 22 Desember 2015.

Semarang, Januari 2016
Pembimbing

Sukmawati Nur Endah, S.Si, M.Kom
NIP. 197805022005012002

ABSTRAK

Fenomena banjirnya informasi seiring berkembangnya teknologi dapat ditangani dengan kajian-kajian yang ada pada *information retrieval* (IR). Namun permasalahan pada IR semakin berkembang karena adanya perbedaan bahasa antar negara. Oleh karena itu dibutuhkan suatu aplikasi *cross-language information retrieval* (CLIR) yang mampu mengoptimalkan hasil pencarian dokumen lintas bahasa. Aplikasi mesin pencari dokumen ini dapat melakukan penerjemahan pada masukan *query* lalu dihitung kemiripannya terhadap dokumen yang memiliki bahasa yang berbeda. Nilai kemiripan ini dapat ditentukan berdasarkan nilai similaritas dokumen terhadap *query* dengan *Vector Space Model*. Sepuluh dokumen dengan nilai kemiripan tertinggi akan diambil dan dihitung nilai ketepatannya (*precision*). Berdasarkan hasil evaluasi dengan menggunakan 8 *user* dan 10 *query*, didapatkan nilai *Mean Average Precision* (MAP) keseluruhan sebesar 0.912 atau sekitar 91%.

Kata kunci : mesin pencari dokumen, CLIR, *vector space model*

ABSTRACT

The information flooding phenomenon as technology growth can be handled by existing studies on information retrieval (IR). However, the problem in IR keep on growing because of language differences between countries. Therefore, it needs a cross-language information retrieval (CLIR) application which is able to optimize cross-language document search result. This document search engine application can translate the query input, then calculate its similarity to the document that has a different language. The similarity value can be determined based on the similarity of the document to the query with Vector Space Model. Ten documents with the highest similarity value will be taken and its value of accuracy (precision) will be calculated. Based on the evaluation results by using 8 users and 10 queries, it was discovered that its overall Mean Average Precision (MAP) is 0.912 or about 91%.

Keyword: document search engine, CLIR, *vector space model*

KATA PENGANTAR

Segala puji syukur bagi Tuhan Yang Maha Esa atas karunia-Nya yang diberikan kepada penulis sehingga penulis dapat menyelesaikan penulisan laporan tugas akhir yang berjudul “Aplikasi Mesin Pencari Dokumen *Cross Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*”. Laporan tugas akhir ini disusun sebagai salah satu syarat untuk memperoleh gelar sarjana strata satu pada Jurusan Ilmu Komputer/ Informatika Fakultas Sains dan Matematika Universitas Diponegoro Semarang.

Dalam penyusunan laporan ini penulis banyak mendapat bimbingan dan bantuan dari berbagai pihak. Untuk itu, pada kesempatan ini penulis mengucapkan rasa hormat dan terima kasih kepada:

1. Ragil Saputra, S.Si, M.Cs selaku Ketua Jurusan Ilmu Komputer/Informatika
2. Helmie Arif Wibawa, S.Si, M.Cs selaku Koordinator Tugas Akhir
3. Sukmawati Nur Endah, S.Si, M.Kom selaku dosen pembimbing
4. Semua pihak yang telah membantu kelancaran dalam penyusunan tugas akhir, yang tidak dapat penulis sebutkan satu persatu.

Penulis menyadari bahwa dalam laporan ini masih banyak kekurangan baik dari segi materi ataupun dalam penyajiannya karena keterbatasan kemampuan dan pengetahuan penulis. Oleh karena itu, kritik dan saran sangat penulis harapkan. Semoga laporan ini dapat bermanfaat bagi pembaca dan penulis pada umumnya.

