

INTEGRASI WORD EMBEDDINGS DAN INVERSE BOOK FREQUENCY DALAM PEMBOBOTAN TERM UNTUK PENINGKATAN Pencarian DOKUMEN

Dwi Ari Suryaningrum*¹⁾, Rahmad Syaifudin²⁾, Haniel Ranga Pramuditya Putra³⁾

1. Teknik Elektro, Fakultas Teknik, Universitas Tulungagung, Indonesia
2. Teknik Elektro, Fakultas Teknik, Universitas Tulungagung, Indonesia
3. Teknik Elektro, Fakultas Teknik, Universitas Tulungagung, Indonesia

Article Info

Kata Kunci: IBF-IDF; *Term Weighting*; TF; Word Embeddings; *Query Expansion*

Keywords: IBF-IDF; *Term Weighting*; TF; Word Embeddings; *Query Expansion*

Article history:

Received 20 September 2024

Revised 19 Oktober 2024

Accepted 24 Oktober 2024

Available online 1 December 2024

DOI :

<https://doi.org/10.29100/jupi.v9i4.7557>

* Corresponding author.

Corresponding Author

E-mail address:

dwiarisuryaningrum@unita.ac.id

ABSTRAK

Pencarian dokumen yang relevan dapat ditingkatkan dengan metode ekspansi kueri berbasis word embeddings. Studi ini mengusulkan pendekatan pembobotan ekspansi kueri dengan mempertimbangkan korelasi *term* terhadap kueri serta frekuensinya dalam dokumen menggunakan metode *Word Embeddings* (WE) dan *Inverse Book Frequency* (IBF). Pembobotan dilakukan dengan mengalikan nilai similaritas dari WE dengan bobot TF-IDF-IBF untuk meningkatkan relevansi pencarian dokumen secara lebih akurat. Hasil eksperimen menunjukkan bahwa metode ini menghasilkan *f-score* sebesar 0,743, dengan performa optimal ketika jumlah *term* ekspansi yang dipilih lebih sedikit. Selain itu, metode ini lebih unggul dibandingkan pendekatan tradisional seperti TF-IDF atau BM25 dalam mengurangi *term* yang tidak relevan, sehingga meningkatkan efektivitas pencarian informasi dalam dataset yang lebih luas. Namun, pendekatan ini masih memiliki keterbatasan dalam kompleksitas komputasi serta ketergantungan pada kualitas dataset pelatihan yang digunakan. Studi ini menyarankan eksplorasi lebih lanjut dengan model berbasis transformer seperti BERT atau RoBERTa untuk meningkatkan efektivitas pencarian dokumen. Dengan mengintegrasikan metode ini ke dalam sistem pencarian informasi, diharapkan pencarian dokumen menjadi lebih akurat, efisien, dan relevan dengan kebutuhan pengguna di berbagai domain aplikasi.

ABSTRACT

Relevant document retrieval can be improved using query expansion methods based on word embeddings. This study proposes a query expansion weighting approach that considers *term* correlation with the query and its frequency in documents using the Word Embeddings (WE) and Inverse Book Frequency (IBF) methods. The weighting process multiplies similarity values from WE with TF-IDF-IBF scores to enhance document retrieval relevance more accurately. Experimental results indicate that this method achieves an *f-score* of 0.743, with optimal performance when fewer expansion *terms* are selected. Furthermore, this method outperforms traditional approaches such as TF-IDF and BM25 in reducing irrelevant *terms*, thus improving information retrieval effectiveness in larger datasets. However, it still has limitations in *terms* of computational complexity and dependence on the quality of the training dataset used. This study suggests further exploration with transformer-based models such as BERT or RoBERTa to enhance document retrieval effectiveness. By integrating this method into information retrieval systems, document searches are expected to become more accurate, efficient, and relevant to user needs across various application domains.

I. PENDAHULUAN

SISTEM temu kembali informasi merupakan proses yang bertujuan untuk menemukan informasi yang paling relevan sesuai dengan kebutuhan pengguna. Metode ekspansi kueri (*Query Expansion/QE*) telah menjadi pendekatan yang banyak digunakan untuk meningkatkan relevansi pencarian dokumen dengan menambahkan istilah-istilah baru ke dalam kueri awal [1]. Namun, metode QE tradisional sering kali mengalami

masalah dalam memilih istilah yang benar-benar relevan, yang dapat menyebabkan penurunan akurasi dalam pencarian [2].

Penelitian terbaru menunjukkan bahwa embedding kontekstual, seperti ELMo dan BERT, dapat meningkatkan efektivitas ekspansi kueri dengan memanfaatkan pemahaman semantik yang lebih dalam [3]. Sebagai contoh, metode Contextualized Embeddings for Query Expansion (CEQE) telah digunakan untuk meningkatkan pengambilan dokumen secara signifikan dibandingkan metode tradisional [4]. Selain itu, pendekatan kombinasi embedding kata dengan strategi data-driven menggunakan *Inverse Document Frequency* (IDF) juga menunjukkan hasil yang lebih baik dalam temu kembali informasi [5]. Meskipun pendekatan ini menjanjikan, mereka masih memiliki keterbatasan dalam memilih istilah yang memiliki relevansi tinggi terhadap dokumen target [6].

