



Tesis - KI142502

**PERBAIKAN METODE REKOMENDASI DISKUSI  
PEMROGRAMAN DENGAN NORMALISASI  
IDENTIFIER MENGGUNAKAN LINGUA::IDSPLITTER**

Nanang Fakhur Rozi  
5112201045

DOSEN PEMBIMBING

Daniel Oranova Siahaan, S.Kom., M.Sc., P.D.Eng.  
Fajar Baskoro, S.Kom., M.T.

PROGRAM MAGISTER

BIDANG KEAHLIAN REKAYASA PERANGKAT LUNAK

JURUSAN TEKNIK INFORMATIKA

FAKULTAS TEKNOLOGI INFORMASI

INSTITUT TEKNOLOGI SEPULUH NOPEMBER

SURABAYA

2016



Thesis - KI142502

**IMPROVEMENT OF PROGRAMMING DISCUSSION  
RECOMMENDER METHOD BASED ON  
LINGUA::IDSPLITTER IDENTIFIER NORMALIZATION**

Nanang Fakhur Rozi  
5112201045

SUPERVISOR

Daniel Oranova Siahaan, S.Kom., M.Sc., P.D.Eng.  
Fajar Baskoro, S.Kom., M.T.

MAGISTERPROGRAM

THE EXPERTISE FIELD OF SOFTWARE ENGINEERING  
DEPARTMENT OF INFORMATICS  
FACULTY OF INFORMATION TECHNOLOGY  
INSTITUT TEKNOLOGI SEPULUH NOPEMBER  
SURABAYA  
2016

## LEMBAR PENGESAHAN

Tesis disusun untuk memenuhi salah satu syarat memperoleh gelar  
Magister Komputer (M.Kom.)

di

Institut Teknologi Sepuluh Nopember Surabaya

oleh:

Nanang Fakhur Rozi

NRP. 5112201045

Dengan judul :

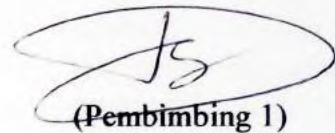
Perbaikan Metode Rekomendasi Diskusi Pemrograman dengan  
Normalisasi *Identifier* Menggunakan Lingua::IdSplitter

Tanggal Ujian : 24-6-2016

Periode Wisuda : 2015 Genap

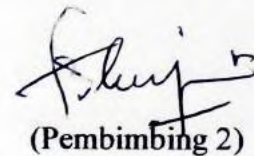
Disetujui oleh:

Daniel Oranova Siahaan, S.Kom., M.Sc., PD.Eng.  
NIP. 197411232006041001



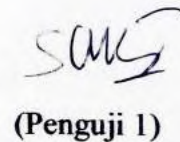
(Pembimbing 1)

Fajar Baskoro, S.Kom., M.T.  
NIP. 197404031999031002



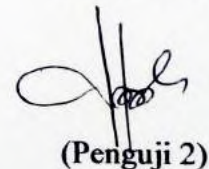
(Pembimbing 2)

Dr. Ir. Siti Rochimah, M.T.  
NIP. 196810021994032001



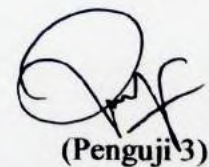
(Penguji 1)

Sarwosri, S.Kom., M.T.  
NIP. 197608092001122001



(Penguji 2)

Rizky Januar Akbar, S.Kom., M.Eng.  
NIP. 198701032014041001



(Penguji 3)



Direktur Program Pasca Sarjana,

Prof. Dr. Djauhar Manfaat, M.Sc., Ph.D.  
NIP. 196012021987011001

*Halaman ini sengaja dikosongkan*

# PERBAIKAN METODE REKOMENDASI DISKUSI PEMROGRAMAN DENGAN NORMALISASI *IDENTIFIER* MENGUNAKAN LINGUA::IDSPLITTER

Nama : Nanang Fakhrrur Rozi  
NRP : 5112201045  
Pembimbing I : Daniel Oranova Siahaan, S.Kom., M.Sc., P.D.Eng.  
Pembimbing II : Fajar Baskoro, S.Kom., M.T.

## ABSTRAK

Situs tanya-jawab Stack Overflow telah sering digunakan sebagai acuan oleh *programmer*. Informasi atau solusi dalam proses pengembangan perangkat lunak dapat dicari dengan bantuan mesin pencari pada situs. Namun, perbedaan dalam gaya penulisan, terutama pada penulisan *identifier* program, sering menyebabkan rekomendasi (pencarian) menjadi tidak sesuai dengan kebutuhan *programmer*. Beberapa *programmer* menulis *identifier* dalam bentuk singkatan sementara yang lain tidak sehingga menurunkan kinerja rekomendasi. Penelitian ini mengadopsi Lingua::IdSplitter untuk menormalkan *identifier* pada data diskusi Stack Overflow. Proses normalisasi dilakukan dengan memisahkan *identifier* yang terdiri atas komposisi term serta memperluas singkatan yang ada pada *identifier* ke bentuk penuh. Hasil penelitian menunjukkan bahwa normalisasi *identifier* menggunakan Lingua::IdSplitter hanya mampu meningkatkan performa sistem rekomendasi ketika *identifier* dengan unsur singkatan mendominasi pada data diskusi.

Kata kunci: *sistem rekomendasi, identifier, normalisasi, stack overflow*

*Halaman ini sengaja dikosongkan*

# **IMPROVEMENT OF PROGRAMMING DISCUSSIONRECOMMENDER METHODBASED ONLINGUA::IdSPLITTER IDENTIFIER NORMALIZATION**

Name : Nanang Fakhrur Rozi  
NRP : 5112201045  
Supervisor I : Daniel Oranova Siahaan, S.Kom., M.Sc., P.D.Eng.  
Supervisor II : Fajar Baskoro, S.Kom., M.T.

## ***ABSTRACT***

Stack Overflow as a question and answer (Q&A) site has often been used as reference by programmers. Information or solutions in the software development process can be searched with the help of search engine on the site. However, differences in writing style, especially on the program identifier writing, often causing the recommendation (search) become incompatible with the programmers need. Some programmers write identifier in abbreviated form while another are not so that it lowers the recommendation performance. This research adopted the Lingua::IdSplitter to normalize the identifier on Stack Overflow discussion. The normalization process is done by separating the identifier comprising the composition of the term as well as expand the existing abbreviation on the identifier to the full form. The results showed that the identifier normalization using Lingua::IdSplitter only able to improve the system performance when the abbreviated term within identifierdominate in the dataset.

Keywords:*recommender system, identifier, normalization, stack overflow*

*Halaman ini sengaja dikosongkan*



## DAFTAR ISI

HALAMAN JUDUL .....	i
LEMBAR PENGESAHAN.....	iii
ABSTRAK.....	v
<i>ABSTRACT</i> .....	vii
KATA PENGANTAR .....	ix
DAFTAR ISI .....	xi
DAFTAR GAMBAR .....	xiii
DAFTAR TABEL.....	xv
BAB 1 PENDAHULUAN .....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	5
1.3 Batasan Masalah .....	5
1.4 Tujuan Penelitian.....	5
1.5 Manfaat Penelitian .....	5
1.6 Kontribusi Penelitian .....	6
BAB 2 KAJIAN PUSTAKA DAN DASAR TEORI .....	7
2.1 Sistem Rekomendasi.....	7
2.1.1 Rekomendasi Berbasis Konten.....	7
2.1.2 Rekomendasi Berbasis Metode Kolaboratif.....	7
2.1.3 Rekomendasi Berbasis Metode Hibrida.....	8
2.2 <i>Identifier</i> pada Kode Program .....	8
2.3 Sistem Rekomendasi dalam Rekayasa Perangkat Lunak.....	9
2.4 Praproses pada Sistem Rekomendasi.....	12
2.4.1 Tokenisasi, Filter Simbol, dan <i>Case Folding</i> .....	13
2.4.2 Normalisasi dengan <i>Lingua::IdSplitter</i> .....	13
2.4.3 Penghapusan Stop Word .....	16
2.4.4 Stemming .....	16
2.5 Perekomendasi Berdasarkan Similaritas Data .....	17
2.6 Stack Overflow Sebagai Sumber Rekomendasi .....	18
2.7 Evaluasi Sistem Rekomendasi.....	20

BAB 3 METODE PENELITIAN .....	21
3.1 Data Masukan .....	22
3.2 Filter Simbol.....	24
3.3 Tokenisasi.....	25
3.4 Normalisasi Term dengan Lingua::IdSplitter dan <i>Case Folding</i> .....	26
3.5 Penghapusan <i>Stop Word</i> .....	27
3.6 Stemming .....	28
3.7 Pembobotan TF-IDF .....	28
3.8 Cosine <i>Similarity</i> .....	29
3.9 Evaluasi Hasil.....	29
BAB 4 HASIL DAN PEMBAHASAN .....	31
4.1 Proses Rekomendasi .....	31
4.2 Hasil Pengujian.....	35
4.2.1 Pengujian dengan Kode Program Terkait <i>Thread Synchronization</i> .....	35
4.2.2 Pengujian dengan Kode Program Terkait <i>Accessing MySQL</i> ....	37
4.2.3 Pengujian dengan Kode Program Terkait <i>Bubble Sort</i> .....	39
4.2.4 Pengujian dengan Kode Program Terkait <i>Draw Rectangle</i> .....	40
4.3 Uji Coba Menggunakan Data Alternatif .....	42
4.4 Ringkasan Hasil Uji Coba .....	43
BAB 5 KESIMPULAN.....	45
5.1 Kesimpulan .....	45
5.2 Saran .....	45
DAFTAR PUSTAKA .....	47
Lampiran 1 Kamus Kata pada Lingua::IdSplitter.....	49
Lampiran 2 Stop List.....	51
Lampiran 3 Detail Data Uji .....	55
Lampiran 4 Hasil Uji Coba Menggunakan Normalisasi Lingua::IdSplitter .....	57
Lampiran 5 Hasil Uji Coba Tanpa Normalisasi Lingua::IdSplitter.....	59
BIODATA PENULIS .....	61

## DAFTAR TABEL

Tabel 2.1	Ringkasan Perbandingan Penelitian yang Diacu .....	12
Tabel 2.2	Himpunan Kandidat Hasil Split Berdasarkan Skor Tertinggi.....	15
Tabel 3.1	Spesifikasi Data Uji .....	30

*Halaman ini sengaja dikosongkan*

## DAFTAR GAMBAR

Gambar 2.1	Hasil Tokenisasi, Filter Simbol, dan Case Folding pada Data Teks .....	13
Gambar 2.2	Mekanisme Pemisahan dan Ekspansi Identifier .....	14
Gambar 2.3	Otomata untuk Identifier timesort.....	15
Gambar 2.4	Penentuan Kesamaan Dokumen (Vektor) Berdasarkan Sudut yang Terbentuk .....	18
Gambar 2.5	Ilustrasi Diskusi pada Stack Overflow .....	19
Gambar 3.1	Desain Sistem Rekomendasi dengan Normalisasi Lingua::IdSplitter .....	21
Gambar 3.2	Query untuk Mengambil Pertanyaan dari Stack Overflow .....	22
Gambar 3.3	Query untuk Mengambil Jawaban dari Stack Overflow .....	22
Gambar 3.4	Data Masukan Ketika Dibaca Oleh Sistem .....	23
Gambar 3.5	Query Sistem Berupa Kode Program Java .....	24
Gambar 3.6	Hasil Filter Simbol.....	25
Gambar 3.7	Hasil Tokenisasi.....	26
Gambar 3.8	Hasil Sebelum dan Sesudah Pemrosesan dengan Lingua::IdSplitter .....	26
Gambar 3.9	Format Eksekusi Lingua::IdSplitter .....	27
Gambar 3.10	Hasil Sebelum dan Sesudah Penghapusan Stop Word.....	27
Gambar 3.11	Hasil Sebelum dan Sesudah Stemming .....	28
Gambar 3.12	Contoh Daftar Term Beserta Bobotnya dalam Bentuk Matriks Term-Dokumen.....	29
Gambar 4.1	Contoh Kode Program Sebelum Diekstrak Kata Kuncinya .....	32
Gambar 4.2	Kata Kunci Hasil Ekstraksi dengan Memanfaatkan Normalisasi Identifier .....	32
Gambar 4.3	Kata Kunci Hasil Ekstraksi Tanpa Normalisasi Identifier .....	33
Gambar 4.4	Proses Rekomendasi.....	34
Gambar 4.5	Perbandingan Precision untuk Topik Thread Synchronization .....	35
Gambar 4.6	Perbandingan Recall untuk Topik Thread Synchronization.....	36

Gambar 4.7	Perbandingan Precision untuk Topik Accessing MySQL .....	37
Gambar 4.8	Perbandingan Recall untuk Topik Accessing MySQL .....	38
Gambar 4.9	Perbandingan Precision untuk Topik Bubble Sort.....	39
Gambar 4.10	Perbandingan Recall untuk Topik Bubble Sort .....	40
Gambar 4.11	Perbandingan Precision untuk Topik Draw Rectangle .....	40
Gambar 4.12	Perbandingan Recall untuk Topik Draw Rectangle .....	41
Gambar 4.13	Ringkasan Perbandingan Uji Coba Data Alternatif .....	42

# BAB 1

## PENDAHULUAN

### 1.1 Latar Belakang

Perusahaan perangkat lunak yang baik seharusnya dapat menyediakan *deliverable* dengan waktu yang singkat dan dengan kualitas yang baik pula. Untuk mencapai hal tersebut, dibutuhkan manajemen serta kerja sama yang baik dari setiap *stakeholder*, antara lain *programmer* yang menjadi tulang punggung terciptanya perangkat lunak. Dengan demikian, *programmer* dengan performa yang baik merupakan suatu keniscayaan dalam menyelesaikan pengerjaan perangkat lunak.

Cegielskidan Hall (2006) menyebutkan bahwa performap*programmer* dapat diukur berdasarkan kepribadian, kemampuan kognitif, serta tingkat kepercayaan terhadap nilai teoretis (*theoretical value belief*). *Programmer* yang selalu mencari pembuktian kebenaran terhadap hasil kerjanya, tidak asal-asalan dalam menyediakan solusi dalam bentuk kode program, cenderung memiliki performayang baikdalam jangka panjang. Hal tersebut memaksa *programmer*, terutama *programmer* pemula, untuk menginvestasikan waktunya dalam meningkatkan kemampuan terkait pemahaman algoritma, cara memprogram, informasi-informasi pendukung, atau hal-hal lain yang berkaitan dengan perangkat lunak yang akan dibangun. Mendapatkan informasi dari anggota tim atau teman kerja menjadi salah satu caranya (LaToza dkk., 2006). Namun, hal itu akan membebani anggota lain untuk juga meluangkan waktunya yang pada akhirnya akanberimbas pada perlambatan waktu *delivery*.

Guna menghindari hal tersebut, cara lain pun ditempuh. Seiring perkembangan web 2.0, informasi di internet pun semakin beragam, informasi-informasi bergunaterkait proses pengembangan perangkat lunak yang ditenagai oleh komunitas pun tersedia, antara lain yang berbentuk situs tanya-jawab (Anderson dkk., 2012). Media tersebut pada akhirnya banyak digunakan oleh para pengembang perangkat lunak untuk mendapatkan solusi terhadap masalah-masalah

yang sedang dihadapi. Media tersebut dapat berupa sebagian dari *platform* media sosial seperti grup di Facebook atau situs khusus seperti StackOverflow.

Stack Overflow merupakan situs tanya-jawab populer yang berfokus pada bidang pemrograman. Diskusi yang terjadi pada Stack Overflow adalah diskusi yang sangat aktif dan terpercaya. Anderson dkk. (2012) menyebutkan bahwa lebih dari 90% pertanyaan yang dilontarkan oleh penanya dijawab tidak lebih dari 12 menit oleh pengguna lainnya, penanya relatif mendapatkan kepuasan atas jawaban tersebut, serta kualitas jawaban yang juga selalu mendapat perhatian dari membebernya.

*Programmer* dapat mengakses rekomendasi informasi tersebut melalui kolom pencarian yang tersedia pada situs. Mengingat mesin pencari pada Stack Overflow juga diadopsi dari mesin pencarian teks yang umum digunakan, maka berbagai penyesuaian pun perlu dilakukan, sehingga hasil pencarian yang diberikan juga belum tentu sesuai dengan rekomendasi yang dibutuhkan pengguna. Hal tersebut umumnya diakibatkan oleh sulitnya merumuskan *query* yang tepat sebagai masukan sistem pencarian, inkonsistensi penggunaan terminologi, serta berbagai kesulitan lainnya (Manning, 2008). Oleh karena itu, penelitian-penelitian tetap dilakukan untuk menyempurnakan hasil pencarian tersebut.

Cordeiro dkk. (2012) melakukan penelitian tentang sistem rekomendasi berbasis konteks. Penelitian tersebut dilatarbelakangi oleh ketidakmampuan mesin pencari secara umum untuk memberikan rekomendasi diskusi Stack Overflow yang baik apabila *query* yang diberikan berupa informasi *stack trace*. Informasi tersebut biasanya muncul ketika terjadi eksepsi saat menjalankan program Java. Informasi *stack trace* tersebut diasumsikan oleh peneliti sebagai konteks permasalahan *programmer*. Dalam memberikan rekomendasi, sistem akan mengekstrak *stack trace* yang muncul untuk kemudian dicari strukturnya berdasarkan nama eksepsi, nama *method*, beserta referensinya. Struktur tersebut selanjutnya akan dijadikan sebagai acuan *query* yang akan dibangkitkan. Rekomendasi yang dihasilkan memiliki performa yang lebih baik dibandingkan dengan hasil rekomendasi dengan *query* yang berbasis kata kunci maupun hasil rekomendasi dari mesin pencari pada situs Stack Overflow sendiri.



