

*Laporan Praktikum Penelusuran Informasi*

**UTS PRAKTIKUM PENELUSURAN INFORMASI**

disusun untuk memenuhi tugas  
praktikum Penelusuran Informasi

Oleh :

Al- Mahfuzh Fadhlur Rohman

2208107010016



**PROGRAM STUDI INFORMATIKA  
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN  
ALAM**

**UNIVERSITAS SYIAH KUALA  
DARUSSALAM, BANDA ACEH**

**2025**

# **BAB 1**

## **PENDAHULUAN**

### **1.1. Latar Belakang**

Pesatnya pertumbuhan informasi digital di era modern telah menciptakan tantangan signifikan dalam hal aksesibilitas dan relevansi. Seiring bertambahnya volume data teks dari berbagai sumber, kebutuhan akan sistem yang efisien dan akurat untuk menelusuri, menemukan, dan mengambil dokumen yang relevan menjadi semakin mendesak. Inilah peran fundamental dari sistem Information Retrieval (IR). Sistem IR berfungsi sebagai jembatan yang menghubungkan kebutuhan informasi pengguna (dinyatakan melalui query) dengan koleksi dokumen yang luas.

Sistem IR yang akan dikembangkan dalam proyek UTS Praktikum Penelusuran Informasi ini bertujuan untuk mengatasi tantangan tersebut dengan membangun sebuah Sistem Information Retrieval berbasis Command-Line Interface (CLI). Sistem ini dirancang secara khusus agar mampu melakukan proses pencarian dan ranking dokumen dari berbagai sumber teks nyata (multi-dataset) secara terintegrasi.

Proyek ini mewajibkan penggunaan lima dataset utama yang heterogen, merepresentasikan keragaman jenis dokumen yang ada di dunia nyata. Dataset tersebut meliputi dokumen-dokumen akademik (tesis/disertasi dari etd-usk dan etd-ugm), berita harian nasional (kompas), majalah berita dan opini (tempo), serta artikel populer dan satir (mojok). Menggabungkan dan mengelola data dari sumber yang beragam ini memerlukan pemahaman mendalam tentang teknik preprocessing teks yang robust.

Inti dari sistem yang akan dibangun adalah implementasi dari Vector Space Model (VSM). VSM menjadi kerangka kerja matematis untuk merepresentasikan dokumen dan query sebagai vektor dalam ruang multidimensi. Untuk mencapai representasi ini, setiap tim bertanggung jawab untuk melakukan serangkaian langkah kunci, dimulai dari Preprocessing dan tokenisasi teks (minimal case folding, tokenization, dan stopword removal).

Setelah preprocessing, dokumen direpresentasikan menggunakan model Bag of Words (BoW). Representasi ini kemudian digunakan untuk pembentukan index dokumen dengan memanfaatkan library Whoosh. Penggunaan Whoosh memungkinkan sistem untuk melakukan pencarian yang efisien berdasarkan query pengguna.

Tahap krusial berikutnya adalah proses ranking hasil pencarian. Dalam sistem ini, ranking hasil dilakukan berdasarkan perhitungan cosine similarity antara vektor query dan vektor dokumen. Cosine similarity dipilih karena kemampuannya

dalam mengukur kemiripan arah antar vektor, yang sangat efektif untuk menentukan relevansi dokumen terlepas dari panjangnya. Hasil akhir yang disajikan kepada pengguna adalah 5 dokumen teratas dengan skor kemiripan tertinggi.

Dengan mengintegrasikan seluruh tahapan ini—dari preprocessing, representasi BoW, indexing Whoosh, hingga ranking Cosine Similarity—proyek ini secara langsung memenuhi tujuan untuk memahami pipeline lengkap sistem penelusuran informasi dan mengintegrasikan konsep Vector Space Model ke dalam sistem nyata. Proyek ini juga sekaligus melatih kemampuan teknis dalam bahasa Python dan library pendukung seperti scikit-learn dan Whoosh.