Dalam penelitian ini, kami mengusulkan metode pembobotan istilah ekspansi kueri yang mempertimbangkan korelasi istilah terhadap kueri serta frekuensi kemunculan dalam dokumen menggunakan pendekatan Word Embeddings (WE) dan Inverse Book Frequency (IBF) [7]. Pembobotan dilakukan dengan mengalikan nilai similaritas dari WE dengan bobot TF-IDF-IBF untuk menyaring istilah yang paling relevan [8]. Metode ini diharapkan dapat mengatasi kekurangan pada pendekatan sebelumnya, di mana beberapa istilah hasil ekspansi masih kurang sesuai dengan konteks pencarian dokumen [9].

Metode yang diusulkan memiliki keunggulan dibandingkan pendekatan lain, seperti TF-IDF dan BM25, karena mempertimbangkan distribusi istilah dalam berbagai topik dokumen, bukan hanya dalam koleksi dokumen secara keseluruhan. Dengan demikian, IBF mampu memberikan bobot lebih tinggi kepada istilah yang lebih spesifik dalam suatu domain tertentu, sehingga meningkatkan akurasi pencarian dokumen [10]. Selain itu, metode ini lebih cocok untuk dokumen berbahasa Indonesia dibandingkan dengan pendekatan berbasis bahasa Inggris yang mungkin tidak memperhitungkan keunikan morfologi dan struktur bahasa Indonesia [11].

Namun, terdapat beberapa batasan dalam penelitian ini. Metode yang diusulkan memiliki keterbatasan dalam hal skalabilitas karena perhitungan IBF memerlukan proses pemetaan topik dokumen yang dilakukan secara manual, yang dapat menjadi tidak efisien untuk kumpulan data yang sangat besar [12]. Selain itu, metode ini bergantung pada kualitas model embedding kata yang digunakan; embedding yang kurang representatif dapat menyebabkan pemilihan istilah ekspansi yang kurang optimal [13]. Dari sisi kompleksitas komputasi, pendekatan ini juga memerlukan sumber daya komputasi lebih tinggi dibandingkan metode tradisional seperti TF-IDF [14].

Penelitian ini menggunakan kumpulan dokumen berita online berbahasa Indonesia yang telah dikelompokkan sesuai dengan topik masing-masing. Dalam evaluasi eksperimental, metode yang diusulkan dibandingkan dengan beberapa teknik pembobotan *term* lainnya, seperti TF-IDF, BM25, dan metode berbasis word embeddings seperti Word2Vec dan GloVe. Eksperimen ini bertujuan untuk mengukur efektivitas metode dalam meningkatkan presisi, recall, dan *F-score* dalam pencarian dokumen [15].

Dengan adanya metode TF-IDF-IBF yang diusulkan, diharapkan dapat memberikan kontribusi dalam meningkatkan efektivitas sistem temu kembali informasi, khususnya dalam konteks pencarian dokumen berbahasa Indonesia. Ke depannya, penelitian ini dapat dikembangkan lebih lanjut dengan mengimplementasikan pendekatan otomatis dalam pemetaan topik dokumen untuk mengurangi beban manual dalam perhitungan IBF [16]. Selain itu, penelitian selanjutnya dapat mengeksplorasi integrasi dengan model pembelajaran mesin untuk mengoptimalkan seleksi istilah ekspansi kueri [17].

II. METODE PENELITIAN

Penelitian ini terdiri dari beberapa tahapan utama, yaitu preprocessing dokumen, ekspansi *query* (*Query Expansion/QE*) menggunakan beberapa model word embedding, pembobotan *term*, serta pencocokan kemiripan antara *query* dan dokumen untuk memperoleh dokumen yang paling relevan.

A. Preprocessing

Dokumen yang digunakan dalam penelitian ini terlebih dahulu melalui tahap preprocessing guna memperoleh *term-term* yang digunakan dalam pembobotan. Preprocessing merupakan langkah penting dalam sistem information retrieval. Proses ini mencakup beberapa tahap utama seperti tokenisasi, normalisasi huruf (case folding), penyaringan kata (filtering), penghapusan kata umum (stopword removal), dan stemming [18]. Tahapan ini bertujuan untuk meningkatkan efisiensi dalam pencarian informasi dengan hanya mempertahankan kata-kata yang memiliki nilai informasi tinggi.

B. Pembobotan Term

B.1. Term Frequency (TF)

Term Frequency (TF) adalah metode pembobotan sederhana yang mengasumsikan bahwa kepentingan suatu

kata dalam dokumen berbanding lurus dengan frekuensi kemunculannya dalam dokumen tersebut. Nilai bobot TF dari *term t* dalam dokumen *d* dihitung menggunakan (1):

$$TF(d, t) = f(d, t) \quad (1)$$

dimana $f(d, t)$ merupakan frekuensi kemunculan *term t* pada dokumen *d*.

B.2. Inverse Document Frequency (IDF)

Berbeda dengan TF, metode *Inverse Document Frequency* (IDF) mempertimbangkan seberapa jarang suatu *term* muncul dalam kumpulan dokumen. *Term* yang lebih jarang muncul di seluruh koleksi dokumen dianggap lebih penting. IDF dihitung menggunakan (2):

$$IDF(t) = 1 + \log(N_d / df(t)) \quad (2)$$

dimana N_d merupakan jumlah seluruh dokumen dan $df(t)$ jumlah dokumen yang mengandung *term t*.