Semarang, Januari 2016

Penulis,

Abdul Maghfir Zakir

24010311130054

DAFTAR ISI

HALAMAN PERNYATAAN KEASLIAN SKRIPSI.....	ii
HALAMAN PENGESAHAN	iii
ABSTRAK	v
ABSTRACT	vii
KATA PENGANTAR.....	vii
DAFTAR ISI	viii
DAFTAR GAMBAR.....	x
DAFTAR TABEL	xii
DAFTAR LAMPIRAN	xiii
BAB I PENDAHULUAN	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah.....	2
1.3. Tujuan dan Manfaat	2
1.4. Ruang Lingkup.....	3
1.5. Sistematika Penulisan	3
BAB II TINJAUAN PUSTAKA.....	5
2.1. Information Retrieval.....	5
2.1.1. <i>Cross-Language</i> Information Retrieval	6
2.2. Arsitektur <i>Information Retrieval</i>	7
2.2.1. Tokenisasi	8
2.2.2. Penghapusan <i>Stopword</i>	8
2.2.3. Stemming	9
2.2.4. Pembobotan.....	11
2.3. Proses Pengembangan Perangkat Lunak	12
2.3.1. <i>System/Information Engineering</i>	12
2.3.2. <i>Analysis</i>	13
2.3.2.1. Pemodelan Fungsional.....	13
2.3.3. <i>Design</i>	17
2.3.3.1. Perancangan Struktur.....	17
2.3.3.2. Perancangan Fungsional	17
2.3.3.3. Perancangan Antarmuka.....	17
2.3.4. <i>Code</i>	17
2.3.5. <i>Test</i>	17
2.4. PHP	18

2.5. MySQL	19
2.6. <i>Vector Space Model</i>	19
2.7. Evaluasi.....	22
2.7.1. <i>Precision dan Recall</i>	23
2.7.2. <i>Average Precision</i>	24
2.7.3. <i>Mean Average Precision</i>	25
BAB III ANALISIS DAN PERANCANGAN	26
3.1. Analisis	26
3.1.1. Deskripsi Umum Aplikasi.....	26
3.1.2. Kebutuhan Fungsional dan Non Fungsional	27
3.1.3. Pemodelan Fungsional	28
3.1.3.1. Context Diagram (CD)	28
3.1.3.2. <i>Diagram Flow Diagram (DFD)</i>	29
3.2. Perancangan Sistem	30
3.2.1. Perancangan Struktur Data.....	31
3.2.2. Perancangan Antarmuka	33
3.2.3. Perancangan Fungsi	39
BAB IV IMPLEMENTASI DAN PENGUJIAN	59
4.1. Implementasi.....	59
4.1.1. Lingkungan Implementasi Sistem.....	59
4.1.2. Implementasi Data	59
4.1.3. Implementasi Antarmuka	62
4.1.4. Implementasi Fungsi	66
4.2. Pengujian.....	66
4.2.1. Rencana Pengujian	67
4.2.1.1. Rencana Pengujian Fungsional Sistem.....	67
4.2.1.2. Rencana Evaluasi <i>Vector Space Model</i>	67
4.2.2. Pengujian dan Hasil Uji	69
4.2.2.1. Pengujian dan Hasil Uji Fungsional Sistem	70
4.2.2.2. Evaluasi dan Hasil Evaluasi <i>Vector Space Model</i>	70
4.2.3. Analisa Hasil Evaluasi	72
BAB V PENUTUP	73
5.1. Kesimpulan	73
5.2. Saran	73
DAFTAR PUSTAKA.....	74
LAMPIRAN	76

DAFTAR GAMBAR

Gambar 2.3. Model <i>Waterfall</i>	12
Gambar 2.4. Cosine dari θ didapatkan dari d_j dan q (Yates R, 1999).....	20
Gambar 3.1. Gambaran Umum Aplikasi.....	26
Gambar 3.2. <i>Context Diagram</i> Aplikasi Mesin Pencari Dokumen <i>Cross Language</i> Bahasa Inggris – Bahasa Indonesia Menggunakan <i>Vector Space Model</i>	29
Gambar 3.3. DFD Aplikasi Mesin Pencari Dokumen <i>Cross-Language</i> Bahasa Inggris – Bahasa Indonesia Menggunakan <i>Vector Space Model</i>	30
Gambar 3.4. Desain Antarmuka Halaman Utama Aplikasi	34
Gambar 3.5. Desain Antarmuka Halaman Utama Aplikasi.....	34
Gambar 3.6. Desain Antarmuka Halaman <i>Log in</i>	34
Gambar 3.7. Desain Antarmuka Halaman <i>Upload</i> Dokumen	35
Gambar 3.8. Desain Antarmuka Halaman Gagal <i>Upload</i> Dokumen.....	36
Gambar 3.9. Desain Antarmuka Halaman <i>Stemming</i> Dokumen.....	36
Gambar 3.10. Desain Antarmuka Halaman <i>Stemming</i> Dokumen Berhasil	37
Gambar 3.11. Desain Antarmuka Halaman Hasil Pencarian Dokumen.....	38
Gambar 3.12. Desain Antarmuka Halaman Detail Proses Pencarian	38
Gambar 3.13. Desain Antarmuka Halaman Detail Proses Pencarian	39
Gambar 3.15. <i>Flowchart</i> Proses Fungsi-Fungsi pada <i>Admin</i>	40
Gambar 3.16. <i>Flowchart</i> Proses Fungsi-Fungsi pada <i>User</i>	54
Gambar 3.17. <i>Flowchart</i> Proses <i>Login</i> Aplikasi	41
Gambar 3.18. <i>Flowchart</i> Proses <i>Upload</i> Dokumen.....	42
Gambar 3.19. <i>Flowchart</i> Proses <i>Preprocess</i> Dokumen	43
Gambar 3.20. <i>Flowchart</i> Proses <i>Stemming</i> Dokumen.....	44
Gambar 3.21. <i>Flowchart</i> Proses Porter Stemmer.....	45
Gambar 3.22. <i>Flowchart</i> Lanjutan Proses Porter Stemmer	46
Gambar 3.23. <i>Flowchart</i> Proses Nazief dan Andriani.....	51
Gambar 3.24. <i>Flowchart</i> Pengindeksan Dokumen	52
Gambar 3.25. <i>Flowchart</i> Proses Pembobotan	53
Gambar 3.26. <i>Flowchart</i> Proses Pencarian Dokumen.....	55
Gambar 3.27. <i>Flowchart Preprocess Query</i>	56
Gambar 3.28. <i>Flowchart</i> Proses <i>Vector Space Model</i>	57
Gambar 4.1. Struktur Tabel <i>eng_ind</i> pada MySQL.....	60
Gambar 4.2. Struktur Tabel <i>ind_eng</i> pada MySQL.....	60
Gambar 4.3. Struktur Tabel <i>Admin</i> pada MySQL.....	60