Namun, penelitian ini hanya akan memiliki hasil yang maksimal apabila data diskusi yang dimaksud mengandung informasi *stack trace* pula.

Penelitian lain yang telah dilakukan adalah Seahawk (Ponzanelli dkk., 2013). Seahawk merupakan sistem rekomendasi data diskusi StackOverflow yang mampu memberikan hasil rekomendasi dengan *query* yang dibangkitkan secara otomatis berdasarkan kode program yang sedang dibuka pada editor kode dalam IDE Eclipse. Kode program tersebut diasumsikan sebagai teks bahasa alami sehingga *query* yang dibangkitkan akan menghasilkan kumpulan kata kunci yang diperoleh dari proses tokenisasi dan beberapa praproses lain terhadap kode program. Peningkatan data diskusi dilakukan dengan memanfaatkan Apache Solr dengan model pembobotan TF-IDF. Pengujian dilakukan menggunakan kode program Java yang telah disiapkan sebelumnya. Hasilnya, Seahawk secara umum dapat merekomendasikan diskusi yang berguna bagi *programmer*. Namun, penelitian ini masih memiliki kekurangan. Hasil rekomendasi menjadi buruk ketika kata kunci yang dibangkitkan mengandung unsur singkatan seperti `strCmp`, `usrId`, atau lainnya. Padahal, pembentukan term yang baik merupakan kunci sukses untuk mendapatkan informasi yang tepat dari literatur yang tersimpan (Krauthammer dan Nenadic, 2004).

Penelitian lain terkait sistem rekomendasi yaitu penelitian Arwan dkk. (2015) yang menghasilkan JECO (Java Example Code). Penelitian ini mengusulkan sistem rekomendasi kode program pada data diskusi Stack Overflow dengan pemodelan topik LDA (Latent Dirichlet Allocation). *Query* yang digunakan berupa kata kunci yang dimasukkan secara manual oleh pengguna. Data diskusi dan *query* kemudian dicari topik beserta proporsinya menggunakan LDA. Berdasarkan topik tersebut, *query* dan data diskusi dihitung kesamaannya sehingga didapatkan hasil rekomendasi dengan memanfaatkan *Cosine Similarity*. Hasilnya, JECO mampu merekomendasikan kode program Java ke *programmer* dengan rata-rata *precision* sebesar 48% dan rata-rata *recall* sebesar 58%. Namun, term singkatan juga tidak diperhatikan pada penelitian ini baik pada *query* maupun data diskusi. Hal tersebut tentu akan mempengaruhi *recall* pada proses rekomendasi.

Unsur singkatan biasa muncul pada *identifier* kode program. Seorang *programmer* terkadang lebih suka menuliskan `string` dengan `str` saja pada *identifier*-nya. Sementara *programmer* lain lebih suka menulis *identifier* dalam bentuk lengkap dibandingkan dengan bentuk singkatan. Sebuah penelitian yang dilakukan oleh Carvalho dkk. (2015) mengembangkan `Lingua::IdSplitter`. Algoritma tersebut memiliki kemampuan untuk memisah *identifier* yang umumnya terdiri atas komposisi term serta mengekspansi term singkatan menjadi term lengkap, baik yang menggunakan model penulisan *all lowercase* maupun *camel case*. `Lingua::IdSplitter` dikembangkan mengingat *identifier* merupakan salah satu sumber informasi yang cukup relevan untuk memahami sebuah program. Selain itu, *identifier* yang lengkap cenderung lebih memudahkan program untuk dipahami (Lawrie dkk., 2007). Cara pemisahan *identifier* yang dilakukan `Lingua::IdSplitter` dilakukan dengan dua tahap, yakni *hard split* dan *soft split*. *Hard split* dilakukan untuk memisahkan *identifier* berdasarkan karakter tertentu seperti tanda garis bawah (*underscore*) atau berdasarkan tata penulisan *camel case*. Sedangkan *soft split* akan memisah *identifier* yang tidak terlalu terlihat pemisahannya. Hal ini dilakukan dengan bantuan kamus kata untuk mendapatkan himpunan kandidat katanya. Kandidat kata tersebut selanjutnya akan diperingkat berdasarkan otomatis yang dibangun berdasarkan himpunan kandidat kata. Himpunan kata dengan skor tertinggi akan menjadi keluaran `Lingua::IdSplitter`. Penelitian ini membandingkan algoritma yang diusulkan dengan algoritma-algoritma lain yang sejenis. Hasilnya, algoritma ini mampu menghasilkan term yang telah dipisah dan diekspansi dengan tingkat *f-measure* sebesar 90%, tidak jauh berbeda dengan hasil algoritma pembanding lainnya.

Penelitian ini mengadopsi `Lingua::IdSplitter` sebagai algoritma untuk menormalisasi *identifier* pada sistem rekomendasi data diskusi dari Stack Overflow yang dibangun dengan masukan berupa kode program Java. Dengan term yang baik, maka hasil rekomendasi akan menjadi lebih baik pula. Selain itu, model data yang digunakan adalah model data Bag-of-Words (BOW) yang telah diboboti menggunakan algoritma TF-IDF. Model tersebut dipilih mengingat penggunaan LDA untuk memodelkan topik masih memiliki kekurangan karena

tidak adanya proses seleksi atau pemilihan term penting (Hu, 2009) sehingga akan mengakibatkan ketidaksesuaian proporsi topik terhadap dokumen.

## 1.2 Rumusan Masalah

Rumusan masalah dalam penelitian ini yaitu:

1. Bagaimana membangun sistem rekomendasi diskusi pemrograman dengan masukan berupa kode program?
2. Bagaimana hasil metode rekomendasi diskusi pemrograman dapat diperbaiki dengan normalisasi *identifier* menggunakan `Lingua::IdSplitter`?

## 1.3 Batasan Masalah

Agar penelitian yang dilakukan dapat lebih terarah dan sesuai dengan tujuan, maka perlu dibatasi ruang lingkup dalam penelitian ini. Adapun batasan masalah dalam penelitian ini yaitu:

1. Masukan yang dapat diproses adalah kode program Java dengan *identifier* yang dibentuk dari kata-kata bahasa Inggris.
2. Data diskusi yang diambil adalah data diskusi yang memiliki label Java saja mengingat data masukan pada proses rekomendasi juga berupa kode program dengan bahasa Java.

## 1.4 Tujuan Penelitian

Tujuan penelitian ini yaitu membangun dan memperbaiki sistem rekomendasi diskusi pemrograman dengan melakukan normalisasi *identifier* menggunakan `Lingua::IdSplitter`.

## 1.5 Manfaat Penelitian

Manfaat dari penelitian ini yaitu dihasilkannya sistem rekomendasi diskusi pemrograman yang relevan yang mampu mendukung para *programmer* dalam menemukan solusi-solusi dari Stack Overflow terkait permasalahan pengembangan perangkat lunak.

## **1.6 Kontribusi Penelitian**

Kontribusi penelitian ini adalah perbaikan term indeks dan *query* pada sistem rekomendasi diskusi pemrograman dengan melakukan normalisasi *identifier* menggunakan `Lingua::IdSplitter`.

## **BAB 2**

### **KAJIAN PUSTAKA DAN DASAR TEORI**

#### **2.1 Sistem Rekomendasi**

Sistem rekomendasi merupakan metode yang mampu memberikan rekomendasi dengan memprediksi nilai sebuah item bagi pengguna untuk kemudian mempresentasikan item dengan nilai prediksi tertinggi sebagai hasil rekomendasi. Sistem ini awalnya merupakan sebuah mekanisme *information filtering*, yakni bertujuan menyaring informasi sebagai akibat membludaknya informasi di internet. Terdapat tiga pendekatan yang bisa digunakan dalam membangun sistem rekomendasi, yaitu rekomendasi berbasis konten, rekomendasi berbasis metode kolaboratif, dan rekomendasi berbasis metode hibrida (Adomavicius dan Tuzhilin, 2005).

##### **2.1.1 Rekomendasi Berbasis Konten**

Jenis sistem rekomendasi ini memiliki akar keilmuan yang berasal dari ilmu temu kembali informasi (*information retrieval*). Pada rekomendasi berbasis konten, pengguna memasukkan beberapa informasi ke dalam sistem, dalam hal ini disebut *query*, misalnya permasalahan yang sedang dihadapi. Sistem kemudian mengubah *query* tersebut ke model tertentu untuk kemudian digunakan oleh sistem dalam komputasi untuk menyaring begitu banyak informasi yang tersedia akan menghasilkan rekomendasi untuk disajikan kepada pengguna.

##### **2.1.2 Rekomendasi Berbasis Metode Kolaboratif**

Jenis sistem rekomendasi ini juga sering disebut *collaborative filtering system*. Model rekomendasi ini memanfaatkan data yang dimiliki oleh pengguna lain untuk membangun model. Artinya preferensi dari pengguna lain yang memiliki relasi ke pengguna yang ingin mendapat rekomendasi akan mempengaruhi hasil rekomendasinya.

### 2.1.3 Rekomendasi Berbasis Metode Hibrida

Jenis sistem rekomendasi ini menggabungkan karakteristik yang ada pada sistem rekomendasi berbasis konten dengan sistem rekomendasi berbasis metode kolaboratif. Model rekomendasi ini dikembangkan guna membantu menghindari batasan-batasan yang ada pada kedua jenis sistem rekomendasi tersebut.

Pada penelitian ini, model sistem rekomendasi yang digunakan adalah model rekomendasi berbasis konten. Model kolaboratif tidak digunakan mengingat model kolaboratif membutuhkan informasi dari pengguna lain yang dalam hal ini tidak tersedia. Selain itu, pemilihan model juga didasarkan pada penelitian-penelitian terdahulu yang digunakan sebagai acuan peneliti.

## 2.2 Identifier pada Kode Program

Menurut Gosling dkk. (2015), *identifier* dalam ranah pemrograman, khususnya bahasa pemrograman Java, adalah rangkaian karakter Java dan digit Java dengan panjang tak hingga. Syaratnya, karakter pertama pada *identifier* harus berupa karakter Java. Karakter Java sendiri merupakan karakter alfabet (non-simbol), yakni karakter-karakter yang biasa digunakan untuk merangkai kata, baik yang terdapat dalam kumpulan karakter latin maupun non-latin. Contoh karakter non-latin dapat ditemukan pada tulisan-tulisan berbahasa Arab, Cina, Jepang, Korea, atau lainnya. Sementara itu, digit Java merupakan angka 0–9, baik dalam bentuk karakter latin maupun non-latin. Sebagai tambahan, karakter-karakter yang digunakan oleh Java didasarkan pada standar pengodean *unicode*.

*Identifier* pada java biasa digunakan untuk menamai elemen-elemen di dalam kode program, seperti *class*, variabel, atau *method* (Vohra dkk., 2015). Pada umumnya, penulisannya mengacu pada konvensi yang dibuat oleh Sun Microsystems (1997), yakni menggunakan aturan *camel-case* atau *underscore* seperti pada penulisan variabel `myWidth` dan `MAX_WIDTH`. Selain itu, penulisan *identifier* bersifat *case-sensitive*, yakni membedakan huruf besar dan kecil, misalnya membedakan antara variabel `var2a` dan `var2A`. Namun, dalam penelitian ini, konsep sensitivitas karakter tersebut diabaikan. Langkah tersebut

diambil mengingat perbedaan *case* tidak akan membedakan makna atau konsep *identifier*.

Mengingat *identifier* umumnya berupa gabungan kata yang mengakibatkan panjangnya *identifier*, maka ada *programmeryang* menggunakan unsur singkatan dalam membentuk *identifier*. Singkatan-singkatan tersebut antara lain singkatan-singkatan yang umum digunakan seperti HTML, URL, atau singkatan-singkatan dalam bahasa tertentu.

### **2.3 Sistem Rekomendasi dalam Rekayasa Perangkat Lunak**

Sistem rekomendasi dalam rekayasa perangkat lunak atau sering disebut *Recommendation Systems for Software Engineering (RSSE)* difungsikan untuk membantu pengembang dalam berbagai aktivitas, mulai dari sugesti informasi, penggunaan ulang kode program, hingga penulisan laporan bug yang efektif (Robillard dkk., 2010). Dengan semakin kompleksnya proses pengembangan perangkat lunak, maka sistem semacam ini sangat dibutuhkan oleh para pengembang, terlebih untuk pengembang dengan tingkat pengalaman yang masih rendah. Guna mendukung hal tersebut, berbagai penelitian telah dilakukan.

Salah satu topik RSSE yang telah diteliti yaitu sistem rekomendasi informasi, khususnya berupa data forum diskusi (tanya-jawab). Penelitian terhadap data diskusi dilakukan mengingat informasi tersebut merupakan salah satu sumber informasi yang penting dan layak untuk dipertimbangkan sebagai sumber rekomendasi. Kelayakan tersebut didasarkan pada penelitian Anderson dkk. (2012) yang menyatakan bahwa berbagai situs tanya-jawab yang ditenagai oleh komunitas saat ini dapat dijadikan sebagai sumber informasi mengingat mekanisme reputasi ataupun mekanisme-mekanisme lain yang mendukung kekuatan dan validitas situs tanya-jawab tersebut. Studi kasus yang digunakan pada penelitian tersebut adalah situs tanya-jawab Stack Overflow.

Penelitian ini disusun dengan berlandaskan pada penelitian-penelitian terdahulu yang sudah pernah dilakukan. Penelitian-penelitian tersebut adalah penelitian yang berkaitan dengan topik sistem rekomendasi data diskusi pemrograman yang berasal dari Stack Overflow yang memanfaatkan model rekomendasi berbasis konten. Sistem rekomendasi tersebut dirancang untuk

membantu para pengembang, khususnya *programmer*, dalam menyelesaikan permasalahan-permasalahan ketika mengembangkan program.

Cordeiro dkk. (2012) telah melakukan penelitian tentang sistem rekomendasi berbasis konteks. Konteks permasalahan program dilihat dari informasi *stack trace* yang umumnya muncul saat terjadi eksepsi ketika program dijalankan. Informasi tersebut merupakan informasi yang penting karena berkaitan erat dengan permasalahan yang sedang terjadi. Akan tetapi, mesin pencari yang umum digunakan belum mampu menghasilkan rekomendasi solusi atas permasalahan program jika masukannya berupa informasi *stack trace*. Dalam mendapatkan hasil rekomendasi, informasi *stack trace* yang muncul akan diekstrak sehingga didapatkan informasi eksepsi beserta referensinya yang dibentuk dalam sebuah struktur. Struktur tersebut selanjutnya akan diubah menjadi representasi leksikal berdasarkan nama eksepsi dan nama referensi menggunakan aturan penulisan *camel case*. Sebelum diubah ke dalam representasi leksikal, informasi referensi tersebut akan diboboti terlebih dahulu guna mengetahui referensi-referensi yang paling relevan terhadap eksepsi. Lima term dengan bobot tertinggi akan digunakan sebagai *query* pada sistem rekomendasi. Hasil penelitian menunjukkan bahwa pembentukan *query* yang berasal dari informasi *stack trace* ini mampu menghasilkan rekomendasi yang lebih baik dibandingkan dengan mesin pencarian yang umum digunakan serta fasilitas pencarian yang telah disediakan oleh Stack Overflow dengan masukan berupa informasi *stack trace*. Adapun proses pengindeksannya, penelitian ini memanfaatkan Apache Lucene. Lebih jauh, penelitian ini masih memiliki kelemahan, yakni hasil rekomendasi hanya terbatas pada data-data diskusi yang mengandung informasi *stack traces* saja. Dengan demikian, informasi yang direkomendasikan juga akan menjadi sangat sempit, padahal terkadang permasalahan-permasalahan yang terjadi membutuhkan dukungan pengetahuan yang lebih luas. Selain itu, karena proses representasi leksikalnya menggunakan aturan *camel case*, maka performa rekomendasi juga akan menurun apabila ditemukan *method* atau referensi yang tidak menggunakan aturan *camel case*.

Penelitian berikutnya adalah penelitian yang dilakukan oleh Ponzanelli dkk. (2013). Penelitian tersebut mengembangkan sebuah sistem rekomendasi



dengan konteks permasalahan yang didasarkan pada kode program yang sedang terbuka pada IDE Eclipse. Sintaks-sintaks pada program tersebut selanjutnya akan diekstrak sehingga didapatkan kumpulan term. Kumpulan term yang telah didapatkan direpresentasikan dalam bentuk vektor terbobot berdasarkan pembobotan TF-IDF. Sebelum dibentuk representasi vektornya, akan dilakukan praproses pada term-term yang telah diperoleh. Beberapa praproses yang dilakukan antara lain penghapusan *stop word*, *stemming*, filter sinonim, dan pengubahan term ke bentuk *lower case*. Sepuluh term dengan frekuensi tertinggi akan dijadikan sebagai *query* untuk menghasilkan rekomendasi. Penelitian ini mampu memberikan rekomendasi terutama untuk informasi-informasi yang bersifat spesifik. Namun, penelitian ini memiliki kekurangan. Jika term hasil ekstraksi mengandung unsur singkatan, hasil rekomendasinya menjadi buruk. Kelemahan lainnya yaitu ketidakmampuan sistem untuk merekomendasikan informasi yang relevan ketika kode program tidak menggunakan standar-standar yang sering digunakan seperti *framework* atau implementasi buruk dari pengembang.