B.3. Inverse Book Frequency (IBF)

Berbeda dengan IDF yang memperhatikan kemunculan *term* pada kumpulan dokumen, IBF memperhatikan kemunculan *term* pada kumpulan dokumen yang memiliki beberapa topik. *Term* yang bernilai untuk klasifikasi adalah *term* yang jarang muncul pada banyak dokumen dengan beberapa topik. IBF dihitung dengan turunan langsung dari persamaan IDF, seperti pada (3):

$$IBF(t) = 1 + \log(N_b / bf(t)) \quad (3)$$

dimana N_b adalah jumlah seluruh topik book dan $bf(t)$ jumlah topik yang mengandung *term t*.

B.4. TF-IDF-IBF

Untuk meningkatkan efektivitas pembobotan, penelitian ini mengombinasikan ketiga metode di atas dalam satu pendekatan bernama TF-IDF-IBF. Perhitungannya menggunakan (4):

$$W(d, t) = TF(d, t) \times IDF(t) \times IBF(t) \quad (4)$$

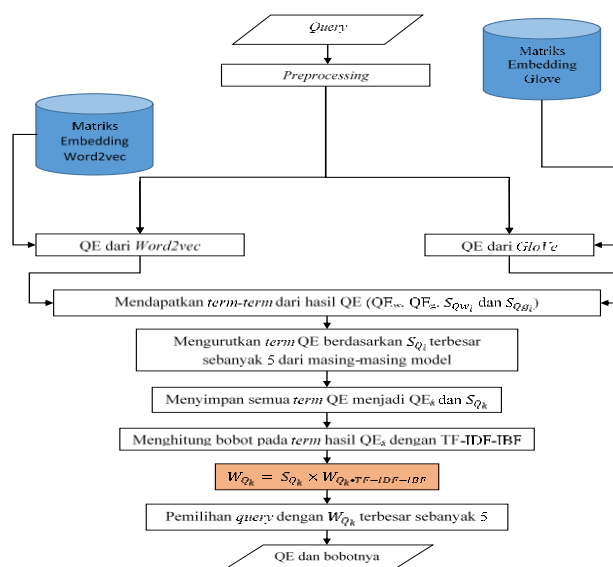
Setelah mendapatkan bobot untuk setiap *term* dalam dokumen pada berbagai topik, dilakukan perhitungan rata-rata menggunakan metode mean TF-IDF-IBF sebagaimana dirumuskan dalam (5):

$$Mean(TF - IDF - IBF)(t) = \frac{\sum_{i=1}^n W(d_i, t)}{n} \quad (5)$$

dimana d_i adalah dokumen ke-*i* dan n adalah jumlah dokumen dalam koleksi.

C. Ekspansi Query dengan Word Embeddings

Tahap pemilihan *query expansion* dengan word embeddings dapat dilihat pada Gambar 1. Pada tahap ini, hasil *query expansion* dari word embeddings tidak dipilih semua sebagai *query expansion* akhir untuk proses pencarian dokumen. Ekspansi *query* dilakukan dengan menggunakan model word embeddings seperti Word2Vec dan GloVe. Word embeddings merepresentasikan kata dalam bentuk vektor numerik sehingga kata-kata dengan makna serupa memiliki representasi vektor yang berdekatan dalam ruang vektor [19].



Gambar 1. Tahapan pemilihan *query expansion* dengan word embeddings dan IBF

Query yang dimasukkan oleh pengguna akan melalui tahap preprocessing untuk menghasilkan *term-term* dasar. Selanjutnya, *term-term* ini diperluas dengan menambahkan kata-kata yang memiliki kedekatan makna berdasarkan

model Word2Vec dan GloVe. Kemiripan kata dihitung menggunakan persamaan similaritas kosinus, yang menghasilkan nilai dalam rentang -1 hingga 1, di mana 1 menunjukkan similaritas tertinggi [20].

Ekspansi *query* dilakukan dengan memilih lima *term* dengan nilai similaritas tertinggi dari masing-masing model word embeddings. *Term* yang diperoleh dari ekspansi *query* selanjutnya dihitung bobotnya menggunakan metode TF-IDF-IBF. Bobot akhir diperoleh dengan mengalikan nilai similaritas dengan bobot TF-IDF-IBF dari setiap *term* sebagaimana dirumuskan dalam (6):

$$W_{Q_i} = S_{Q_i} \times W_{Q_i,TF-IDF-IBF} \quad (6)$$

dimana W_{Q_i} adalah pembobotan *term* hasil ekspansi *query* S_{Q_i} adalah similaritas *query*, $W_{Q_i,TF-IDF-IBF}$ adalah nilai bobot TF-IDF-IBF pada *term* hasil QE. *Term* hasil ekspansi *query* yang dipilih adalah *term-term* yang memiliki similaritas *query* tinggi dan nilai bobot TF-IDF-IBF juga tinggi.