Gambar 4.4.	Struktur Tabel Kata Dasar pada MySQL	61
Gambar 4.5.	Struktur Tabel <i>Stopwords_eng</i> pada MySQL	61
Gambar 4.6.	Struktur Tabel <i>Stopwords_ind</i> pada MySQL.....	61
Gambar 4.7.	Struktur Folder Korpus	61
Gambar 4.8.	Antarmuka Halaman Utama.....	62
Gambar 4.9.	Antarmuka Halaman <i>Log In</i>	63
Gambar 4.10.	Antarmuka Halaman <i>Upload</i> Dokumen	63
Gambar 4.11.	Antarmuka Gagal <i>Upload</i> Dokumen	64
Gambar 4.12.	Antarmuka <i>Stemming</i> Dokumen.....	64
Gambar 4.13.	Antarmuka <i>Stemming</i> Dokumen Berhasil	65
Gambar 4.14.	Antarmuka Hasil Pencarian Dokumen	65
Gambar 4.15.	Antarmuka Detail Proses Pencarian.....	66
Gambar 4.16.	Antarmuka Detail Proses Pencarian.....	66
Gambar 4.17.	Grafik <i>Average Precision</i>	71
Gambar 4.18.	Grafik <i>Mean Average Precision</i> Setiap <i>User</i>	71

DAFTAR TABEL

Tabel 2.1.	Contoh SRS	13
Tabel 2.2	Tabel Penomoran DFD	14
Tabel 2.3	Tabel Notasi DFD	16
Tabel 2.4.	Nilai bobot kata dari contoh <i>Vector Space Model</i>	21
Tabel 2.5.	Perhitungan <i>Recall</i> dan <i>Precision</i>	23
Tabel 3.1.	Kebutuhan Fungsional	27
Tabel 3.2.	Kebutuhan Non Fungsional.....	27
Tabel 3.3.	Struktur Tabel Kamus Inggris – Indonesia.....	31
Tabel 3.4.	Struktur Tabel Kamus Indonesia – Inggris.....	31
Tabel 3.5.	Struktur Tabel Kata Dasar	31
Tabel 3.6.	Struktur Tabel <i>Admin</i>	32
Tabel 3.7.	Struktur Tabel <i>Stopwords</i> Inggris.....	32
Tabel 3.8.	Struktur Tabel <i>Stopwords</i> Inggris.....	32
Tabel 3.9.	Struktur Folder Korpus	33
Tabel 4.1.	Rencana Pengujian Fungsional Sistem	67
Tabel 4.2.	Daftar Dokumen yang Digunakan dalam Pengujian.....	68
Tabel 4.3.	Daftar Penggunaan <i>Query</i>	69
Tabel 4.4.	Hasil Perhitungan <i>Average Precision</i>	70
Tabel 4.5.	Hasil Mean Average Precision setiap user	71
Tabel 7.1	Sampel Daftar Data <i>Stopwords</i>	76
Tabel 7.2	Sampel Daftar Data Kata Dasar	78
Tabel 7.3	Sampel Daftar Data Kamus	80
Tabel 7.4.	Deskripsi dan Hasil Uji <i>Log In Admin</i>	82
Tabel 7.5.	Deskripsi dan Hasil Uji <i>Upload</i> Dokumen.....	83
Tabel 7.6.	Deskripsi dan Hasil Uji <i>Stemming</i> Dokumen.....	84
Tabel 7.7.	Deskripsi dan Hasil Uji Pengindeksan Dokumen.....	84
Tabel 7.8.	Deskripsi dan Hasil Uji Pencarian Dokumen.....	85
Tabel 7.9.	Deskripsi dan Hasil Uji <i>Download</i> Dokumen.....	85
Tabel 7.10.	Deskripsi dan Hasil Uji Menampilkan Detail Proses Pencarian	86
Tabel 7.11.	Hasil Perhitungan <i>Precision</i> dan <i>Recall</i>	87