Penelitian lainnya yaitu penelitian yang dilakukan oleh Arwan dkk. (2015). Masukan dalam sistem rekomendasi pada penelitian ini berupa *query* yang dimasukkan pengguna secara manual. Model data yang digunakan berbeda dengan kedua penelitian sebelumnya. Penelitian ini memanfaatkan pemodelan topik untuk merepresentasikan data diskusi yang akan direkomendasikan. Algoritma pemodelan topik yang digunakan adalah Latent Dirichlet Allocation (LDA). Setiap dokumen memungkinkan untuk memiliki lebih dari satu topik bahasan. Tiap topik yang melekat pada dokumen akan memiliki proporsi yang berbeda, satu topik bisa memiliki porsi yang lebih besar dibandingkan topik lainnya. Porsi topik yang melekat pada dokumen tersebut yang pada akhirnya akan dijadikan sebagai fitur dokumen yang juga dirupakan dalam model ruang vektor sebagaimana dua penelitian sebelumnya. Penentuan porsi topik tersebut juga akan dilakukan pada *query* pengguna, sehingga pada akhirnya dapat ditentukan kedekatan *query* dengan data diskusi yang ada menggunakan perhitungan Cosine Similarity. Hasil rekomendasi mampu menghasilkan rata-rata *precision* sebesar 48% dan rata-rata *recall* sebesar 58%.

Pemodelan topik digunakan pada penelitian ini untuk merepresentasikan dokumen dengan tujuan menangani permasalahan kesamaan topik antara satu term dengan term lain yang terkadang tidak dapat diidentifikasi oleh model vektor yang berbasis TF-IDF. Akan tetapi, pemodelan topik juga masih memiliki kekurangan antara lain tidak adanya pemilihan term penting yang mengakibatkan ketidaktepatan proporsi topik terhadap dokumen. Untuk memperjelas, Tabel 2.1 akan memberikan ringkasan singkat tentang ketiga penelitian yang telah diambil sebagai dasar dalam penelitian ini.

Tabel 2.1 Ringkasan Perbandingan Penelitian yang Diacu

Komponen Penelitian	Penulis		
	Cordeiro dkk. (2012)	Ponzanelli dkk. (2013)	Arwan dkk. (2015)
<b>Judul</b>	<i>Context-Based Recommendation to Support Problem Solving in Software Development</i>	<i>Leveraging Crowd Knowledge for Software Comprehension and Development</i>	<i>Source Code Retrieval on StackOverflow Using LDA</i>
<b>Model Data</b>	Model vektor (bobot TF-IDF)	Model vektor (bobot TF-IDF)	Model vektor (proporsi topik LDA)
<b>Praproses</b>	Praproses bawaan Apache Lucene	Praproses bawaan Apache Solr	Praproses umum mesin pencari dan <i>identifier splitting</i>
<b>Simpulan</b>	Hasil lebih baik dibandingkan dengan mesin pencari umum dan fasilitas pencarian Stack Overflow	Hasil rekomendasi secara umum baik	Rata-rata <i>precision</i> 48% dan rata-rata <i>recall</i> 58%

Ketiga penelitian tersebut belum ada yang melakukan normalisasi *identifier* atau secara khusus mengekspansi *identifier* baik pada data maupun *query* yang akan mengakibatkan turunnya *recall* hasil rekomendasi. Hal tersebut juga mengingatkan akan pentingnya pembentukan term yang baik (Krauthammer dan Nenadic, 2004).

## 2.4 Praproses pada Sistem Rekomendasi

Praproses merupakan salah satu hal yang lumrah dilakukan sebelum sebuah data diproses lebih lanjut. Praproses ini digunakan untuk menghapus informasi-informasi yang tidak berguna ataupun menormalisasi data agar data menjadi bersih atau sesuai format serta siap untuk diproses. Mengingat data yang

akan dibahas mengandung kode program, maka praproses yang dilakukan salah satunya harus mempertimbangkan pola penulisan kode program (Mahmoud dan Niu, 2011).

#### 2.4.1 Tokenisasi, Filter Simbol, dan *Case Folding*

Jika diberikan kumpulan karakter dan pendefinisian unit dokumen, tokenisasi adalah proses pemotongan kumpulan karakter tersebut menjadi potongan-potongan yang disebut token. Gambar 2.1 merupakan contoh dari tokenisasi. Lebih dari itu, duplikasi token yang telah ditiadakan (token unik) sering juga disebut sebagai term pada bidang ilmu temu kembali informasi (Manning dkk., 2008). Mengiringi tokenisasi, penghapusan simbol atau data non-tekstual biasanya juga akan dilakukan. Pada kode program, simbol ini dapat berupa tanda plus (+), minus (-), modulus (%), atau simbol lain seperti @ yang biasa digunakan untuk memberi anotasi pada kode program (Mahmoud dan Niu, 2011). Sebagai tambahan, pengubahan huruf besar ke huruf kecil (*case folding*) juga biasa dilakukan untuk menyamaratakan karakter ke bentuk huruf kecil saja.

Masukan
even though backquote (') is not a standard quote character in
Keluaran
Even though backquote is not a standard quote character in

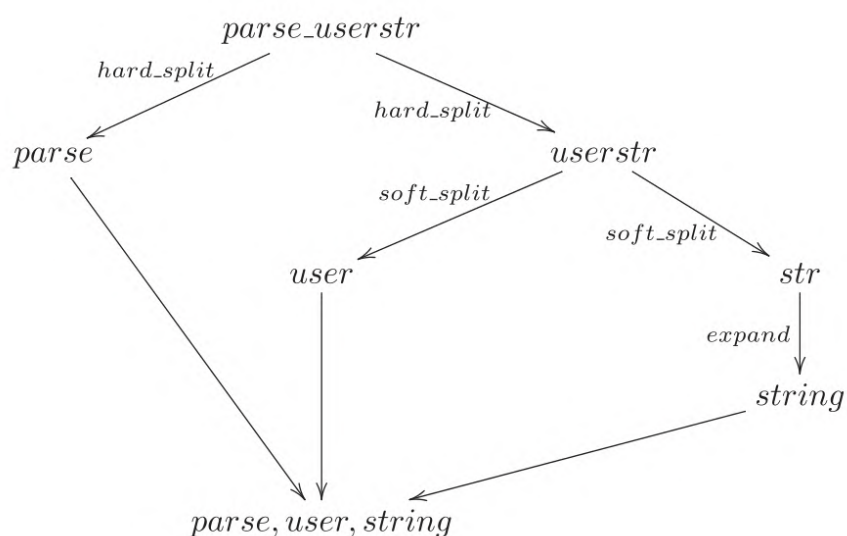
Gambar 2.1 Hasil Tokenisasi, Filter Simbol, dan Case Folding pada Data Teks

#### 2.4.2 Normalisasi dengan `Lingua::IdSplitter`

`Lingua::IdSplitter` merupakan sebuah algoritma yang dikembangkan oleh Carvalho dkk. (2015). Algoritma tersebut memiliki kemampuan untuk

memisahkan *identifier* kode program yang umumnya terdiri atas komposisi term serta mengekspansi term singkatan menjadi term lengkap, baik yang menggunakan model penulisan *all lowercase* maupun *camel case*. `Lingua::IdSplitter` dikembangkan mengingat *identifier* merupakan salah satu sumber informasi yang cukup relevan dalam memahami sebuah program. Selain itu, *identifier* yang lengkap cenderung lebih memudahkan program untuk dimengerti (Lawrie dkk., 2007).

Cara pemisahan *identifier* pada `Lingua::IdSplitter` dilakukan melalui dua tahap, yakni *hard split* dan *soft split*. *Hard split* dilakukan untuk memisahkan *identifier* berdasarkan karakter tertentu seperti tanda garis bawah (*underscore*) atau berdasarkan aturan penulisan *camel case*. Sementara itu, *soft split* akan memisahkan *identifier* yang tidak terlalu terlihat (tanda) pemisahannya. Ilustrasi proses pemisahan dan ekspansi ini dapat dilihat pada Gambar 2.2.



Gambar 2.2 Mekanisme Pemisahan dan Ekspansi *Identifier* (Carvalho dkk., 2015)

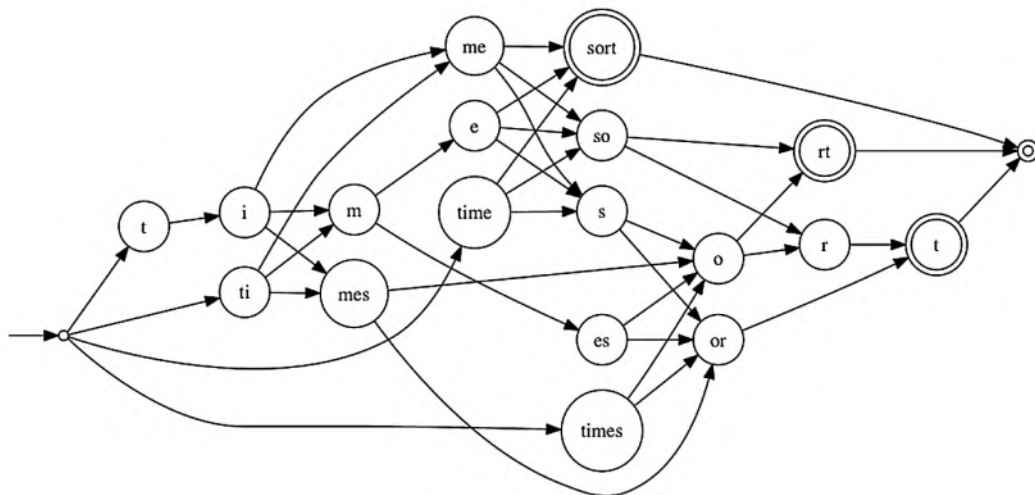
*Soft split* dilakukan salah satunya dengan bantuan kamus kata (akronim dan singkatan) untuk mendapatkan himpunan kandidat katanya. Kamus kata yang digunakan dapat dilihat pada Lampiran 1. Apabila kandidat kata tersebut tidak ditemukan dalam kamus, maka `Lingua::IdSplitter` akan memotong *identifier* dengan jumlah karakter tertentu untuk kemudian dicek validitas katanya. Kandidat kata tersebut selanjutnya akan diperingkat berdasarkan otomatis yang dibangun

berdasarkan himpunan kandidat kata. Himpunan kata dengan skor tertinggi akan menjadi keluaran `Lingua::IdSplitter`. Ilustrasi hasil penskoran himpunan kandidat kata dapat dilihat pada Tabel 2.2. Sedangkan otomata yang dibangun dapat dilihat pada Gambar 2.3.

Tabel 2.2 Himpunan Kandidat Hasil Split Berdasarkan Skor Tertinggi

Kandidat <i>Split</i>	Skor
{time, sort}	1,4400
{ti, me, sort}	0,1920
{time, so, rt}	0,1920
{times, o, rt}	0,1500
{...}	...

Sumber: Carvalho dkk. (2015)



Gambar 2.3 Otomata untuk Identifier *timesort* (Carvalho dkk., 2015)

Otomata pada Gambar 2.3 digunakan untuk menghitung skor kandidat hasil pemisahan. Berdasarkan penghitungan skor kandidat kata pada Tabel 2.2, didapatkan bahwa kandidat term *time+sort* adalah term yang paling baik untuk menjadi hasil pemisahan atas term *timesort*.

Adapun normalisasi `Lingua::IdSplitter` yang digunakan dalam penelitian ini merupakan program berbahasa Perl yang telah tersedia secara bebas. Program ini dapat difungsikan sebagai pustaka pemrograman (Perl) atau dijalankan sebagai program eksternal melalui antarmuka baris perintah (*command line interface*) yang tersedia.

### 2.4.3 Penghapusan Stop Word

Penghapusan *stopword* adalah proses penghapusan term yang tidak memiliki arti, tidak relevan, atau term yang sangat umum. Term yang diperoleh dari tahap tokenisasi dicek dalam suatu *stoplist*, apabila sebuah kata masuk di dalam *stop list* maka kata tersebut tidak akan diproses lebih lanjut. Sebaliknya, apabila kata tidak termasuk di dalam *stop list* maka kata tersebut akan masuk pada proses berikutnya. *Stop list* tersimpan dalam suatu dokumen yang dimuat saat pemrosesan term. *Stop list* yang digunakan diambil dari MySQL (2006) yang memuat sebanyak 571 kata bahasa Inggris. Adapun daftar lengkapnya dapat dilihat pada Lampiran 2.

### 2.4.4 Stemming

*Stemming* dilakukan atas asumsi bahwa kata-kata yang memiliki stem yang sama memiliki makna yang serupa sehingga pengguna tidak keberatan untuk memperoleh dokumen-dokumen yang di dalamnya terdapat kata-kata dengan stem yang sama dengan *query*-nya. Teknik-teknik *stemming* dapat dikategorikan menjadi tiga, yaitu berdasarkan aturan bahasa tertentu, berdasarkan kamus, dan berdasarkan kemunculan bersama.

Proses ini memiliki dua tujuan. Dalam hal efisiensi, *stemming* mengurangi jumlah kata-kata unik dalam indeks sehingga mengurangi kebutuhan ruang penyimpanan untuk indeks dan mempercepat proses pencarian. Dalam hal keefektifan, *stemming* meningkatkan *recall* dengan mengurangi bentuk-bentuk kata menjadi bentuk dasarnya atau *stem*-nya. Sehingga, dokumen-dokumen yang menyertakan suatu kata dalam berbagai bentuknya memiliki kecenderungan yang sama untuk ditemukembali. Hal tersebut tidak akan diperoleh jika tiap bentuk kata disimpan secara terpisah dalam indeks. Adapun *stemmer* yang paling umum digunakan terhadap dokumen berbahasa Inggris adalah Porter Stemmer yang berbasis aturan bahasa Inggris (Manning, 2008). Meski masih memiliki kekurangan, namun metode *stemming* ini adalah metode yang cukup ringan dengan hasil yang cenderung lebih baik dibanding metode lain yang sejenis. Sedangkan hasil *stemming* berdasarkan kemunculan bersama akan sangat dipengaruhi oleh koleksi kata yang didapatkan dari dokumen. Sehingga, kata yang

tidak pernah muncul pada dokumen tidak akan bisa diperoleh *stem*-nya (Jivani, 2011). *Stemming* berdasarkan kamus tidak digunakan mengingat hasil berupa stem sudah mencukupi pada studi kasus yang diambil tanpa mengubah stem ke bentuk lema (kata valid pada kamus).

## 2.5 Rekomendasi Berdasarkan Similaritas Data

Perekomendasi baru bisa dilakukan setelah praproses dilakukan. Dengan masukan berupa daftar term, kemudian term tersebut diboboti dengan metode pembobotan umum TF-IDF. Algoritma ini menggabungkan dua konsep penghitungan bobot, yaitu frekuensi kemunculan kata di dalam sebuah dokumen tertentu atau biasa disebut TF (*term frequency*) dan invers frekuensi dokumen, yaitu frekuensi dokumen yang mengandung kata tersebut atau biasa disebut IDF (*inverse document frequency*). Bobot TF dan IDF kemudian digabungkan dengan cara dikalikan sebagaimana pada rumus (2.1) untuk mendapatkan bobot komposit berdasarkan pertimbangan kedua kriteria tersebut (Manning, 2008).

$$\text{TF-IDF}_{t,d} = \text{TF}_{t,d} \times \text{IDF}_t \quad (2.1)$$

dengan:

$\text{TF}_{t,d}$  = frekuensi kemunculan term  $t$  pada dokumen  $d$

$\text{IDF}_t$  = invers frekuensi dokumen untuk term  $t$  yang dihitung menggunakan rumus (2.2)

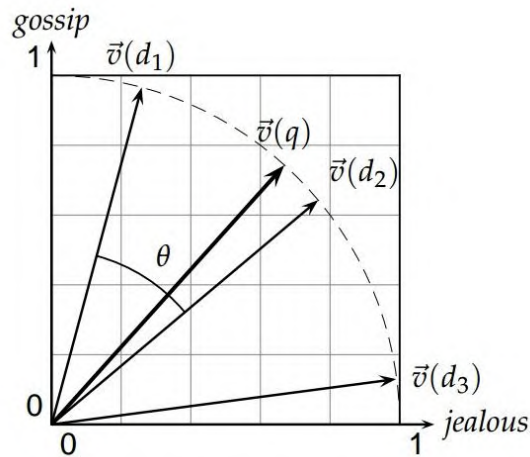
$$\text{IDF}_t = 1 + \log \frac{N}{\text{DF}_t + 1} \quad (2.2)$$

dengan:

$N$  = total dokumen pada koleksi/korpus

$\text{DF}_t$  = jumlah dokumen yang mengandung term  $t$

Pembobotan setiap term yang telah dilakukan akan direpresentasikan dalam ruang vektor (*vector space model*) dalam bentuk matriks term-dokumen sehingga setiap dokumen akan berbentuk sebagai sebuah fitur vektor. Vektor-vektor tersebut yang kemudian akan digunakan untuk menghitung kemiripan antardokumen dengan pendekatan *cosine similarity*. Ilustrasi *cosine similarity* dapat dilihat pada Gambar 2.4.