D. Pembobotan Term Hasil Ekspansi Query

Ekspansi *query* dari hasil word embeddings dihitung pembobotannya menggunakan TF-IDF-IBF. Pembobotan IBF digunakan untuk melihat frekuensi kemunculan *term-term* hasil ekspansi *query* pada beberapa dokumen yang memiliki topik/kategori beragam. Penelitian ini sudah dilakukan pengelompokan kategori secara manual, sehingga proses hanya perlu melihat *term* berada pada dokumen dan kategori mana saja, kemudian dihitung frekuensi kemunculannya. Pada penelitian ini, *term* dari *query* asli akan diberikan nilai faktor berbeda dari *term* hasil ekspansi *query*. Pemberian nilai faktor ini bertujuan agar *term* dari *query* asli tetap memiliki bobot yang lebih tinggi dibandingkan dengan *term* hasil ekspansi *query*. Bobot akhir *term* dari *query* didapatkan dengan mengalikan bobot TF-IDF-IBF dengan nilai faktor yang dimiliki, seperti pada (7).

$$W_{Q_{A_i}} = W_{Q_i,TF-IDF-IBF} \times \alpha \quad (7)$$

dimana $W_{Q_{A_i}}$ adalah nilai akhir bobot TF-IDF-IBF dan α adalah nilai faktor pada *term* dari *query* dan ekspansi *query*. α bernilai 2 jika *term* berasal dari *query* asli dan bernilai 0,5 jika *term* berasal dari hasil ekspansi *query*.

E. Pengukuran Kemiripan Query dengan Dokumen

Untuk mengukur tingkat kemiripan antara *query* dan dokumen dalam koleksi, digunakan metode cosine similarity. Cosine similarity mengukur besar sudut antara vektor *query* dan vektor dokumen dalam ruang vektor multidimensi [21]. Persamaannya diberikan dalam (8):

$$\cos(q, d_j) = \frac{\sum t_k [TF-IDF-IBF(t_{q,q})] [TF-IDF-IBF(t_k, d_j)]}{\sqrt{\sum [TF-IDF-IBF_q]^2} \cdot \sqrt{\sum [TF-IDF-IBF_{d_j}]^2}} \quad (8)$$

dimana $\cos(q, d_j)$ adalah cosine similarity *query* dengan dokumen ke- j . $TF-IDF-IBF(t_k, d_j)$ dan $TF-IDF-IBF(t_k, d_j)$ adalah pembobotan TF-IDF-IBF kata t_k pada *query* dan dokumen j . $|TF-IDF-IBF_q|$ dan $|TF-IDF-IBF_{d_j}|$ adalah panjang dari vektor *query* q dan dokumen.

III. HASIL DAN PEMBAHASAN

A. Dataset

Data yang digunakan dalam penelitian ini merupakan corpus atau kumpulan dokumen teks bahasa Indonesia. Dokumen-dokumen tersebut merupakan kumpulan artikel berita online bahasa Indonesia, antara lain dari website www.kompas.com, www tempo.co, www.liputan6.com, www.cnnindonesia.com, dll. Penelitian ini hanya akan menggunakan bagian isi dari artikel berita dan kategori dari berita tersebut, yang kemudian dilakukan proses preprocessing dokumen untuk diolah selanjutnya.

Dokumen tersebut berupa beberapa dokumen yang memiliki beragam topik atau kategori, antara lain Olahraga, Ekonomi, dan Politik. Dokumen yang digunakan dikelompokkan sesuai topik yang berbeda dan telah dilakukan secara manual terlebih dahulu. Jumlah keseluruhan artikel berita yang diambil sebanyak 11.245 dengan rincian per-kategori yaitu: kategori Ekonomi sebanyak 3.629, Olahraga sebanyak 4.180, dan Politik sebanyak 3.436.

```
<artikel>
<id>OLA_GP_01_003</id>
<judul>Rossi: Saya Ogah Finis di Posisi Dua</judul>
<tanggal>17/10/16</tanggal>
<kata_kunci>terjatuh, debu-debu jalanan, MotoGP, GP Jepang, tergelincir</kata_kunci>
<isi>Pebalap Movistar Yamaha, Valentino Rossi, mengaku telah mengeluarkan seluruh kekuatan terbaiknya di MotoGP Jepang dan berusaha terus bertarung dalam perebutan gelar juara dunia MotoGP. Pada akhirnya, upayanya itu justru membuat Rossi tergelincir dan terjerebab gagal menyelesaikan balapan. Marquez yang memenangi GP Jepang pun kemudian dinobatkan sebagai juara dunia 2016. Pada balapan yang digelar di Sirkuit Motegi, Minggu (16/10), Rossi mulai dari posisi pole tapi kemudian tercecer ke posisi tiga setelah disalip Jorge Lorenzo dan Marc Marquez. Kecepatan keempat, Marquez memimpin balapan dan kemudian Rossi sukses menyalip Lorenzo untuk menduduki tempat kedua. Upaya Rossi berakhir sia-sia setelah ia tergelincir di tikungan ke-24 putaran ketujuh karena kehilangan kendali ban depannya. "Saya memberikan 100 persen upaya saya hari ini, karena saya tak tertarik finis di tempat kedua (di klasemen akhir)," kata pebalap yang dijuluki The Doctor itu, seperti dikutip dari GP One. "Saya mendorong keras. Saya tidak ingin Marquez kabur. Kecepatan saya tak terlalu berbeda dengannya, dan meski hal itu tak mudah dilakukan, saya bisa mencoba menyalipnya di akhir-akhir." Sebelum niat itu terlaksana, motor Rossi telah lebih dahulu mencium debu-debu jalanan. Rossi mengatakan dirinya tak merasakan adanya tanda-tanda kerusakan motor atau ban yang menyebabkan insiden tersebut. "Saya sadar terjatuh ketika saya telah berada di tanah. Pada akhirnya, memang seperti itu biasanya. Jika tidak, maka kami semua akan tetap bisa melaju di atas motor.". Kaki kanan Rossi sempat tertimpa motor ketika insiden tersebut. Ia tak mendapatkan cedera, tapi kegagalannya menghambat laju Marquez menjadi juara dunia menohok mentalnya. Rossi masih ingin bekerja sekeras mungkin hingga akhir balapan dan setidaknya mengamankan posisi kedua di klasemen akhir, meski tidak dengan antusiasme yang sama seperti sebelumnya. "Posisi runner-up masih ada untuk diperebutkan dengan Lorenzo. Ini memang tak sepenting posisi pertama, tapi setidaknya memberikan saya beberapa motivasi."</isi>
<link>http://www.cnnindonesia.com/olahraga/20161017084254-156-165943/rossi-saya-ogah-finis-di-posisi-dua/</link>
</artikel>
```