DAFTAR LAMPIRAN

Lampiran 1. Sampel Daftar Data <i>Stopwords</i>	76
Lampiran 2. Sampel Daftar Data Kata Dasar	78
Lampiran 3. Sampel Daftar Data Kamus.....	80
Lampiran 4. Pengujian dan Hasil Uji Fungsional Sistem.....	82
Lampiran 5. Hasil Evaluasi <i>Vector Space Model</i>	87
Lampiran 6. <i>Source Code</i> Fungsi Pencarian Dokumen.....	91

BAB I

PENDAHULUAN

1.1. Latar Belakang

Jumlah informasi yang tersedia di internet saat ini sangat banyak dan terus bertambah setiap saat. Hal ini dikarenakan perkembangan teknologi yang dimanfaatkan sangat membantu dalam penyebaran informasi. Seiring dengan peningkatan jumlah dan keragaman informasi yang beredar, pengguna semakin sulit mendapatkan informasi yang diinginkan. Kebutuhan penggunapun mulai bergeser dari yang dulunya mencari informasi secara kuantitatif menjadi kualitatif. Salah satu bidang yang membahas tentang pencarian informasi yaitu *information retrieval* (IR). Fenomena banjirnya informasi dapat diselesaikan dengan kajian-kajian yang ada pada IR agar pengguna dapat menemukan informasi yang sesuai keinginan.

Permasalahan pada IR semakin berkembang karena adanya perbedaan bahasa dari berbagai negara. Misalnya dilakukan pencarian informasi dengan memasukkan *query* dengan bahasa tertentu, maka hasil yang diterima pengguna hanya dokumen berisi informasi yang ditulis dengan bahasa tersebut. Sehingga hasil pencarian tidak dapat memberikan hasil yang maksimum untuk pengguna.

Peningkatan hasil pencarian agar lebih optimal pada *search engine* ini dapat diselesaikan dengan *cross-language information retrieval* (CLIR). CLIR pada dasarnya memiliki metode yang sama dengan IR, hanya saja CLIR memiliki tahap penerjemahan pada *query* atau pada dokumen korpus. Kedua metode penerjemahan tersebut memiliki kekurangan masing-masing. Metode dengan menerjemahkan dokumen korpus akan memiliki proses yang lama dan tidak praktis sedangkan untuk menerjemahkan *query* menjadi sulit pada saat *query* bersifat ambigu.

Penerapan CLIR ini terbagi lagi menjadi beberapa model, antara lain *Probabilistic Model*, *Set-theoretic Model* dan *Algebraic Model*. *Probabilistic Model* contohnya pada penerapan *Teorema Bayes*, penerapan *Set-theoretic Models* pada *Standard Boolean* dan *Extended*, dan *Algebraic Model* contohnya adalah *Vector Space Model*. Dari tiga model yang telah disebutkan, *Algebraic Model* dengan contoh *Vector Space Model* adalah model yang paling sederhana dalam pencarian kata. *Vector Space Model* telah terbukti memiliki efektifitas dalam pencarian kata

dengan menampilkan hasil pencariannya berdasar kemiripan *vector query* dan *vector* dokumen (Bari & Saputra, 2011).