Gambar 2.4 Penentuan Kesamaan Dokumen (Vektor) Berdasarkan Sudut yang Terbentuk (Manning, 2008)

Tiap anak panah pada Gambar 2.4 mengilustrasikan dokumen yang direpresentasikan dalam bentuk vektor. Semakin kecil sudut yang terbentuk antara dua vektor, maka semakin mirip kedua dokumen tersebut. Adapun rumus yang digunakan adalah:

$$sim(d_1, d_2) = \frac{\vec{d}_1 \cdot \vec{d}_2}{\|\vec{d}_1\| \|\vec{d}_2\|} \quad (2.3)$$

dengan pembilang berupa operasi *dot product* antara vektor pertama dan vektor kedua. Sedangkan penyebut berupa jarak Euclidian antara kedua vektor.

## 2.6 Stack Overflow Sebagai Sumber Rekomendasi

Stack Overflow merupakan situs tanya-jawab populer yang berfokus pada bidang pemrograman. Diskusi yang terjadi pada Stack Overflow adalah diskusi yang sangat aktif dan terpercaya. Anderson dkk. (2012) menyebutkan bahwa lebih dari 90% pertanyaan yang dilontarkan oleh penanya dijawab tidak lebih dari 12 menit oleh pengguna lainnya, penanya relatif mendapatkan kepuasan atas jawaban tersebut, serta kualitas jawaban yang didapat juga selalu mendapat perhatian dari membeinya. Ilustrasi diskusi pada Stack Overflow dapat dilihat pada Gambar 2.5.

Diskusi pada Stack Overflow umumnya berupa data teks. Selain itu, pengguna umumnya mencantumkan potongan kode program pada pertanyaan



ataupun jawaban manakala dibutuhkan. Dan terkadang, pengguna juga mencantumkan gambar yang mampu mendukung pemahaman terhadap pertanyaan ataupun jawaban pengguna. Selain bertanya dan menjawab, pengguna Stack Overflow diberi fasilitas untuk melakukan *vote* terhadap pertanyaan maupun jawaban. Pertanyaan yang mendapat banyak *vote* dapat dikategorikan sebagai pertanyaan yang penting dan merupakan permasalahan yang sering terjadi. Sedangkan jawaban dengan *vote* yang banyak menandakan jawaban tersebut merupakan jawaban yang sangat bisa dipertimbangkan sebagai solusi akhir yang dapat diterima oleh penanya. *Vote* pada Gambar 2.5 dirupakan dengan ikon segitiga menghadap ke atas (*vote up*) dan ke bawah (*vote down*).

The screenshot shows a Stack Overflow question with the following content:

**How do I turn off text antialiasing in this Java function?**

▲ Basically I would like to turn off antialias in the following:

2 ▼ ★ 1

```
public BufferedImage createText(String text) {
    //create image
    BufferedImage image = new BufferedImage(95, 20,
        BufferedImage.TYPE_INT_ARGB);
    Graphics2D graphics = (Graphics2D) image.getGraphics();

    //set background
    graphics.setColor(Color.white);
    graphics.fillRect(0, 0, 95, 20);

    //draw text
    graphics.setRenderingHint(RenderingHints.KEY_ANTIALIASING,
        RenderingHints.VALUE_ANTIALIAS_OFF);

    graphics.setColor(Color.black);
    graphics.setFont(new Font("volter", Font.PLAIN, 9));
    graphics.drawString(text, 0, 10);

    return image;
}
```

but it's not working, here is something this function generates:

testing

I just want black and white to be used, nothing else so it's important I get antialias disabled!

java antialiasing

share improve this question

edited Dec 18 '11 at 6:15  
 palacsint 14.8k ● 5 ● 45 ● 77

asked Dec 18 '11 at 4:49  
 Macmee 4,826 ● 3 ● 20 ● 55

add a comment

1 Answer

active oldest votes

Gambar 2.5 Ilustrasi Diskusi pada Stack Overflow (Zorychta, 2011)

## 2.7 Evaluasi Sistem Rekomendasi

Karena model sistem rekomendasi yang diteliti adalah rekomendasi berbasis konten, maka sistem dapat dievaluasi dengan menggunakan perhitungan *precision* dan *recall* (Manning, 2008) berdasarkan *query* yang dimasukkan. *Precision* merupakan rasio jumlah dokumen relevan yang dikembalikan terhadap jumlah seluruh hasil yang dikembalikan. Rumusnya sebagaimana berikut.

$$\text{Precision} = \frac{\#(\text{dokumen relevan yang dikembalikan})}{\#(\text{dokumen yang dikembalikan})} \quad (2.4)$$

Sedangkan *Recall* merupakan rasio jumlah dokumen relevan yang dikembalikan terhadap jumlah seluruh dokumen yang relevan. Rumusnya adalah sebagai berikut.

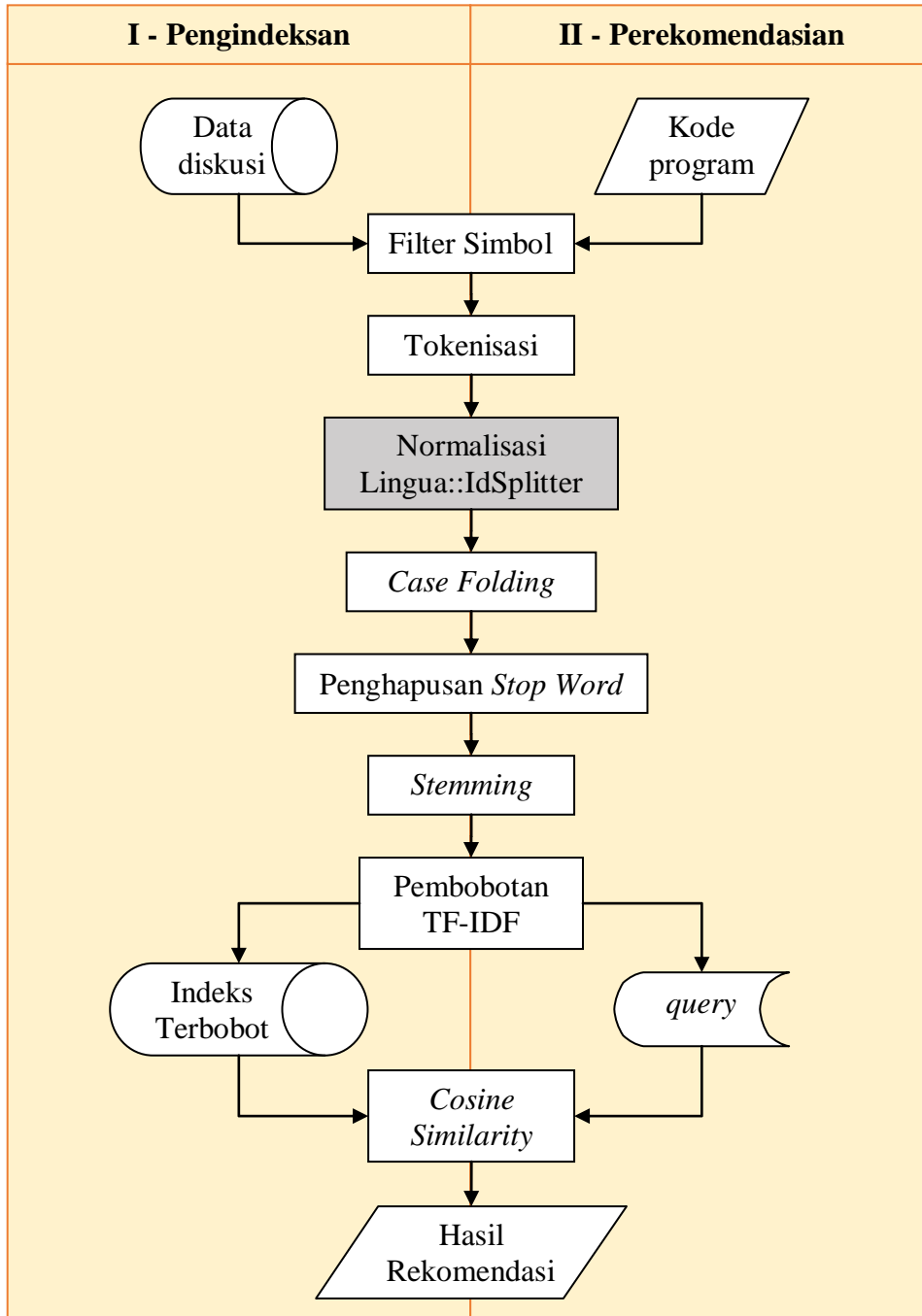
$$\text{Recall} = \frac{\#(\text{dokumen relevan yang dikembalikan})}{\#(\text{dokumen yang relevan})} \quad (2.5)$$

Kedua pengukuran tersebut akan memiliki jangkauan nilai mulai 0 hingga 1 (*floating point*) atau sering pula dirupakan dalam bentuk persentase dengan jangkauan nilai mulai 0 hingga 100.

### BAB 3

## METODE PENELITIAN

Penelitian ini menggunakan kerangka proses sebagaimana yang ditunjukkan pada Gambar 3.1.



Gambar 3.1 Desain Sistem Rekomendasi dengan Normalisasi Lingua::IdSplitter

Sebagaimana tampak pada Gambar 3.1, hasil rekomendasi didapatkan dari dua tahap yaitu tahap pengindeksan yang hanya dijalankan sebanyak satu kali dan tahap perkomendasi data diskusi. Tahap pengindeksan harus dilakukan terlebih dahulu agar sistem rekomendasi dapat dijalankan.

### 3.1 Data Masukan

Masukan pada penelitian ini terbagi menjadi dua jenis, yakni masukan untuk tahap pengindeksan dan masukan untuk tahap perkomendasi. Masukan untuk tahap pengindeksan berupa teks diskusi (pertanyaan dan jawaban) yang diambil dari Stack Overflow. Data diambil dengan bantuan Stack Exchange Data Explorer berdasarkan *query* SQL berikut.

```
SELECT * FROM Posts WHERE Id IN (PostIds) ORDER BY Id
```

Gambar 3.2 *Query* untuk Mengambil Pertanyaan dari Stack Overflow

*PostIds* yang menjadi nilai operator *IN* pada *query* Gambar 3.2 didapatkan dari kumpulan *id post* yang digunakan pada penelitian Arwan dkk. (2015). Kumpulan *id post* tersebut merupakan data pertanyaan dengan label (*tag*) Java yang berjumlah 153 data. Data tersebut terdiri atas 34 data tentang topik *sorting*, 34 data tentang topik *database*, 32 data tentang topik *text file*, 34 data tentang topik *graphics*, dan 19 data tentang topik *thread*. Selain data pertanyaan, jawaban atas pertanyaan tersebut juga diambil. *Query* yang digunakan adalah sebagai berikut.

```
SELECT * FROM Posts WHERE ParentId IN (PostIds) ORDER BY Id
```

Gambar 3.3 *Query* untuk Mengambil Jawaban dari Stack Overflow

*Query* pada Gambar 3.3 mirip dengan *query* pada Gambar 3.2. Yang membedakan keduanya hanya pada kolom yang digunakan untuk memfilter data. Kolom yang digunakan yakni *ParentId*, sementara untuk *query* data pertanyaan menggunakan kolom *Id*. Data yang telah dihasilkan dari *query* tersebut kemudian diunduh dalam bentuk *file* CSV untuk kemudian diimpor secara lokal ke basis data MySQL.

Data yang digunakan hanya data pada kolom *Body* yang berisi pertanyaan atau jawaban dari pengguna. Ilustrasinya dapat dilihat pada Gambar 3.4.

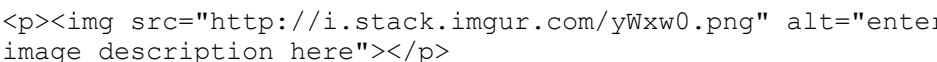
**Body (Pertanyaan), ID Post: 8549891**

How do I turn off text antialiasing in this Java function?

Basically I would like to turn off antialias in the following:

```
public BufferedImage createText(String text) {  
    //create image  
    BufferedImage image = new BufferedImage(95, 20,  
        BufferedImage.TYPE_INT_ARGB);  
    Graphics2D graphics = (Graphics2D) image.getGraphics();  
  
    //set background  
    graphics.setColor(Color.white);  
    graphics.fillRect(0, 0, 95, 20);  
  
    //draw text  
    graphics.setRenderingHint(RenderingHints.KEY_ANTIALIASING,  
        RenderingHints.VALUE_ANTIALIAS_OFF);  
  
    graphics.setColor(Color.black);  
    graphics.setFont(new Font("volter", Font.PLAIN, 9));  
    graphics.drawString(text, 0, 10);  
  
    return image;  
}
```

but it's not working, here is something this function generates:

 enter image description here

I just want black and white to be used, nothing else so it's important I get antialias disabled!

**Body (Jawaban), ID Post: 8549970, ParentID: 8549891**

Try this instead:

```
graphics.setRenderingHint(RenderingHints.KEY_TEXT_ANTIALIASING,  
    RenderingHints.VALUE_TEXT_ANTIALIAS_OFF);
```

Gambar 3.4 Data Masukan Ketika Dibaca Oleh Sistem

Ilustrasi yang ditunjukkan pada Gambar 3.4 merupakan data masukan pada tahap pengindeksan. Kedua jenis informasi tersebut (pertanyaan dan jawaban-jawaban) akan digabung menjadi satu untuk mewakili satu data diskusi. Dapat dilihat bahwa informasi yang akan diindeks terdiri atas tag-tag HTML, hal tersebut

mengingat data diskusi di Stack Overflow berupa dokumen HTML. Gambar tersebut merupakan representasi fisik web Stack Overflow yang telah dijelaskan pada subbab 2.6. Dapat dilihat pula bahwa data tidak hanya berupa kata-kata (bahasa natural) dari pengguna, namun juga berupa potongan kode program yang biasanya diapit dalam tag `<code>`.

Adapun masukan untuk tahap perkomendasi merupakan potongan kode program yang dimasukkan oleh pengguna. Ilustrasinya dapat dilihat pada Gambar 3.3.

```
...
public void paint(Graphics g) {

    Graphics2D g2 = (Graphics2D) g;
    g2.setRenderingHint(RenderingHints.KEY_TEXT_ANTIALIASING,
        RenderingHints.VALUE_TEXT_ANTIALIAS_ON);

    g.drawString("This is my string", 200, 200);
}
...
```

Gambar 3.5 Query Sistem Berupa Kode Program Java

Potongan kode program yang dimaksud pada Gambar 3.5 merupakan potongan kode program dalam bahasa Java (\*.java) yang dijadikan masukan pada tahap perkomendasi.

### 3.2 Filter Simbol

Data masukan yang mengandung simbol-simbol non-alfabet seperti titik koma (;), tanda tanya (?), dan lain-lain akan dihapus sehingga yang tersisa hanya karakter alfabet. Sedangkan tanda titik (.) akan diubah menjadi spasi karena tanda titik digunakan oleh Java untuk memisahkan objek atau *class* dengan membernya. Selain itu, karena data diskusi berupa data HTML, maka *tag-tag* HTML pada diskusi juga akan dihilangkan. Selain mengandung *tag* HTML, dokumen web terkadang juga mengandung simbol yang dirupakan dalam bentuk entitas HTML seperti `&lt;` yang ketika ditampilkan di peramban web akan muncul sebagai simbol *kurang dari* (<). Entitas HTML tersebut juga akan dihapus sebelum data diproses lebih lanjut. Gambar yang disertakan pada diskusi juga akan dihapus mengingat tag `<img>` pada data hanya memuat lokasi gambar saja pada server. Ilustrasi data yang telah difilter dapat dilihat pada Gambar 3.4.

```

How do I turn off text antialiasing in this Java function

Basically I would like to turn off antialias in the following

public BufferedImage createTextString text

    create image
    BufferedImage image new BufferedImage
        BufferedImageTYPE_INT_ARGB
    Graphics2D graphics Graphics2D imagegetGraphics

    set background
    graphicssetColorColorwhite
    graphicsfillRect

    draw text
    graphicssetRenderingHintRenderingHintsKEY_ANTIALIASING
        RenderingHintsVALUE_ANTIALIAS_OFF

    graphicssetColorColorblack
    graphicssetFontnew FontvolterFontPLAIN
    graphicsdrawStringtext

    return image

but its not working here is something this function generates

I just want black and white to be usednothing else so its
important I get antialias disabled

Try this instead

GraphicssetRenderingHintRenderingHints KEY_TEXT_ANTIALIASING
        RenderingHintsVALUE_TEXT_ANTIALIAS_OFF

```

Gambar 3.6 Hasil Filter Simbol

Dapat dilihat pada Gambar 3.6, yang akan tersisa pada data hanya karakter-karakter alfabet saja. Semua tag HTML sudah tidak ada. Dengan kata lain, dokumen HTML akan diubah bentuknya menjadi teks *plain*.

### 3.3 Tokenisasi

Datayang telahdifilter kemudian ditokenisasi sehingga akan didapatkan daftar term.Model tokenisasi yang digunakan adalah *word token*, yakni pemisahan term berdasarkan karakter spasi (*white space*) sebagaimana tampak padaGambar 3.7.

...	create	Graphics2D	setColor
-----	--------	------------	----------

The	image	graphics	Color
following	BufferedImage	Graphics2D	white
public	image	image	graphics
BufferedImage	new	getGraphics	fillRect
createText	BufferedImage	set	...
String	BufferedImage	background	
Text	TYPE INT ARGB	graphics	

Gambar 3.7 Hasil Tokenisasi

Dapat dilihat pada Gambar 3.7 tersebut bahwa setelah proses tokenisasi dilakukan, akan didapatkan daftar term yang dapat digunakan untuk proses berikutnya.