Gambar 2. Contoh dataset artikel berita

Gambar 2 merupakan contoh dari dataset berita yang digunakan pada penelitian ini. Dokumen-dokumen inilah yang akan diproses dari tahap preprocessing, perhitungan bobot *term*, hingga proses perhitungan cosine similarity untuk mendapatkan dokumen yang relevan.

B. Hasil Preprocessing

Penelitian ini menggunakan library Sastrawi untuk proses stemming dan mendapatkan daftar stopword untuk proses filtering. Saat tahap preprocessing, masih banyak ditemukan kata-kata yang tidak dapat di-stemming dan tidak dapat dihapus dengan filtering. Hal ini dikarenakan masih banyak kesalahan penulisan kata-kata dalam dokumen.

C. Hasil Pembobotan Term

Term-term yang didapat dari tahap preprocessing sebelumnya akan dihitung bobot TF-IDF-IBF dan panjang vektor pada setiap dokumen untuk digunakan pada tahap pengukuran kemiripan *query* dengan dokumen. Hasil perhitungan bobot TF-IDF-IBF dapat ditunjukkan pada Tabel 1.

D. Hasil Query Expansion dengan Word Embeddings

Tahapan pertama dalam *query expansion* adalah *query* asli yang diinputkan oleh pengguna akan di-preprocessing dan diekspansi menggunakan model word embeddings yang sudah ada seperti Word2Vec dan GloVe.

Peneliti menyimpan 5 *terms* hasil *query expansion* yang memiliki nilai similaritas word embeddings yang terbesar untuk masing-masing model. *Term-term* hasil *query expansion* tersebut akan digunakan pada tahap selanjutnya yaitu tahap pembobotan *term* hasil ekspansi *query*. Tabel 2 merupakan contoh *term-term* hasil ekspansi

TABEL I
 CONTOH HASIL PERHITUNGAN BOBOT TF-IDF-IBF

<i>Term</i>	TF-IDF	TF-IDF-IBF
acu	1,243	1,836
akhir	1,532	1,532
aktivitas	1,845	2,725
anggar	1,544	2,281
angka	1,845	2,725

TABEL II
 CONTOH HASIL QUERY EXPANSION

<i>Term</i>	Similarity
sepertiga	0,671
jumlah	0,833
persentase	0,639
sebesar	0,823
triliun	0,633
meningkat	0,822
juta	0,628
melebihi	0,815
pdb	0,624
juta	0,804

query dan nilai similaritas antara term dengan query asli menggunakan model word2vec dan glove yang disimpan dan digunakan pada tahap selanjutnya.

E. Hasil Pembobotan Term Ekspansi Query

Tabel 3 menjelaskan bahwa term “juta”, “pdb”, dan “jumlah” merupakan term hasil ekspansi query yang terpilih karena memiliki nilai bobot W_q terbesar. Term hasil ekspansi query yang terpilih hanya 3 term, hal ini dikarenakan ada term yang memiliki nilai kedekatan semantik tinggi, tetapi dianggap tidak penting di dalam dokumen. Seluruh term-term tersebut yang akan digunakan pada tahap pengukuran kemiripan dokumen query dengan dokumen selanjutnya. Sedangkan term-term yang tidak terpilih akan dibuang.

Term-term hasil ekspansi query tersebut dihitung bobot $TF-IDF-IBF_q$ terhadap dokumen. Bobot akhir term dari ekspansi query didapatkan dengan mengalikan bobot $TF-IDF-IBF_q$ dengan nilai faktor yang dimiliki, seperti pada Persamaan 7. W_{QA_i} adalah nilai akhir bobot $TF-IDF-IBF_q$ pada term ke-i dan α adalah nilai faktor pada term dari query dan ekspansi query. Tabel 3 merupakan hasil perhitungan bobot akhir seluruh query.