Penelitian terdahulu yang menggunakan metode ini antara lain Fatkhul Amin dalam publikasinya yang berjudul Sistem Temu Kembali Informasi dengan Pemeringkatan Metode *Vector Space Model* (Amin, 2013). Penelitian ini membahas tentang mesin pencari dokumen menggunakan *Vector Space Model* dengan hasil evaluasi presisi mencapai 99%. Selain itu dari publikasi Andre Hesel, Ricky Dwiputra dan Hendra yang berjudul Studi dan Evaluasi Kinerja Model-Model *Information Retrieval* Berbasis Dokumen Teks menyimpulkan bahwa dalam memilih model *information retrieval* sebaiknya menggunakan *Vector Space Model*. Hal ini dikarenakan pada penelitian diketahui bahwa *Vector Space Model* memberikan *Mean Average Precision* yang lebih besar yaitu 71,25% dibanding dengan *Boolean Model* dan *Latent Semantic Indexing (LSI)* (Hesel, et al., 2013).

Berdasarkan uraian diatas maka akan dibangun sebuah aplikasi mesin pencari dokumen *cross language* Bahasa Inggris – Bahasa Indonesia Menggunakan *vector space model*.

1.2. Rumusan Masalah

Berdasarkan latar belakang yang telah diuraikan, maka dapat dirumuskan masalah yaitu bagaimana membuat aplikasi mesin pencari dokumen *cross language* Bahasa Inggris – Bahasa Indonesia Menggunakan *vector space model*.

1.3. Tujuan dan Manfaat

Penelitian Tugas Akhir ini bertujuan untuk menghasilkan suatu aplikasi yang dapat mencari dokumen dengan kemampuan *cross-language* Bahasa Inggris – Bahasa Indonesia Menggunakan *vector space model*.

Aplikasi mesin pencari ini diharapkan bisa bermanfaat dalam menemukan dokumen yang relevan dengan Menggunakan *query* Bahasa Inggris untuk menemukan dokumen yang relevan dalam Bahasa Indonesia sekalipun, begitupun sebaliknya.

1.4. Ruang Lingkup

Ruang lingkup dari pengembangan Aplikasi Mesin Pencari Dokumen *Cross Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model* adalah sebagai berikut:

1. Penerapan CLIR yang digunakan adalah penerjemahan pada *query*.
2. Bahasa yang digunakan adalah Bahasa Inggris dan Bahasa Indonesia.
3. Dokumen diolah berasal dari korpus lokal, bukan menggunakan *web crawling*.
4. Dokumen yang diolah menggunakan korpus menggunakan format file txt.
5. Dokumen yang ter-*retrieve* dapat diunduh dengan format pdf.
6. Evaluasi metode menggunakan minimal 20 dokumen jurnal Bahasa Indonesia dan 20 dokumen jurnal Bahasa Inggris.
7. Pengembangan aplikasi menggunakan model proses pengembangan perangkat lunak waterfall.
8. Pembangunan aplikasi menggunakan Bahasa Pemrograman PHP dan *Database Management System* MySQL.

1.5. Sistematika Penulisan

Sistematika penulisan yang digunakan dalam tugas akhir ini dibagi menjadi dalam beberapa pokok bahasan, yaitu:

BAB I PENDAHULUAN

Bab ini membahas latar belakang, rumusan masalah, tujuan dan manfaat, ruang lingkup dan sistematika penulisan dalam pembuatan tugas akhir mengenai Aplikasi Mesin Pencari Dokumen *Cross-Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*.

BAB II TINJAUAN PUSTAKA

Bab ini menyajikan tinjauan pustaka yang berhubungan dengan topik tugas akhir. Tinjauan pustaka yang digunakan dalam penyusunan tugas akhir ini meliputi *information retrieval*, arsitektur *information retrieval*, proses pengembangan perangkat lunak, PHP, MySQL, *Vector Space Model*, dan evaluasi.

BAB III ANALISIS DAN PERANCANGAN

Bab ini menyajikan mengenai pembahasan tahapan dari model pengembangan perangkat lunak waterfall yang meliputi tahap analisis dan perancangan dari Aplikasi Mesin Pencari Dokumen *Cross-Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*.

BAB IV IMPLEMENTASI DAN PENGUJIAN

Bab ini menyajikan mengenai pembahasan tahapan dari model pengembangan perangkat lunak waterfall yang meliputi tahap implementasi dan pengujian dari Aplikasi Mesin Pencari Dokumen *Cross-Language* Bahasa Inggris – Bahasa Indonesia Menggunakan *Vector Space Model*.

BAB V PENUTUP

Bab ini merupakan kesimpulan dari bab-bab yang telah dibahas sebelumnya serta saran untuk pengembangan penelitian tugas akhir lebih lanjut.