### 3.4 Normalisasi Term dengan Lingua::IdSplitter dan Case Folding

Term-term hasil tokenisasi akan dinormalisasi dengan Lingua::IdSplitter. Normalisasi ini akan mengubah *identifier* yang umumnya mengandung gabungan kata serta mengekspansi singkatan menjadi kata lengkap. Term yang memiliki variasi huruf besar dan huruf kecil selanjutnya akan diubah semuanya ke huruf kecil. Hal ini dimaksudkan untuk menyamaratakan semua karakter ke bentuk huruf kecil. Ilustrasinya dapat dilihat pada Gambar 3.8.

...	...
The	the
following	following
public	public
BufferedImage	buffered
createText	image
String	create
...	text
graphics	string
setColor	...
Color	graphics
white	set
graphics	color
fillRect	white
...	graphics
	fill
	rectangle
	...
Sebelum	Sesudah

Gambar 3.8 Hasil Sebelum dan Sesudah Pemrosesan dengan Lingua::IdSplitter

Dapat dilihat pada Gambar 3.8 bahwa term `BufferedImage` akan dapat berubah menjadi `buffered` dan `image`, `createText` dapat berubah menjadi `create` dan `text`, `setColor` dapat berubah menjadi `set` dan `color`, serta term `fillRect` juga dapat berubah menjadi `fill` dan `rectangle`.



Mengingat `Lingua::IdSplitter` merupakan modul Perl, maka pada penelitian ini, `Lingua::IdSplitter` akan difungsikan sebagai program eksternal yang akan dieksekusi pada aplikasi berbasis Java yang dikembangkan dengan bantuan objek `Runtime` dan `Process`. Adapun format perintahnya dapat dilihat pada Gambar 3.9.

```
perl id-splitter termTarget
```

Gambar 3.9 Format Eksekusi `Lingua::IdSplitter`

Sebagaimana tampak pada Gambar 3.9, `term` yang akan dinormalisasi (`termTarget`) menjadi parameter `Lingua::IdSplitter(id-splitter)` yang pada akhirnya akan menghasilkan  $\geq 1$  term hasil normalisasi.

### 3.5 Penghapusan *Stop Word*

Daftar term yang telah diperoleh selanjutnya difilter kembali untuk menghilangkan term-term yang merupakan term umum (*stop word*) seperti term *the*, *a*, *or*, dan lain-lain sebagaimana tampak pada Gambar 3.10.

<pre>... the following public buffered image create text string ... graphics set color graphics fill rectangle ...</pre>	<pre>... following public buffered image create text string ... graphics set color graphics fill rectangle ...</pre>
Sebelum	Sesudah

Gambar 3.10 Hasil Sebelum dan Sesudah Penghapusan *Stop Word*

Sebagaimana tampak pada Gambar 3.10, term *the* akan dihapus dari daftar term karena termasuk salah satu term dalam *stop list* sebagaimana tercantum dalam Lampiran 2.

### 3.6 Stemming

Daftar term yang telah difilter dengan *stop word* selanjutnya akan diubah lagi bentuknya dengan proses *stemming* untuk mengurangi variasi kata yang memiliki bentuk dasar (stem). Ilustrasi proses *stemming* dapat dilihat pada Gambar 3.11.

<pre> ... following public buffered image create text string ... graphics set color white graphics fill rectangle ... </pre>	<pre> ... follow public buffer imag creat text string ... graphic set color white graphic fill rectangl ... </pre>
Sebelum	Sesudah

Gambar 3.11 Hasil Sebelum dan Sesudah *Stemming*

Karakter yang disorot dengan warna merah pada Gambar 3.11 adalah karakter yang akan dihapus ketika dilakukan *stemming*.

### 3.7 Pembobotan TF-IDF

Term yang telah didapatkan selanjutnya dihitung frekuensi kemunculannya serta diboboti dengan menggunakan TF-IDF berdasarkan rumus (2.1). Dengan demikian, akan didapatkan daftar term unik beserta bobotnya terhadap data diskusi. Hasil pembobotan ini akan menghasilkan fitur vektor untuk tiap-tiap data diskusi. Adapun ilustrasinya dapat dilihat pada Gambar 3.12.  $D_1$ ,  $D_2$ ,  $D_3$  yang ditunjukkan pada Gambar 3.12 merepresentasikan data diskusi dengan ID 1, 2, 3. Term *color*, *white*, dan seterusnya merupakan term unik yang berasal dari hasil *stemming*. Sedangkan nilai yang berada pada tabel merupakan pemetaan bobot term terhadap data diskusi.

	$D_1$	$D_2$	$D_3$	...
...	...	...	...	...

<i>color</i>	0,0006	0,0001	0,0003	...
<i>white</i>	0	0	0	...
<i>graphic</i>	1,0344	1,0344	13,3255	...
<i>fill</i>	0	0	0	...
<i>rectangl</i>	0,0003	0	0	...
...	...	...	...	...

Gambar 3.12 Contoh Daftar Term Beserta Bobotnya dalam Bentuk Matriks Term-Dokumen

### 3.8 Cosine Similarity

Penghitungan *Cosinesimilarity* digunakan untuk membandingkan antara *query* yang dimasukkan dengan data diskusi yang sudah tersedia dalam bentuk indeks terbobot (vektor). Selanjutnya hasil perbandingan akan diperingkat berdasarkan tingkat kemiripannya menggunakan rumus (2.3).

### 3.9 Evaluasi Hasil

Evaluasi pada penelitian ini dilakukan dengan tujuan untuk mengetahui kualitas perbaikan metode rekomendasi yang memanfaatkan proses normalisasi *identifier* menggunakan *Lingua::IdSplitter* melalui penghitungan nilai *precision* dan *recall* berdasarkan rumus (2.4) dan (2.5). Nilai-nilai tersebut akan dibandingkan dengan nilai *precision* maupun *recall* yang dihasilkan oleh sistem rekomendasi yang tidak memanfaatkan proses normalisasi *identifier*. Mengingat data yang digunakan berasal dari penelitian Arwan dkk. (2015), maka hasil penelitian ini juga akan dilihat keselarasannya dengan penelitian yang telah dilakukan oleh Arwan dkk. tersebut. Selain itu, guna mengetahui kemampuan *Lingua::IdSplitter* dalam menormalisasi *identifier*, akan digunakan pula data diskusi alternatif sejumlah 80 buah (20 diskusi/topik) dengan *identifier* singkatan yang lebih mendominasi dibandingkan dengan singkatan pada data diskusi yang digunakan pada penelitian Arwan dkk. (2015). Dominasi *identifier* singkatan tersebut muncul terutama sebagai *identifier-identifier* yang merupakan *identifier* atau term kunci topik.

Data uji (*query*) masukan pada sistem ini adalah kode program Java yang berjumlah 32 buah. Kode program tersebut diambil secara acak dari *repository* GitHub dengan spesifikasi sebagaimana tampak pada Tabel 3.1.

Tabel 3.1 Spesifikasi Data Uji

<b>Topik Kode Program</b>	<b>Hasil Diskusi yang Diharapkan</b>	<b>Jml. <i>Query</i></b>
<i>Thread synchronization</i>	Diskusi yang membahas tentang sinkronisasi <i>thread</i>	8
<i>AccessingMySQL</i>	Diskusi yang membahas tentang pengaksesan data pada MySQL	8
<i>Bubble sort</i>	Diskusi yang membahas tentang metode <i>bubble sort</i>	8
<i>Draw rectangle</i>	Diskusi yang membahas tentang cara menggambar persegi panjang	8
<b>Total <i>Query</i></b>		<b>32</b>

Setiap kode program yang dimasukkan, diharapkan mampu menghasilkan rekomendasi diskusi sesuai topik dari kode program sebagaimana yang telah disajikan pada Tabel 3.1. Detail kode-kode program yang dimaksud dapat dilihat pada Lampiran 3. Sedangkan untuk hasil rekomendasi yang akan dihitung performanya adalah 15 data dengan nilai bobot dokumen tertinggi pertama sebagaimana yang dilakukan oleh Ponzanelli (2013).

## BAB 4

### HASIL DAN PEMBAHASAN

#### 4.1 Proses Rekomendasi

Proses rekomendasi dilakukan dengan cara melakukan *query* terhadap sistem menggunakan kode program yang telah didapatkan. *Query* yang dimaksud adalah kumpulan kata kunci yang diperoleh dari kode program yang telah diproses sebelumnya sesuai metode yang digunakan. Gambar 4.1 merupakan contoh kode program terkait topik *graphics* yang digunakan saat proses rekomendasi.

```
/*
 * File: DrawRectangle.java
 * -----
 * This program allows users to create rectangles on the canvas
 * by clicking and dragging with the mouse.
 */

import java.awt.event.*;
import acm.graphics.*;
import acm.program.*;

/** This class allows users to drag rectangles on the canvas */
public class DrawRectangle extends GraphicsProgram {

    /** Runs the program */
    public void run() {
        addMouseListeners();
    }

    /** Called on mouse press to record the starting coordinates */
    public void mousePressed(MouseEvent e) {
        startX = e.getX();
        startY = e.getY();
        lastX = e.getX();
        lastY = e.getY();
        gobj = getElementAt(startX, startY);
        if (gobj == null) {
            currentRect = new GRect(startX, startY, 0, 0);
            currentRect.setFilled(true);
            add(currentRect);
        }
    }

    /** Called on mouse drag to reshape the current rectangle */
    public void mouseDragged(MouseEvent e) {
        if (gobj != null) {
            gobj.move(e.getX() - lastX, e.getY() - lastY);
            lastX = e.getX();
        }
    }
}
```

```

        lastY = e.getY();
    } else {
        double x = Math.min(e.getX(), startX);
        double y = Math.min(e.getY(), startY);
        double width = Math.abs(e.getX() - startX);
        double height = Math.abs(e.getY() - startY);
        currentRect.setBounds(x, y, width, height);
    }
}

/* Private state */
private GRect currentRect; /* The current rectangle */
private double startX; /* The initial mouse X position */
private double startY; /* The initial mouse Y position */
private GObject gobj;
private double lastX;
private double lastY;
}

```

Gambar 4.1 Contoh Kode Program Sebelum Diekstrak Kata Kuncinya

Contoh kode program yang tampak pada Gambar 4.1 merupakan kode program Java yang mengandung berbagai unsur seperti *class*, objek, variabel, *method*, komentar, dan lainnya. Kode program tersebut kemudian diekstrak untuk mendapatkan kata-kata kunci dengan menggunakan metode yang sama dengan metode yang digunakan ketika proses pengindeksan. Hasil ekstraksi kata kunci, baik yang memanfaatkan normalisasi *identifier* dan yang tidak, dapat dilihat pada Gambar 4.2 dan Gambar 4.3.

```

file draw rectangl java program user creat rectangl canva click
drag mous import java awt event import acm graphic import acm
program class user drag rectangl canva public class draw
rectangl extend graphic program run program public void run add
mous listen call mous press record start coordin public void
mous press mous event start start object element start start
object null current rectangl rectangl start start current
rectangl set fill true add current rectangl call mous drag
reshap current rectangl public void mous drag mous event object
null object move doubl math min start doubl math min start doubl
width math ab start doubl height math ab start current rectangl
set bound width height privat state privat rectangl current
rectangl current rectangl privat doubl start initi mous posit
privat doubl start initi mous posit privat object object privat
doubl privat doubl

```

Gambar 4.2 Kata Kunci Hasil Ekstraksi dengan Memanfaatkan Normalisasi *Identifier*

```

file drawrectangl java program user creat rectangl canva click
drag mous import java awt event import acm graphic import acm
program class user drag rectangl canva public class drawrectangl
extend graphicsprogram run program public void run
addmouselisten call mous press record start coordin public void
mousepress mouseev startx getx starti geti lastx getx lasti geti
gobj getelementat startx starti gobj null currentrect grect
startx starti currentrect setfil true add currentrect call mous
drag reshap current rectangl public void mousedrag mouseev gobj
null gobj move getx lastx geti lasti lastx getx lasti geti doubl
math min getx startx doubl math min geti starti doubl width math
ab getx startx doubl height math ab geti starti currentrect
setbound width height privat state privat grect currentrect
current rectangl privat doubl startx initi mous posit privat
doubl starti initi mous posit privat gobject gobj privat doubl
lastx privat doubl lasti

```

Gambar 4.3 Kata Kunci Hasil Ekstraksi Tanpa Normalisasi *Identifiser*

Dapat dilihat pada Gambar 4.2, kata kunci yang dihasilkan melalui proses normalisasi *identifiser* yakni sebanyak 142 term (48 term unik). Sedangkan yang tampak pada Gambar 4.3, kata kunci hasil ekstraksi tanpa normalisasi sebanyak 146 term (60 term unik). Dapat dilihat pula, kumpulan kata kunci yang dihasilkan melalui proses normalisasi cenderung lebih mirip bahasa alami dibandingkan dengan kumpulan kata kunci tanpa normalisasi. Hal tersebut disebabkan karena *identifiser* yang umumnya terdiri atas gabungan kata sudah terpisah menjadi kata-kata tunggal. Selain itu, term singkatan juga dapat diubah ke bentuk term lengkapnya. Kata kunci yang telah didapat selanjutnya akan digunakan sebagai *query* untuk mendapatkan data-data diskusi yang relevan dengan kode program tersebut.

Adapun proses rekomendasinya dapat dilihat pada Gambar 4.4. Terlihat pada gambar bahwa proses rekomendasi dapat dilakukan dengan atau tanpa normalisasi dengan memilih mode pada kolom pilihan *Search mode* di aplikasi yang dikembangkan. Kode program yang akan dijadikan sebagai *query* dapat dimasukkan dengan menekan tombol *Select file*. Sistem akan memberikan rekomendasi maksimum sebanyak 15 diskusi. Antara kode program dan hasil rekomendasi dapat diketahui kecocokan topiknya berdasarkan informasi yang ada pada bagian kode program dan keterangan pada tiap-tiap judul diskusi. Selain itu, dapat diketahui pula skor similaritas tiap-tiap diskusi dengan kode program yang dimasukkan.

## Get the recommendations

Search mode:

Showing 15 of 153 results (5 relevant posts).  

RCoon.DrawRectangle.java | Topic: Draw rectangle

```

/*
 * File: DrawRectangle.java
 * -----
 * This program allows users to create rectangles on the canvas
 * by clicking and dragging with the mouse.
 */

import java.awt.event.*;
import acm.graphics.*;
import acm.program.*;

/** This class allows users to draw rectangles on the canvas */

```

[Difference between Graphics object of getGraphics and paintComponent](#)

Score: 1.533

[Circles in circles fading out, creating a target](#)

Score: 1.103

[Rotate an image in java by the specified angle](#)

Score: 1.020

[Mouse Listener Interface and painting](#)

Score: 0.8072

[Compare sorting algorithm](#)

Score: 0.7382, Topic: Bubble sort

[Drawing a nice circle in Java](#)

Score: 0.6618

[how to group graphic object](#)

Score: 0.5905, Topic: Draw rectangle

Gambar 4.4 Proses Rekomendasi

Berdasarkan contoh hasil rekomendasi pada Gambar 4.4, diketahui bahwa diskusi dengan skor tertinggi ternyata tidak terkait dengan topik kode program (*draw rectangle*). Diskusi yang terkait dengan topik tersebut berada pada urutan ke-7 dengan skor 0,5905. Tampak pada gambar pula bahwa diskusi relevan yang direkomendasikan adalah sejumlah 5 diskusi, sementara 10 lainnya tidak relevan.

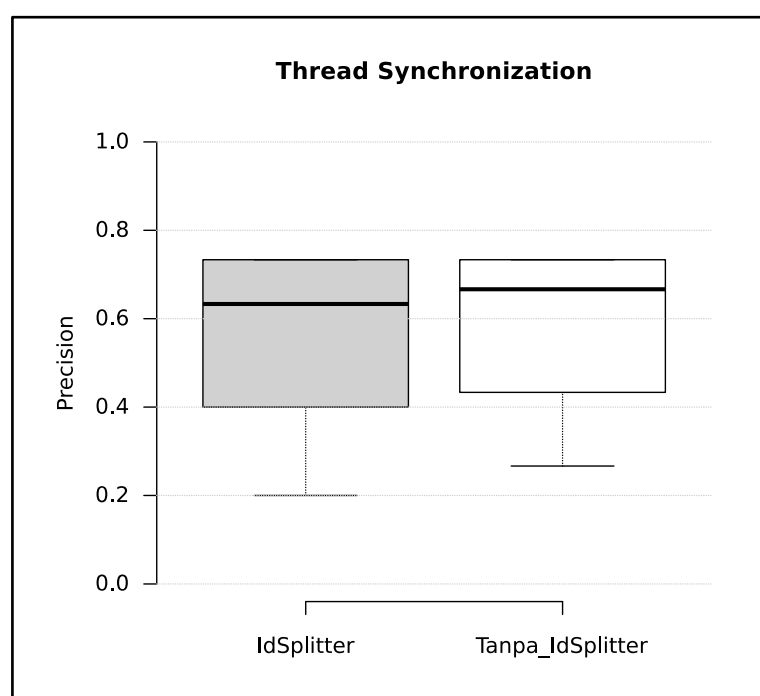


## 4.2 Hasil Pengujian

Sebagaimana telah disebutkan pada subbab 3.9, pengujian akan dilakukan dengan melakukan *query* terhadap sistem menggunakan kode sumber program Java dengan topik yang telah ditentukan. Detail hasil pengujian dapat dilihat pada Lampiran 4 dan Lampiran 5.