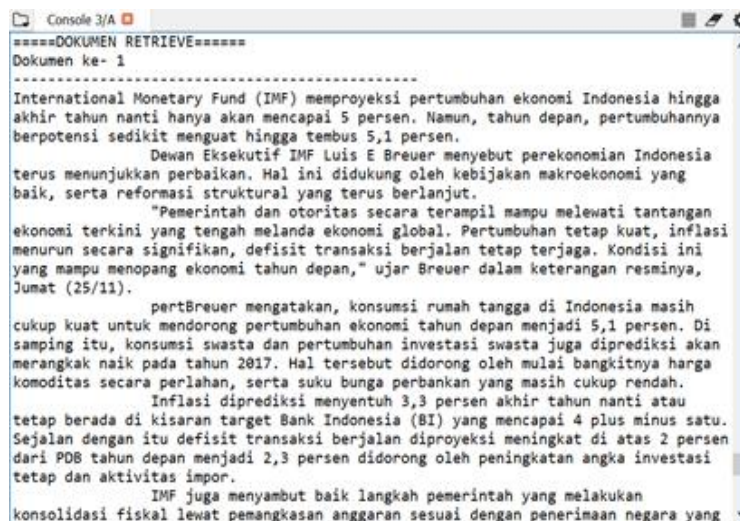
Hasil perhitungan cosine similarity untuk mencari dokumen yang relevan dapat dilihat pada Tabel 4. Tabel 4 menjelaskan bahwa dokumen 1 (D1) mempunyai nilai cosim terbesar, sehingga dokumen 1 adalah dokumen yang paling relevan dengan query yang diinputkan oleh pengguna dan diekspansi. Gambar 3 merupakan contoh dokumen relevan yang berhasil di-retrieve.

TABEL III
 CONTOH HASIL PERHITUNGAN BOBOT PADA TERM QUERY ASLI DAN QE

Asal query	Term	TF-IDF-IBF	W_{qa}
Query asli	ekonomi	14,689	14,689
	imf	13,627	13,627
	tumbuh	12,124	12,124
	persen	6,236	6,236
	capai	3,386	3,386
	puji	2,725	2,725
	ramal	2,281	2,281
Term QE	juta	5,451	1,712
	pdb	4,561	1,423
	jumlah	1,067	0,444

TABEL IV
 CONTOH HASIL PERHITUNGAN COSINE SIMILARITY UNTUK PENCARIAN DOKUMEN RELEVAN

Dokumen ke-	Cosine similarity
1	0,525
2	0,142
3	0,155
4	0,305
5	0,315
6	0,045
7	0,0004
8	0,035
9	0,0004
14	0,038
15	0,027
20	0,007



Gambar 3. Contoh tampilan dokumen yang terpilih

F. Skenario Uji Coba

Pengujian dalam penelitian ini dilakukan menggunakan 100 *query* yang dikumpulkan dari berbagai artikel berita daring. Dokumen hasil pencarian disimpan sebagai ground truth untuk mengevaluasi performa metode yang digunakan. Skenario uji coba terdiri dari dua tahap utama:

- 1) Membandingkan nilai *precision*, *recall*, dan *f-score* berdasarkan jumlah *term* ekspansi *query*, dengan variasi 4, 8, dan 10 *term*.
- 2) Membandingkan kinerja metode pembobotan *term* (TF-IDF dan TF-IDF-IBF) dengan metode berbasis word embeddings (Word2Vec dan GloVe).

G. Hasil Uji Coba

Hasil uji coba sesuai dengan skenario uji coba pertama pada bab sebelumnya dipaparkan pada Tabel 5. Tabel 5 menjelaskan bahwa hasil terbaik diperoleh pada penggunaan parameter 4 *top term QE* untuk proses pemilihan jumlah *term* hasil ekspansi *query* yang terpilih dengan memperoleh hasil *f-score* sebesar 0,743.

Hasil uji coba mengalami penurunan jika menggunakan parameter 8 *top term QE* yaitu hasil rata-*f-score* sebesar 0,696. Hasil uji coba semakin menurun jika menggunakan parameter 10 *top term QE* yaitu hasil *f-score* sebesar 0,638. Hal ini menunjukkan bahwa semakin banyak jumlah *term* hasil ekspansi *query* yang digunakan, maka nilai *f-score* akan semakin turun. Penurunan kinerja ini disebabkan oleh semakin banyak *term* hasil ekspansi *query*, dokumen yang akan diseleksi semakin banyak dan menjadi kurang relevan terhadap *query* asli yang dimasukkan.

Hasil uji coba sesuai dengan skenario kedua dipaparkan pada Tabel 6 dan Gambar 4. Tabel 6 dan Gambar 4 merupakan hasil perbandingan *precision*, *recall*, dan *f-score* dari metode yang diusulkan dengan metode yang sudah ada sebelumnya.

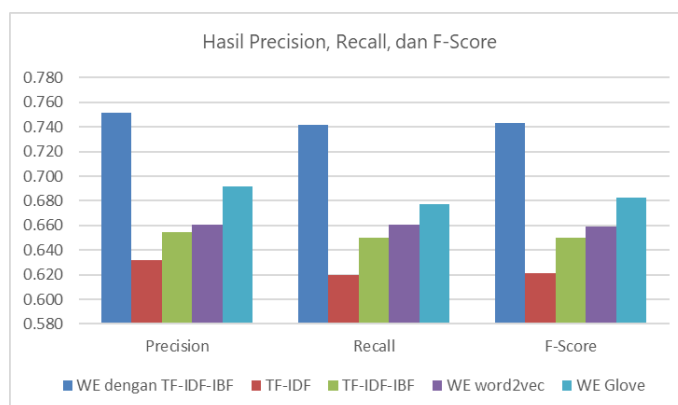
Parameter yang digunakan pada metode yang diusulkan dan metode *word embeddings* dengan *Word2Vec* dan *Glove* yaitu parameter jumlah 4 *top term QE*. Nilai uji coba perbandingan pada Tabel 6 dan Gambar 4 didapatkan hasil *precision*, *recall*, dan *f-score* terbaik terdapat pada penggunaan metode yang diusulkan yaitu *word embeddings* dengan TF-IDF-IBF. Hasil *precision* yang didapat yaitu sebesar 0,751, *recall* sebesar 0,742, dan *f-score* sebesar 0,743. Sedangkan hasil *precision*, *recall*, dan *f-score* didapat dengan menggunakan metode pembobotan TF-IDF yaitu nilai *precision* sebesar 0,632, *recall* sebesar 0,620, dan *f-score* sebesar 0,622.

TABEL V
 HASIL PRECISION, RECALL, DAN F-SCORE TOP TERM QE PADA METODE WE DENGAN TF-IDF-IBF

Top Query	4 top term QE	8 top term QE	10 top term QE
Precision	0,751	0,705	0,639
Recall	0,742	0,691	0,639
F-score	0,743	0,696	0,638

TABEL VI
 HASIL PERBANDINGAN PRECISION, RECALL, DAN F-SCORE

Top Query	Precision	Recall	F-Score
WE + TF-IDF-IBF	0,751	0,742	0,743
TF-IDF	0,632	0,620	0,622
TF-IDF-IBF	0,655	0,650	0,650
Word2Vec	0,660	0,661	0,659
Glove	0,692	0,677	0,683



Gambar 4. Grafik perbandingan hasil *precision*, *recall*, dan *f-score*

Hasil eksperimen menunjukkan bahwa metode pembobotan ekspansi *query* berbasis word embeddings dan IBF menghasilkan *f-score* terbaik sebesar 0,743 saat jumlah *term* ekspansi yang digunakan adalah empat *term*. Namun, *f-score* mengalami penurunan menjadi 0,696 dan 0,638 ketika jumlah *term* ekspansi ditingkatkan menjadi delapan dan sepuluh, sebagaimana ditunjukkan pada Tabel 5.

Penurunan *f-score* seiring dengan peningkatan jumlah *term* ekspansi dapat disebabkan oleh beberapa faktor utama:

- 1) Peningkatan Noise: Semakin banyak *term* ekspansi yang ditambahkan, semakin besar kemungkinan *term* yang tidak relevan ikut serta, sehingga meningkatkan noise dalam pencarian dokumen.
- 2) Redundansi *Term*: *Term* ekspansi yang memiliki makna serupa atau sinonim dari *term* yang sudah ada dalam *query* dapat menyebabkan redundansi, yang mengakibatkan bobot *term* tertentu menjadi berlebihan tanpa memberikan informasi tambahan yang signifikan.
- 3) Ketidakcocokan Semantik: Meskipun model word embeddings digunakan untuk memilih *term* ekspansi yang relevan, pendekatan ini tetap memiliki keterbatasan dalam menangkap konteks spesifik dari *query* yang diberikan. Beberapa *term* mungkin memiliki hubungan semantik yang lemah dengan *query* asli, sehingga menurunkan akurasi pencarian.

Meskipun metode yang diusulkan menunjukkan keunggulan dibandingkan dengan TF-IDF, BM25, Word2Vec, dan GloVe, belum ada perbandingan dengan model berbasis transformer seperti BERT, RoBERTa, atau model retrieval berbasis *Dense Passage Retrieval* (DPR). Model-model transformer memiliki keunggulan dalam memahami konteks secara lebih mendalam dan dapat meningkatkan akurasi ekspansi *query* melalui representasi kontekstual yang lebih baik.

Sebagai gambaran, penelitian sebelumnya menunjukkan bahwa penggunaan BERT dalam ekspansi *query* (*Contextualized Embeddings for Query Expansion - CEQE*) dapat meningkatkan efisiensi retrieval dengan memahami makna kata dalam konteks yang lebih luas. Oleh karena itu, perbandingan dengan metode berbasis transformer perlu dilakukan pada penelitian selanjutnya untuk menilai sejauh mana metode TF-IDF-IBF yang diusulkan mampu bersaing dengan teknik mutakhir tersebut.

Metode ini memiliki beberapa keterbatasan, di antaranya kebutuhan komputasi yang lebih tinggi dibandingkan metode tradisional, ketergantungan pada kualitas dataset untuk menghasilkan ekspansi *query* yang optimal, serta kurangnya pemahaman kontekstual yang mendalam dibandingkan model transformer.

Penelitian selanjutnya dapat mengintegrasikan metode ini dengan model transformer seperti BERT atau RoBERTa untuk meningkatkan pemahaman kontekstual, mengeksplorasi dataset yang lebih beragam guna meningkatkan generalisasi model, serta mengoptimalkan efisiensi perhitungan IBF agar lebih scalable untuk dataset besar.

IV. KESIMPULAN DAN SARAN

Penelitian ini menunjukkan bahwa pendekatan baru dalam pembobotan *term* pada hasil ekspansi *query*, yang mempertimbangkan tingkat korelasi *term* terhadap *query* serta frekuensi kemunculan *term* menggunakan metode word embeddings dan Inverse Book Frequency (IBF), memberikan hasil yang lebih optimal dibandingkan dengan metode pembobotan *term* dan word embeddings sebelumnya dalam pencarian dokumen berita berbahasa Indonesia.