### 4.2.1 Pengujian dengan Kode Program Terkait *Thread Synchronization*

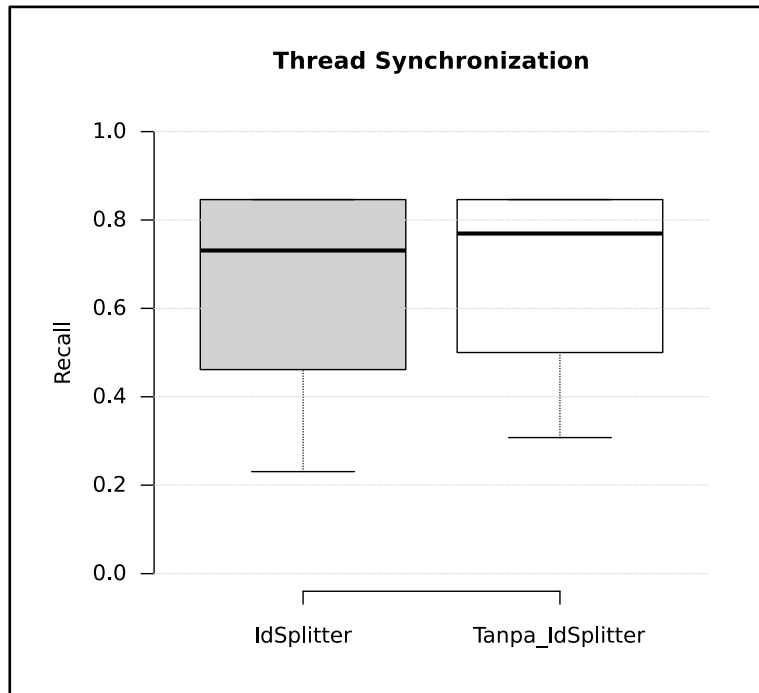
Pengujian terhadap 8 kode program yang terkait dengan topik *thread synchronization* dapat dilihat pada Gambar 4.5 dan Gambar 4.6.



Gambar 4.5 Perbandingan *Precision* untuk Topik *Thread Synchronization*

Berdasarkan grafik pada Gambar 4.5, dapat dilihat bahwa nilai *precision* untuk data yang dinormalisasi dengan `Lingua::IdSplitter` memiliki jangkauan antara 0,2 hingga 0,73 dengan median 0,63 sementara data yang tidak dinormalisasi memiliki jangkauan antara 0,27 hingga 0,73 dengan median 0,67. Sedangkan apabila dilihat dari nilai *recall*, sebagaimana tampak pada Gambar 4.6, diketahui bahwa *recall* yang diperoleh ketika data dinormalisasi yakni antara 0,23 hingga 0,85 dengan median 0,73 sementara data yang tidak dinormalisasi yakni antara 0,31 hingga 0,85 dengan median 0,77. Dengan demikian, dapat disimpulkan

bahwa proses normalisasi tidak dapat meningkatkan nilai *precision* maupun *recall* untuk *query* dengan topik *thread synchronization*.

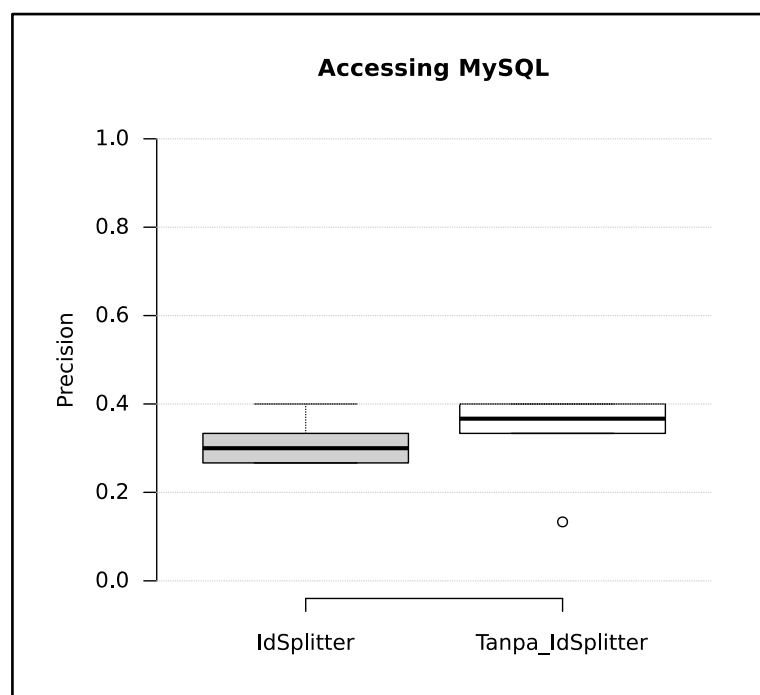


Gambar 4.6 Perbandingan *Recall* untuk Topik *Thread Synchronization*

Hal ini disebabkan oleh tidak ditemukannya *identifier synchronize* yang disingkat menjadi *sync* pada data diskusi sebagaimana yang lumrah digunakan oleh *programmer* dalam menyingkat *identifier*. Oleh karena itu, normalisasi (ekspansi) term tidak berdampak terhadap kasus *query* dengan topik *thread synchronization* ini. Selain itu, sinkronisasi *thread* juga umum menggunakan *exception* dengan jenis `InterruptedException`. Pemisahan *identifier* tersebut menjadi `InterruptedException` justru menurunkan skor dokumen terkait *thread synchronization* karena term `exception` lumrah ditemukan pada topik-topik selainnya. Sementara itu, ekspansi term dalam hal ini tidak terlalu berpengaruh mengingat pada kode program dan data sampel diskusi yang digunakan tidak ada yang menggunakan term `sync` yang biasanya digunakan untuk mewakili term `synchronize` yang merupakan term penting pada topik *thread synchronization*.

#### 4.2.2 Pengujian dengan Kode Program Terkait *Accessing MySQL*

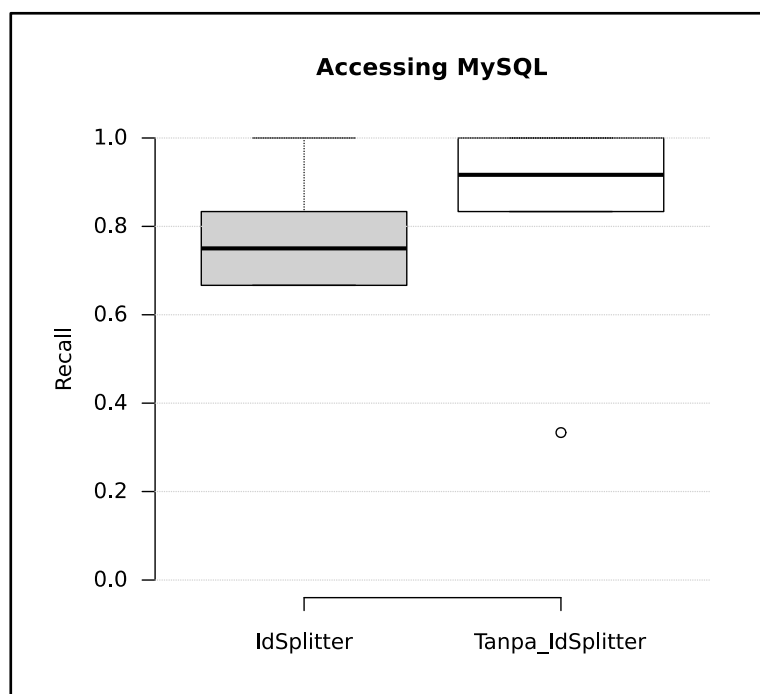
Pengujian terhadap 8 kode program yang terkait dengan topik *accessing MySQL* dapat dilihat pada Gambar 4.7 dan Gambar 4.8.



Gambar 4.7 Perbandingan *Precision* untuk Topik *Accessing MySQL*

Berdasarkan grafik pada Gambar 4.7, dapat dilihat bahwa nilai *precision* untuk data yang dinormalisasi dengan `Lingua::IdSplitter` memiliki jangkauan antara 0,27 hingga 0,4 dengan median 0,3 sementara data yang tidak dinormalisasi memiliki jangkauan antara 0,33 hingga 0,4 dengan median 0,37. Sementara apabila dilihat dari nilai *recall*, sebagaimana tampak pada Gambar 4.8, diketahui bahwa *recall* yang diperoleh ketika data dinormalisasi yakni antara 0,23 hingga 1 dengan median 0,75 sementara data yang tidak dinormalisasi yakni antara 0,83 hingga 1 dengan median 0,92. Dengan demikian, dapat disimpulkan bahwa proses normalisasi juga tidak dapat meningkatkan nilai *precision* maupun *recall* untuk *query* dengan topik *accessing MySQL*. Hal ini antara lain disebabkan oleh hasil ekspansi *identifier jdbc* menjadi `j, database, dan c`. Perluasan tersebut justru membuat direkomendasikannya data-data diskusi lain yang mengandung term `database` meski tidak memiliki term `jdbc` yang menjadi ciri umum bahasan tentang pengaksesan `MySQL`. Selain itu, pemisahan *identifier mysql* menjadi `my`

dan `sql` juga turut menurunkan performa rekomendasi. Pemisahan tersebut justru mengakibatkan hilangnya term penting `mysql` dalam mewakili data diskusi terkait *accessing* MySQL.

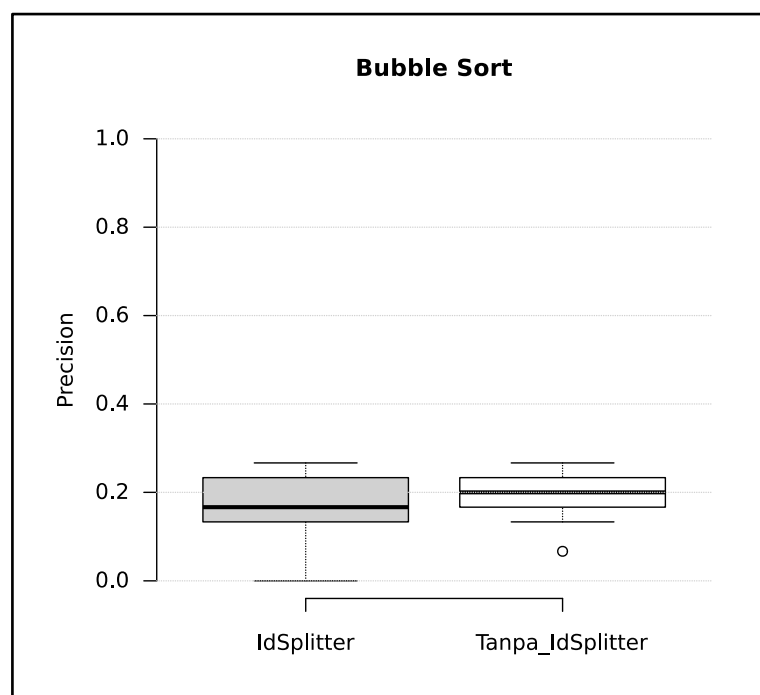


Gambar 4.8 Perbandingan *Recall* untuk Topik *Accessing* MySQL

Normalisasi (pemisahan) dapat memberikan peningkatan hasil rekomendasi untuk *query* kode program `org.embulk.output.MySQLOutputPlugin.java` yang berisi tentang kode program untuk melakukan *insert* ke basis data MySQL secara massal. Tanpa normalisasi, sistem menjadikan *identifier* seperti `batchInsert`, `mysqlBatchInsert`, dan `newBatchInsert` sebagai term penting *query*. Secara logika, term-term tersebut memang tepat mewakili *query*. Namun, hal tersebut tidak didukung oleh sampel data diskusi, tidak ada satupun topik data diskusi yang membahas tentang proses *insert* secara massal. Jika demikian, setidaknya sistem bisa merekomendasikan data diskusi tentang proses *insert* pada MySQL secara umum. Akan tetapi, kenyataannya tidak. Efek tersebut dapat dilihat pada penculan nilai *precision* pada titik 0,13 dan *recall* pada titik 0,33 untuk proses rekomendasi yang tidak menggunakan normalisasi `Lingua::IdSplitter`.

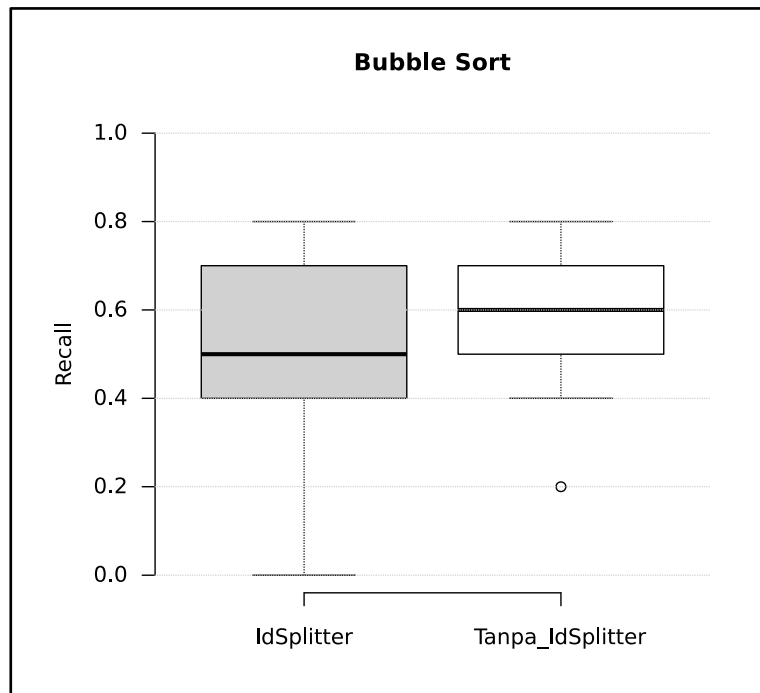
### 4.2.3 Pengujian dengan Kode Program Terkait *Bubble Sort*

Pengujian terhadap 8 kode program yang terkait dengan topik *bubble sort* dapat dilihat pada Gambar 4.9 dan Gambar 4.10.



Gambar 4.9 Perbandingan *Precision* untuk Topik *Bubble Sort*

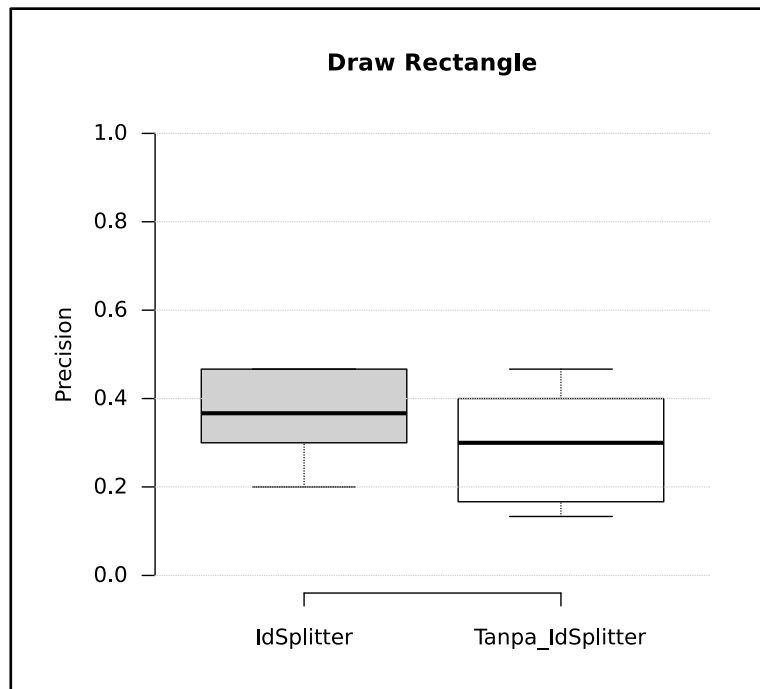
Grafik pada Gambar 4.9 menunjukkan bahwa nilai *precision* untuk data yang dinormalisasi dengan `Lingua::IdSplitter` memiliki jangkauan antara 0 hingga 0,27 dengan median 0,17 sementara data yang tidak dinormalisasi memiliki jangkauan antara 0,13 hingga 0,27 dengan median 0,2. Sementara apabila dilihat dari nilai *recall*, sebagaimana tampak pada Gambar 4.10, diketahui bahwa *recall* yang diperoleh ketika data dinormalisasi yakni antara 0 hingga 0,8 dengan median 0,5 sementara data yang tidak dinormalisasi yakni antara 0,4 hingga 0,8 dengan median 0,6. Dapat disimpulkan bahwa proses normalisasi juga tidak dapat meningkatkan nilai *precision* maupun *recall* untuk *query* dengan topik *bubble sort*. Hal ini antara lain disebabkan oleh hasil pemisahan *identifier* `BubbleSort` menjadi `bubble` dan `sort`. Pemisahan tersebut justru membuat meningkatnya skor data diskusi tentang *sorting* yang lain, yang menggunakan metode selain *bubble sort*. Ekspansi pada topik ini juga tidak terlalu berpengaruh mengingat *identifier* `bubble` maupun `sort` jarang dirupakan dalam bentuk singkat.