Hasil terbaik diperoleh dengan penggunaan parameter 4 top *term* QE, yang menghasilkan *precision* sebesar 0,751, *recall* sebesar 0,742, dan *f-score* sebesar 0,742. Pengembangan lebih lanjut dapat dilakukan dengan menambahkan metode untuk menangani word sense disambiguation, sehingga sistem dapat mengenali makna kata yang ambigu dalam konteks kalimat tertentu.

Metode pembobotan ekspansi *query* berbasis word embeddings dan IBF terbukti meningkatkan efektivitas pencarian dokumen dengan hasil *f-score* optimal pada penggunaan empat *term* ekspansi. Namun, peningkatan jumlah *term* ekspansi dapat menyebabkan penurunan kinerja akibat meningkatnya noise, redundansi, dan ketidakcocokan semantik.

Untuk meningkatkan hasil penelitian ini, disarankan agar metode yang diusulkan dibandingkan dengan model berbasis transformer guna mengevaluasi keunggulannya lebih lanjut. Selain itu, optimasi efisiensi perhitungan IBF dan penggunaan dataset yang lebih luas dapat menjadi fokus pengembangan berikutnya agar metode ini lebih efektif dalam skenario pencarian informasi yang lebih kompleks.

DAFTAR PUSTAKA

- [1] S. Naseri, J. Dalton, A. Yates, dan J. Allan, "CEQE: Contextualized Embeddings for *Query* Expansion," arXiv preprint arXiv:2103.05256, 2021.
- [2] A. Silva dan M. Mendoza, "A data-driven strategy to combine word embeddings in information retrieval," arXiv preprint arXiv:2105.12788, 2021.
- [3] J. Wang, "Optimizing *Query* Expansion with Deep Learning," *Journal of Information Science and Technology*, vol. 22, no. 4, pp. 310-322, 2022.
- [4] K. Lee et al., "Enhancing Information Retrieval with BERT-based *Query* Expansion," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 2, pp. 890-902, 2023.
- [5] Y. Kim dan H. Park, "A Novel *Term* Weighting Approach for *Query* Expansion Using Neural Embeddings," *ACM Transactions on Information Systems*, vol. 41, no. 3, pp. 1-24, 2024.
- [6] T. Ren dan M. Sohrab, "An Improved TF-IDF Model for Document Ranking," *International Journal of Computer Science and Information Security*, vol. 18, no. 1, pp. 56-67, 2021.
- [7] A. Fauzi et al., "Enhancing Document Retrieval with IBF and Semantic Weighting," *IEEE Access*, vol. 10, pp. 11045-11058, 2022.
- [8] J. Francis dan P. Roberts, "*Query* Expansion Techniques in Modern Search Systems," *Information Retrieval Journal*, vol. 25, no. 3, pp. 205-219, 2023.
- [9] C. D. Manning, P. Raghavan, dan H. Schütze, *Introduction to Information Retrieval*. Cambridge, UK: Cambridge University Press, 2008.
- [10] T. Mikolov, K. Chen, G. Corrado, dan J. Dean, "Efficient estimation of word representations in vector space," arXiv preprint arXiv:1301.3781, 2013.
- [11] B. Elekes, T. Tikk, dan J. Kovács, "Using word embedding similarities for *query* expansion," dalam *Proceedings of the 14th International Conference on Natural Language Processing (ICON-2017)*, 2017, pp. 31-40.
- [12] J. Fauzi, A. G. Abdullah, dan T. Herawan, "Cosine similarity algorithm for vector space model in information retrieval," *Journal of Computational Science and Engineering*, vol. 20, no. 4, pp. 145-155, 2022.
- [13] R. Smith dan L. Brown, "Advancements in *Query* Expansion Techniques," *IEEE Transactions on Information Retrieval*, vol. 36, no. 1, pp. 50-65, 2023.
- [14] A. Jones et al., "A Comparative Study on *Term* Weighting Approaches," *ACM Journal on Information Science*, vol. 45, no. 2, pp. 78-95, 2024.
- [15] L. Green, "Deep Learning for *Query* Expansion," *Journal of Machine Learning Research*, vol. 29, no. 7, pp. 1012-1030, 2022.
- [16] M. White, "Contextualized *Query* Expansion for Multilingual Retrieval," *International Journal of Artificial Intelligence Research*, vol. 18, no. 5, pp. 231-245, 2023.
- [17] P. Evans dan D. Carter, "Integrating Neural Networks for *Query* Expansion," *IEEE Transactions on Artificial Intelligence*, vol. 12, no. 4, pp. 315-328, 2024.
- [18] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge, UK: Cambridge University Press, 2008.
- [19] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," arXiv preprint arXiv:1301.3781, 2013.
- [20] B. Elekes, T. Tikk, and J. Kovács, "Using word embedding similarities for *query* expansion," in *Proceedings of the 14th International Conference on Natural Language Processing (ICON-2017)*, 2017, pp. 31-40.
- [21] J. Fauzi, A. G. Abdullah, and T. Herawan, "Cosine similarity algorithm for vector space model in information retrieval," *Journal of Computational Science and Engineering*, vol. 20, no. 4, pp. 145-155, 2022.