Gambar 4.10 Perbandingan *Recall* untuk Topik *Bubble Sort*

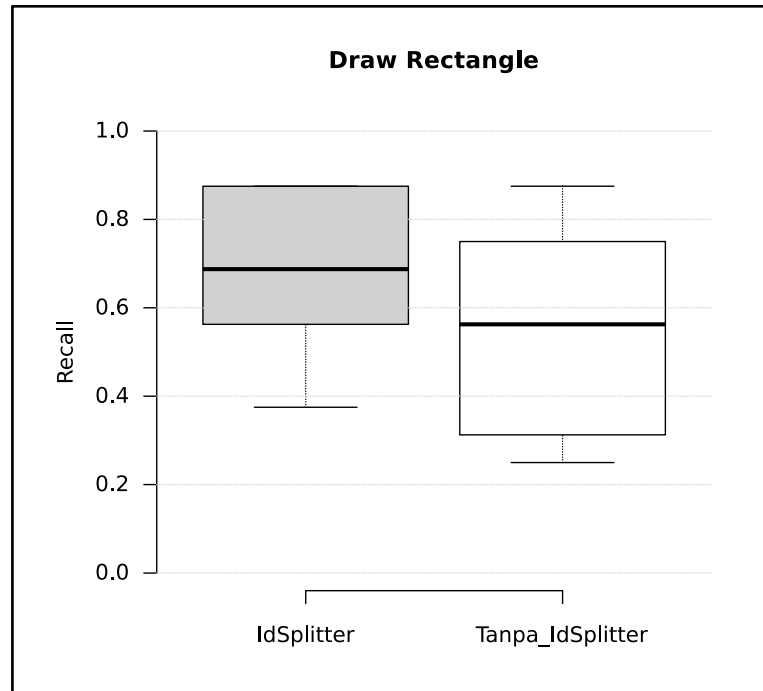
#### 4.2.4 Pengujian dengan Kode Program Terkait *Draw Rectangle*

Pengujian terhadap 8 kode program yang terkait dengan topik *draw rectangle* dapat dilihat pada Gambar 4.11 dan Gambar 4.12.



Gambar 4.11 Perbandingan *Precision* untuk Topik *Draw Rectangle*

Berdasarkan grafik pada Gambar 4.11, dapat dilihat bahwa nilai *precision* untuk data yang dinormalisasi dengan `Lingua::IdSplitter` memiliki jangkauan antara 0,2 hingga 0,47 dengan median 0,37 sementara data yang tidak dinormalisasi memiliki jangkauan antara 0,13 hingga 0,47 dengan median 0,3.



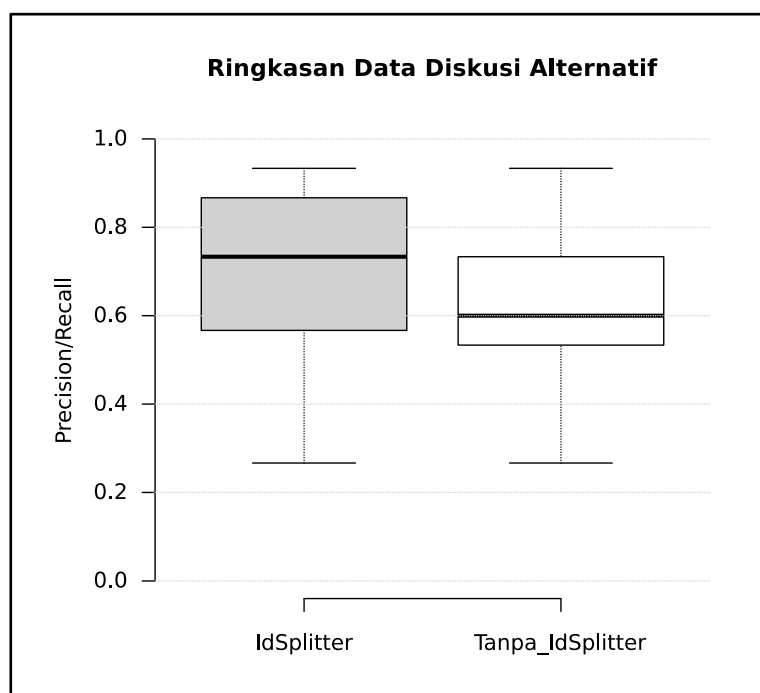
Gambar 4.12 Perbandingan *Recall* untuk Topik *Draw Rectangle*

Sementara apabila dilihat dari nilai *recall*, sebagaimana tampak pada Gambar 4.12, diketahui bahwa *recall* yang diperoleh ketika data dinormalisasi yakni antara 0,38 hingga 0,88 dengan median 0,69 sementara data yang tidak dinormalisasi yakni antara 0,25 hingga 0,88 dengan median 0,56. Dengan demikian, dapat disimpulkan bahwa proses normalisasi dapat meningkatkan nilai median *precision* maupun *recall* masing-masing sebesar 0,07 dan 0,13 untuk *query* dengan topik *draw rectangle*. Hal ini antara lain disebabkan oleh pemisahan *identifier* `DrawFeature` menjadi `Draw` dan `Feature`, `DrawRectangle` menjadi `Draw` dan `Rectangle`, serta `DrawRectangleOptions` menjadi `Draw`, `Rectangle`, dan `Options`. *Identifier* dalam bentuk terpisah lebih banyak ditemui pada data diskusi. Oleh karena itu, proses normalisasi menjadi penting. Selain itu, pemisahan dan ekspansi *identifier* `fillRect` menjadi `fill` dan `Rectangle` serta `drawRect` menjadi `draw` dan `Rectangle` memberikan dampak positif terhadap hasil

rekomendasi, mengingat teks diskusi yang ada lebih banyak menggunakan kata lengkap dibandingkan singkatan.

### 4.3 Uji Coba Menggunakan Data Alternatif

Hasil uji coba menggunakan data alternatif dengan dominasi *identifier* yang mengandung singkatan dapat dilihat pada Gambar 4.13. Pengujian ini menghasilkan nilai-nilai *precision* dan *recall* yang sama, sehingga representasinya hanya ditampilkan dalam satu grafik saja.



Gambar 4.13 Ringkasan Perbandingan Uji Coba Data Alternatif

Berdasarkan grafik pada Gambar 4.13 dapat diketahui bahwa nilai *precision* dan *recall* yang dapat dihasilkan sama, yakni berada pada jangkauan 0,27 hingga 0,93. Namun, median nilai *precision* dan *recall* hasil rekomendasi yang memanfaatkan normalisasi Lingua::IdSplitter lebih tinggi 0,13 dibandingkan dengan tanpa memanfaatkan proses normalisasi.



#### 4.4 Ringkasan Hasil Uji Coba

Berdasarkan hasil pengujian yang telah dijelaskan untuk tiap-tiap topik pada subbab 4.2.1 sampai dengan 4.2.4, diketahui bahwa proses normalisasi dengan `Lingua::IdSplitter` hanya berpengaruh positif terhadap rekomendasi *query* dengan topik *draw rectangle*. Hal tersebut disebabkan oleh keberadaan *identifier<sub>rectangle</sub>* pada data diskusi maupun *query*, baik berupa singkatan maupun term lengkap. Sedangkan berdasarkan uji coba terhadap data alternatif yang telah dipilih, diketahui bahwa normalisasi *identifier* menggunakan `Lingua::IdSplitter` juga berpengaruh positif terhadap hasil rekomendasi dilihat dari nilai median yang dihasilkan. Hal tersebut juga dipengaruhi oleh keberadaan *identifier* pada data diskusi yang ditulis dalam bentuk singkatan.

Adapun apabila dilihat keselarasannya dengan penelitian Arwan dkk. (2015) yang menggunakan model topik dalam merekomendasikan data diskusi, penggunaan model bobot TF-IDF masih cukup bisa diandalkan dengan melihat nilai median *recall* secara umum yang bisa dihasilkan oleh sistem berbasis TF-IDF pada data uji yang sama, yakni sebesar 0,67. Namun demikian, perbaikan masih perlu dilakukan mengingat nilai median *precision* umum yang dihasilkan hanya mampu mencapai 0,33. Hal ini diakibatkan oleh ketidaktepatan pemilihan *identifier-identifier* yang layak untuk dinormalisasi, baik terkait kasus ketidaktepatan pemisahan komposisi term maupun pengeksiansian term penuh dari term singkatan sebagaimana yang telah dijelaskan pada subbab-subbab sebelumnya.

*Halaman ini sengaja dikosongkan*

# Lampiran 1

## Kamus Kata pada Lingua::IdSplitter

### A. Kamus Akronim

acl-----acl	ftp-----ftp	rgb-----rgb
aes-----aes	gdk-----gdk	rhp-----rhp
ansi-----ansi	gid-----gid	rrt-----rrt
api-----api	gmp-----gmp	rtc-----rtc
arf-----arf	gnu-----gnu	rtx-----rtx
arm-----arm	gnu-----gnu	scsi-----scsi
asn-----asn	gtk-----gtk	sdp-----sdp
avp-----avp	hci-----hci	sgc-----sgc
bfd-----bfd	hdr-----hdr	sha-----sha
bhcp-----bhcp	href-----href	smob-----smob
bpf-----bpf	html-----html	smr-----smr
bsd-----bsd	http-----http	sql-----sql
bsd-----bsd	https-----https	ssl-----ssl
bsd-----bsd	ieee-----ieee	stderr-----stderr
bzip-----bzip	ifp-----ifp	stdout-----stdout
cdb-----cdb	inode-----inode	stt-----stt
cdid-----cdid	insn-----insn	stv-----stv
cdrom-----cdrom	ipc-----ipc	svm-----svm
cgi-----cgi	irq-----irq	swf-----swf
cgm-----cgm	iso-----iso	tcl-----tcl
codec-----codec	jit-----jit	tcp-----tcp
conexant-----conexant	jpeg-----jpeg	tipc-----tipc
cpu-----cpu	js-----js	tlv-----tlv
csi-----csi	kld-----kld	tms-----tms
css-----css	mav-----mav	toc-----toc
cst-----cst	mfi-----mfi	tty-----tty
cvs-----cvs	mips-----mips	udp-----udp
dbsm-----dbsm	mmx-----mmx	uid-----uid
dejagnu-----dejagnu	mpn-----mpn	uni-----uni
dll-----dll	mutex-----mutex	uri-----uri
dma-----dma	ncp-----ncp	url-----url
dom-----dom	nntp-----nntp	usb-----usb
dpb-----dpb	nss-----nss	utf-----utf
dsa-----dsa	pci-----pci	vma-----vma
dsp-----dsp	pcm-----pcm	vms-----vms
ebp-----ebp	pcr-----pcr	wais-----wais
elf-----elf	pkix-----pkix	wsa-----wsa
eoas-----eoas	plt-----plt	www-----www
eob-----eob	png-----png	xdr-----xdr
evt-----evt	pop-----pop	xfrm-----xfrm
fdc-----fdc	psp-----psp	xml-----xml
fdt-----fdt	psppire-----psppire	xpt-----xpt
frv-----frv	pwd-----pwd	zip-----zip

Jumlah: 129 kata

## B. Kamus Singkatan (Umum)

abbr -----	abbreviation	eval -----	evaluate	opt -----	option
abbrev -----	abbreviation	evt -----	event	opts -----	options
abs -----	absolute	exec -----	execute	params ----	parameters
addr -----	address	exec -----	execution	paren -----	parenthesis
adr -----	address	exp -----	exponent	perm -----	permission
algo -----	algorithm	expr -----	expression	pg -----	page
alloc -----	allocate	expt -----	exponent	pic -----	picture
alog -----	analog	ext -----	extend	pict -----	picture
alt -----	alternative	fail -----	failure	pntr -----	pointer
arg -----	argument	flg -----	flag	pos -----	position
args -----	arguments	fmt -----	format	prcl -----	protocol
attr -----	attribute	frag -----	fragment	priv -----	private
auth -----	authenticate	frm -----	frame	proc -----	procedure
aux -----	auxiliary	frmt -----	format	proc -----	process
avg -----	average	func -----	function	procs -----	process
blk -----	block	gbl -----	global	prods -----	products
bool -----	boolean	glob -----	global	ptr -----	pointer
buf -----	buffer	graph -----	graphic	rand -----	random
buff -----	buffer	grp -----	group	rect -----	rectangle
byts -----	bytes	hdl -----	handle	ref -----	reference
cal -----	calculator	hdr -----	header	reg -----	register
ch -----	change	help -----	helper	rem -----	remote
char -----	character	hex -----	hexadecimal	rem -----	remove
chg -----	change	hist -----	history	req -----	request
chr -----	character	id -----	identifier	req -----	requirement
clck -----	clock	idx -----	index	rng -----	range
clr -----	clear	imag -----	image	sched -----	scheduler
cmd -----	command	img -----	image	sect -----	sector
cmp -----	compare	indx -----	index	seg -----	segment
cntr -----	counter	info -----	information	sel -----	select
col -----	column	inform -----	information	sem -----	semaphored
comm -----	communication	init -----	initialization	sem -----	semaphore
comp -----	compare	int -----	integer	seq -----	sequence
config -----	configuration	interf -----	interface	sfx -----	suffix
conn -----	connection	intern -----	internal	siz -----	size
connect -----	connection	iter -----	iterator	spec -----	specification
const -----	constant	kernl -----	kernel	src -----	source
constr -----	constraint	lang -----	lang	stat -----	statistic
cont -----	continuation	lbl -----	label	std -----	standard
cpy -----	copy	len -----	length	str -----	string
ctls -----	controls	lib -----	library	struct -----	structure
ctr -----	control	lim -----	limit	succ -----	success
ctrl -----	control	ln -----	link	svc -----	service
cur -----	current	lng -----	long	syml -----	symbol
curr -----	current	lnk -----	link	sys -----	system
cust -----	custom	lst -----	list	sz -----	size
db -----	database	mapp -----	mappings	tab -----	table
dec -----	declaration	matr -----	matrix	targ -----	target
del -----	delete	max -----	maximum	tbl -----	table
dest -----	destination	mem -----	memory	term -----	terminal
dev -----	deviation	met -----	method	tmp -----	temporary
dev -----	device	meth -----	method	trans -----	transparency
dict -----	dictionary	mltpy -----	multiply	trc -----	trace
diff -----	difference	mnt -----	mount	unlck -----	unlock
dir -----	directory	mod -----	modify	updt -----	update
dpcy -----	dependency	mod -----	module	usr -----	user
dyn -----	dynamic	nbr -----	number	val -----	value
edt -----	editing	net -----	network	var -----	variable
elem -----	element	nod -----	node	vec -----	vector
elt -----	element	notif -----	notification	vect -----	vector
enum -----	enumeration	num -----	number	win -----	window
env -----	environment	obj -----	object	wind -----	windows
envp -----	envelope	ok -----	ok	wrd -----	word
ers -----	erase	oper -----	operator		

Jumlah: 191 kata

### C. Kamus Singkatan (*Programming*)

abs -----absolute	exec ----- execute	params --parameters
addr ----- address	fd----- fd	pdf -----pdf
ansi ----- ansi	fifo ----- fifo	ptr-----pointer
arg----- argument	fig ----- figure	rect -----rectangle
args ----- arguments	figs----- figures	regexp---regexp
ascii----- ascii	flg ----- flag	sem -----semantic
avg ----- average	flgs----- flags	snd -----sound
boolean -boolean	func----- function	src -----source
buf----- buffer	funcs --- functions	str-----string
char ----- character	gid ----- gid	strs -----strings
charset -- charset	gnu ----- gnu	struct----structure
chr----- character	hdr ----- header	structs ---structures
chr----- character	hdr ----- header	sym-----symbol
chrs ----- characters	html----- html	syms -----symbols
chrs ----- characters	html----- html	tbl -----table
clr----- clear	ibm----- ibm	tex -----tex
cmd ----- command	init----- initialize	tmp -----temporary/tmp
cmp ----- compare	inode ---- inode	todo -----todo
cobol----- cobol	inode ---- inode	tooltip ----tooltip
const----- constant	int ----- integer	tty -----tty
cpy ----- copy	iso----- iso	tuple-----tuple
cwd----- cwd	iso----- iso	txt -----text
del----- delete	len ----- length	uid -----uid
descr----- description	leng----- length	unix-----unix
dev ----- device	lib ----- library	url-----url
dir----- directory	linux ---- linux	utf8 -----utf8
dirs ----- directories	login ---- login	val -----value
doc ----- document	malloc--- malloc	vals-----values
elem ----- element	mem----- memory	vec -----vector
eof----- eof	misc ---- miscellaneous	ver -----version
eol----- eol	msg----- message	vers -----versions
eor----- eor	msgs --- messages	vga -----vga
eps ----- eps	multi ---- multi	vga -----vga
eps ----- eps	num----- number	vol -----volume
err ----- error	obj ----- object	win-----window
errs----- errors	param --- parameter	

Jumlah: 129 kata

*Halaman ini sengaja dikosongkan*

## Lampiran2

### Stop List

a	available	contains	five	himself
a's	away	corresponding	followed	his
able	awfully	could	following	hither
about	b	couldn't	follows	hopefully
above	be	course	for	how
according	became	currently	former	howbeit
accordingly	because	d	formerly	however
across	become	definitely	forth	i
actually	becomes	described	four	i'd
after	becoming	despite	from	i'll
afterwards	been	did	further	i'm
again	before	didn't	furthermore	i've
against	beforehand	different	g	ie
ain't	behind	do	get	if
all	being	does	gets	ignored
allow	believe	doesn't	getting	immediate
allows	below	doing	given	in
almost	beside	don't	gives	inasmuch
alone	besides	done	go	inc
along	best	down	goes	indeed
already	better	downwards	going	indicate
also	between	during	gone	indicated
although	beyond	e	got	indicates
always	both	each	gotten	inner
am	brief	edu	greetings	insofar
among	but	eg	h	instead
amongst	by	eight	had	into
an	c	either	hadn't	inward
and	c'mon	else	happens	is
another	c's	elsewhere	hardly	isn't
any	came	enough	has	it
anybody	can	entirely	hasn't	it'd
anyhow	can't	especially	have	it'll
anyone	cannot	et	haven't	it's
anything	cant	etc	having	its
anyway	cause	even	he	itself
anyways	causes	ever	he's	j
anywhere	certain	every	hello	just
apart	certainly	everybody	help	k
appear	changes	everyone	hence	keep
appreciate	clearly	everything	her	keeps
appropriate	co	everywhere	here	kept
are	com	ex	here's	know
aren't	come	exactly	hereafter	knows
around	comes	example	hereby	known
as	concerning	except	herein	l
aside	consequently	f	hereupon	last
ask	consider	far	hers	lately
asking	considering	few	herself	later
associated	contain	fifth	hi	latter
at	containing	first	him	latterly

least	nowhere	s	than	unless
less	o	said	thank	unlikely
lest	obviously	same	thanks	until
let	of	saw	thanx	unto
let's	off	say	that	up
like	often	saying	that's	upon
liked	oh	says	thats	us
likely	ok	second	the	use
little	okay	secondly	their	used
look	old	see	theirs	useful
looking	on	seeing	them	uses
looks	once	seem	themselves	using
ltd	one	seemed	then	usually
m	ones	seeming	thence	uucp
mainly	only	seems	there	v
many	onto	seen	there's	value
may	or	self	thereafter	various
maybe	other	selves	thereby	very
me	others	sensible	therefore	via
mean	otherwise	sent	therein	viz
meanwhile	ought	serious	theres	vs
merely	our	seriously	thereupon	w
might	ours	seven	these	want
more	ourselves	several	they	wants
moreover	out	shall	they'd	was
most	outside	she	they'll	wasn't
mostly	over	should	they're	way
much	overall	shouldn't	they've	we
must	own	since	think	we'd
my	p	six	third	we'll
myself	particular	so	this	we're
n	particularly	some	thorough	we've
name	per	somebody	thoroughly	welcome
namely	perhaps	somehow	those	well
nd	placed	someone	though	went
near	please	something	three	were
nearly	plus	sometime	through	weren't
necessary	possible	sometimes	throughout	what
need	presumably	somewhat	thru	what's
needs	probably	somewhere	thus	whatever
neither	provides	soon	to	when
never	q	sorry	together	whence
nevertheless	que	specified	too	whenever
new	quite	specify	took	where
next	qv	specifying	toward	where's
nine	r	still	towards	whereafter
no	rather	sub	tried	whereas
nobody	rd	such	tries	whereby
non	re	sup	truly	wherein
none	really	sure	try	whereupon
noone	reasonably	t	trying	wherever
nor	regarding	t's	twice	whether
normally	regardless	take	two	which
not	regards	taken	u	while
nothing	relatively	tell	un	whither
novel	respectively	tends	under	who
now	right	th	unfortunately	who's



whoever  
whole  
whom  
whose  
why  
will  
willing

wish  
with  
within  
without  
won't  
wonder

would  
would  
wouldn't  
x  
y  
yes

yet  
you  
you'd  
you'll  
you're  
you've

your  
yours  
yourself  
yourselves  
z  
zero

*Halaman ini sengaja dikosongkan*

### Lampiran 3

#### Detail Data Uji

Topik	No.	Nama File	URL Sumber
Thread synchronization	1	com.in28minutes.java.threads.ThreadDeadlock.java	<a href="https://goo.gl/kUv3PT">https://goo.gl/kUv3PT</a>
	2	PingPong.java	<a href="https://goo.gl/pzPrhJ">https://goo.gl/pzPrhJ</a>
	3	IncorrectSynchronizedIncrement.java	<a href="https://goo.gl/jL4tct">https://goo.gl/jL4tct</a>
	4	sleep.practice.com.sleepimpl.Threa.↵ sleep%20Implementation.java	<a href="https://goo.gl/kpBjcc">https://goo.gl/kpBjcc</a>
	5	classes.test.Synchronized.java	<a href="https://goo.gl/IVYHxx">https://goo.gl/IVYHxx</a>
	6	BlockingQueue.java	<a href="https://goo.gl/7b3UYj">https://goo.gl/7b3UYj</a>
	7	java.util.concurrent.Semaphore.TwistedPairVer2.java	<a href="https://goo.gl/s8RjQY">https://goo.gl/s8RjQY</a>
	8	Collision.java	<a href="https://goo.gl/IQ2nvM">https://goo.gl/IQ2nvM</a>
Accessing MySQL	9	ucf.mysql.jdbc.DBEngine.java	<a href="https://goo.gl/Bs2ZD5">https://goo.gl/Bs2ZD5</a>
	10	io.github.web.data.mysql.MySqlJdbcInsert.java	<a href="https://goo.gl/LO1dUX">https://goo.gl/LO1dUX</a>
	11	org.embulk.output.MySQLOutputPlugin.java	<a href="https://goo.gl/ei0lul">https://goo.gl/ei0lul</a>
	12	JdbcMysql.java	<a href="https://goo.gl/iQKBEj">https://goo.gl/iQKBEj</a>
	13	com.github.kristofa.brave.mysql.↵ MySQLStatementInterceptor.java	<a href="https://goo.gl/c2UkH3">https://goo.gl/c2UkH3</a>
	14	org.ops4j.pax.jdbc.mysql.impl.↵ MysqlDataSourceFactory.java	<a href="https://goo.gl/gVcejs">https://goo.gl/gVcejs</a>
	15	com.gmail.johanringmann.blockcounter.MySQL.java	<a href="https://goo.gl/8VPqdV">https://goo.gl/8VPqdV</a>
	16	com.mysql.jdbc.integration.c3p0.↵ MysqlConnectionTester.java	<a href="https://goo.gl/9h5LIX">https://goo.gl/9h5LIX</a>
Bubble sort	17	com.github.pedrovgz.problem74.BubbleSort.java	<a href="https://goo.gl/oJYuE4">https://goo.gl/oJYuE4</a>
	18	LinkedList.java	<a href="https://goo.gl/5azHkx">https://goo.gl/5azHkx</a>
	19	bubbleSort.java	<a href="https://goo.gl/WoZHv">https://goo.gl/WoZHv</a>
	20	spagy.BubbleSort.java	<a href="https://goo.gl/zzbJQ2">https://goo.gl/zzbJQ2</a>
	21	sorting.algorithms.BubbleSort.java	<a href="https://goo.gl/W682RB">https://goo.gl/W682RB</a>
	22	quincy.BubbleSort.java	<a href="https://goo.gl/tYVM7d">https://goo.gl/tYVM7d</a>
	23	jeffheaton.BubbleSort.java	<a href="https://goo.gl/D7Euc4">https://goo.gl/D7Euc4</a>
	24	JakenHerman.Bubblesort.java	<a href="https://goo.gl/j5i1ow">https://goo.gl/j5i1ow</a>
Draw rectangle	25	com.googlecode.lanterna.examples.DrawRectangle.java	<a href="https://goo.gl/C473VJ">https://goo.gl/C473VJ</a>
	26	org.peimari.gleaflet.client.draw.DrawRectangle.java	<a href="https://goo.gl/s9IUce">https://goo.gl/s9IUce</a>
	27	q2.DrawRectangle.java	<a href="https://goo.gl/whiqC4">https://goo.gl/whiqC4</a>
	28	com.mlakhia.draw.shapes.Rectangle.java	<a href="https://goo.gl/7uYWv6">https://goo.gl/7uYWv6</a>
	29	RCoon.DrawRectangle.java	<a href="https://goo.gl/cnkj5f">https://goo.gl/cnkj5f</a>
	30	shapes.rectangle.java	<a href="https://goo.gl/M3B6wa">https://goo.gl/M3B6wa</a>
	31	mrquatsch.Rectangle.java	<a href="https://goo.gl/xUqdJ6">https://goo.gl/xUqdJ6</a>
	32	dcstimm.Rectangle.java	<a href="https://goo.gl/Q6dj69">https://goo.gl/Q6dj69</a>

*Halaman ini sengaja dikosongkan*

## Lampiran 4

### Hasil Uji Coba Menggunakan Normalisasi Lingua::IdSplitter

<b>Topik</b>	<b>#Query</b>	<b>#Relevan</b>	<b>Precision</b>	<b>Recall</b>
<i>Thread Synchronization</i>	1	11	0.7334	0.8462
	2	11	0.7334	0.8462
	3	11	0.7334	0.8462
	4	6	0.4	0.4616
	5	11	0.7334	0.8462
	6	6	0.4	0.4616
	7	3	0.2	0.2308
	8	10	0.6667	0.7693
<i>Accessing MySQL</i>	9	5	0.3334	0.8334
	10	5	0.3334	0.8334
	11	3	0.2	0.5
	12	4	0.2667	0.6667
	13	4	0.2667	0.6667
	14	4	0.2667	0.6667
	15	6	0.4	1
	16	4	0.2667	0.6667
<i>Bubble Sort</i>	17	2	0.1334	0.4
	18	2	0.1334	0.4
	19	0	0	0
	20	4	0.2667	0.8
	21	3	0.2	0.6
	22	3	0.2	0.6
	23	4	0.2667	0.8
	24	2	0.1334	0.4
<i>Draw Rectangle</i>	25	3	0.2	0.375
	26	7	0.4667	0.875
	27	4	0.2667	0.5
	28	5	0.3334	0.625
	29	5	0.3334	0.625
	30	6	0.4	0.75
	31	7	0.4667	0.875
	32	7	0.4667	0.875

Keterangan: Nomor query dapat dilihat pada Lampiran 3.

*Halaman ini sengaja dikosongkan*

## Lampiran 5

### Hasil Uji Coba Tanpa Normalisasi Lingua::IdSplitter

<b>Topik</b>	<b>#Query</b>	<b>#Relevan</b>	<b>Precision</b>	<b>Recall</b>
<i>Thread Synchronization</i>	1	11	0.7334	0.8462
	2	11	0.7334	0.8462
	3	11	0.7334	0.8462
	4	6	0.4	0.4616
	5	11	0.7334	0.8462
	6	6	0.4	0.4616
	7	3	0.2	0.2308
	8	10	0.6667	0.7693
<i>Accessing MySQL</i>	9	5	0.3334	0.8334
	10	5	0.3334	0.8334
	11	3	0.2	0.5
	12	4	0.2667	0.6667
	13	4	0.2667	0.6667
	14	4	0.2667	0.6667
	15	6	0.4	1
	16	4	0.2667	0.6667
<i>Bubble Sort</i>	17	2	0.1334	0.4
	18	2	0.1334	0.4
	19	0	0	0
	20	4	0.2667	0.8
	21	3	0.2	0.6
	22	3	0.2	0.6
	23	4	0.2667	0.8
	24	2	0.1334	0.4
<i>Draw Rectangle</i>	25	3	0.2	0.375
	26	7	0.4667	0.875
	27	4	0.2667	0.5
	28	5	0.3334	0.625
	29	5	0.3334	0.625
	30	6	0.4	0.75
	31	7	0.4667	0.875
	32	7	0.4667	0.875

Keterangan: Nomor query dapat dilihat pada Lampiran 3.

*Halaman ini sengaja dikosongkan*



## **BAB 5**

### **KESIMPULAN**

#### **5.1 Kesimpulan**

Berdasarkan hasil uji coba dan penjelasan yang telah dipaparkan, dapat disimpulkan bahwa:

1. Proses normalisasi *identifier* menggunakan `Lingua::IdSplitter` hanya mampu meningkatkan performasistem rekomendasi manakaladata-data diskusi yang ada mengandung *identifier* dengankomposisi term singkatan.
2. Ketidaktepatan pemilihan *identifier* yang layak untuk dinormalisasi menjadi penyebab utama turunnya performa sistem rekomendasi.
3. Pemanfaatan model TF-IDF masih layak digunakan dalam membangun sistem rekomendasi berdasarkan nilai median *recall* umum yang dihasilkan, yakni sebesar 0,67.

#### **5.2 Saran**

Guna perbaikanproses normalisasi *identifier* pada sistem rekomendasi, perlu dipertimbangkan adanya proses pemilihan *identifier* yang layak untuk dinormalisasi sebelum dijadikan masukan pada proses normalisasi *identifier*.

*Halaman ini sengaja dikosongkan*

## DAFTAR PUSTAKA

- Adomavicius, G. dan Tuzhilin, A., (2005). Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions. *Knowledge and Data Engineering, IEEE Transaction*, 17(6), pp.734-749.
- Anderson, A., Huttenlocher, D., Kleinberg, J., dan Leskovec, J., (2012). Discovering Value from Community Activity on Focused Question Answering Sites: A Case Study of Stack Overflow. *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 850-858). ACM.
- Arwan, A., Rochimah, S., dan Akbar, R.J., (2015). Source Code Retrieval on StackOverflow Using LDA. *Information and Communication Technology (ICoICT), 2015 3rd International Conference* (pp. 295-299). IEEE.
- Carvalho, N.R., Almeida, J.J., Henriques, P.R., dan Varanda, M.J., (2015). From Source Code Identifiers to Natural Language Terms. *Journal of Systems and Software*, 100, pp.117-128.
- Cegielski, C.G. dan Hall, D.J., (2006). What Makes A Good Programmer?. *Communications of the ACM*, 49(10), pp.73-75.
- Cordeiro, J., Antunes, B., dan Gomes, P., (2012). Context-Based Recommendation to Support Problem Solving in Software Development. *Recommendation Systems for Software Engineering (RSSE), 2012 Third International Workshop* (pp. 85-89). IEEE.
- Gosling, J., Joy, B., Steele, G., Bracha, G., dan Buckley, A. (2015). Chapter 3. Lexical Structure, (online), (<https://docs.oracle.com/javase/specs/jls/se8/html/jls-3.html>, diakses 20 Juli 2016).
- Hu, D.J., (2009). Latent Dirichlet Allocation for Text, Images, and Music. *University of California, San Diego*.
- Jivani, A.G., (2011). A Comparative Study of Stemming Algorithms. *Int. J. Comp. Tech. Appl*, 2(6), pp.1930-1938.

- Krauthammer, M. dan Nenadic, G., (2004). Term identification in the biomedical literature. *Journal of biomedical informatics*, 37(6), pp.512-526.
- LaToza, T.D., Venolia, G., dan DeLine, R., (2006). Maintaining Mental Models: A Study of Developer Work Habits. *Proceedings of the 28th international conference on Software engineering* (pp. 492-501). ACM.
- Lawrie, D., Morrell, C., Feild, H., dan Binkley, D., (2007). Effective identifier names for comprehension and memory. *Innovations in Systems and Software Engineering*, 3(4), pp.303-318.
- Mahmoud, A. dan Niu, N., (2011). Source Code Indexing for Automated Tracing. *Proceedings of the 6th International Workshop on Traceability in Emerging Forms of Software Engineering* (pp. 3-9). ACM.
- Manning, C.D., Raghavan, P., dan Schütze, H., (2008). *Introduction to Information Retrieval*. Cambridge: Cambridge University Press.
- MySQL. (2006). *MySQL stopword list*, (online), (<http://dev.mysql.com/doc/refman/5.0/en/fulltext-stopwords.html>, diakses 2 Agustus 2016).
- Ponzanelli, L., Bacchelli, A., dan Lanza, M., (2013). Leveraging Crowd Knowledge for Software Comprehension and Development. *Software Maintenance and Reengineering (CSMR), 2013 17th European Conference* (pp. 57-66). IEEE.
- Robillard, M.P., Walker, R.J., dan Zimmermann, T., (2010). Recommendation Systems for Software Engineering. *Software, IEEE*, 27(4), pp.80-86.
- Sun Microsystems, (1997). *Java Code Conventions*. Sun Microsystems, Inc.
- Vohra, D., Baesens, B., Backiel, A. dan Vanden Broucke, S., (2015). *Beginning Java Programming: The Object-oriented Approach*. John Wiley & Sons.
- Zorychta, D. (18 Desember 2011). *How do I turn off text antialiasing in this Java function?* Stack Overflow, (online), (<http://stackoverflow.com/questions/8549891/how-do-i-turn-off-text-antialiasing-in-this-java-function>, diakses 11 Mei 2016).

## BIODATA PENULIS



Nanang Fakhur Rozi, lahir di Gresik, 9 April 1988. Merupakan anak ke-3 dari 5 bersaudara. Berasal dari desa Kalirejo kecamatan Dukun kabupaten Gresik. Pendidikan yang telah ditempuh antara lain Sekolah Dasar Negeri Kalirejo Gresik, Madrasah Tsanawiyah YKUI Maskumambang Gresik, dan lulus Madrasah Aliyah YKUI Maskumambang Gresik pada tahun 2005, kemudian melanjutkan pendidikan tinggi pada program Diploma III Jurusan Teknologi Informasi di Politeknik Elektronika Negeri Surabaya (PENS) hingga lulus pada tahun 2008. Tahun 2009 kembali melanjutkan studi pada program Diploma IV (Lanjut Jenjang) Jurusan Teknik Informatika di Politeknik Elektronika Negeri Surabaya dan lulus pada tahun 2011. Ketika laporan penelitian ini disusun, penulis sedang menempuh studi akhirnya pada Program Magister Jurusan Teknik Informatika Institut Teknologi Sepuluh Nopember (ITS) Surabaya. Penulis memiliki minat pada bidang rekayasa perangkat lunak serta keterkaitannya dengan bidang kecerdasan buatan.

### Informasi Kontak:

E-mail : [nfrozy@gmail.com](mailto:nfrozy@gmail.com)

Situs Web : <http://nfrozi.web.id>

LinkedIn : <http://linkedin.com/in/nfrozi>