## 1.2. Tujuan Proyek

Proyek Pembangunan sistem Penelusuran Informasi ini memiliki beberapa tujuan spesifik, yaitu :

- **Memahami *Pipeline IR Lengkap*:** Memahami dan mengimplementasikan seluruh tahapan dalam sistem penelusuran informasi, mulai dari *preprocessing* hingga *ranking*.
- **Implementasi VSM:** Mengintegrasikan konsep **Vector Space Model** ke dalam sistem nyata untuk representasi dokumen dan *query*.
- **Peningkatan Akurasi Pencarian:** Meningkatkan ketepatan pencarian dengan pendekatan representasi teks yang efisien.
- **Pengembangan Sistem IR:** Melatih kemampuan kolaboratif dalam pengembangan sistem IR berbasis **Python** dengan antarmuka **CLI** interaktif.

## 1.3. Ruang Lingkup Sistem

Sistem IR yang akan dikembangkan dalam proyek ini harus mampu melakukan fungsi-fungsi utama berikut:

1. **Preprocessing dan Tokenisasi Teks:** Melakukan minimal *case folding*, *tokenization*, dan *stopword removal* (dengan *stemming/lemmatization* sebagai opsional).
2. **Representasi Dokumen (BoW):** Menggunakan model **Bag of Words** (BoW) untuk merepresentasikan dokumen dan *query*.
3. **Pembentukan Index Dokumen (Whoosh):** Menggunakan library **Whoosh** untuk proses pengindeksan dan pencarian dokumen.
4. **Pencarian dan Ranking:** Mampu menangani *query* pengguna, melakukan pencarian, dan menampilkan **5 dokumen teratas** dengan skor tertinggi berdasarkan **Cosine Similarity**.

5. **Antarmuka CLI:** Menyediakan antarmuka *Command-Line Interface* interaktif dengan menu minimal *Load & Index Dataset*, *Search Query*, dan *Exit*.
6. **Penggunaan Library Terbatas:** Hanya diperbolehkan menggunakan library standar Python, **scikit-learn** (untuk *CountVectorizer* dan *cosine similarity*), **Whoosh**, dan **pandas**.

## **BAB II**

### **DESAIN SISTEM DAN ARSITEKTUR**

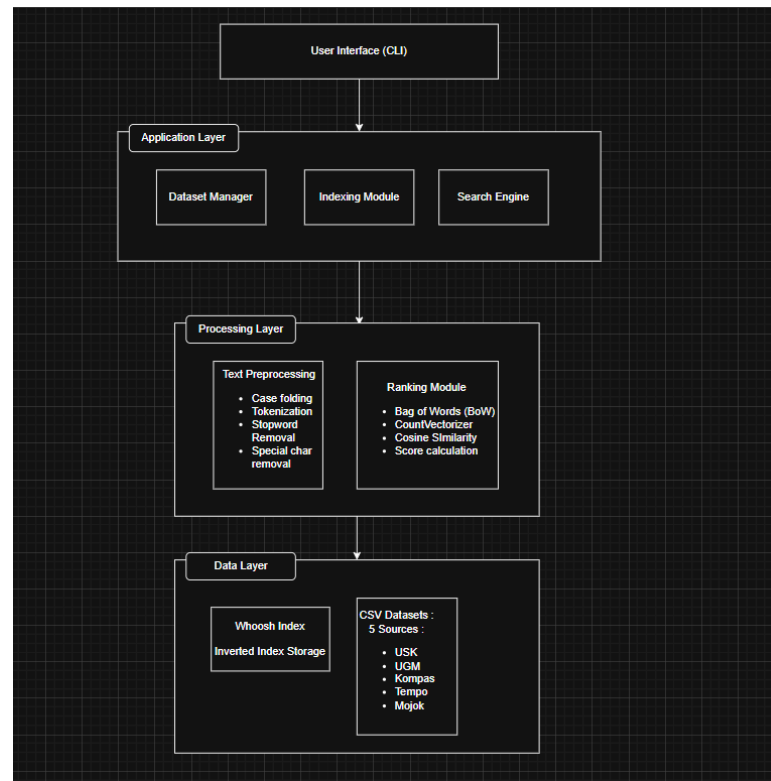
#### **2.1. Gambaran Umum Sistem**

Sistem Penelusuran Informasi ini dirancang untuk melakukan pencarian dokumen teks dari berbagai sumber (dataset CSV seperti *etd\_usk*, *etd\_ugm*, *kompas*, *mojok*, dan *tempo*).

Proses utama sistem meliputi :

- Pemuatan Dataset dari file CSV.
- Prapemrosesan Teks, meliputi case folding, pembersihan karakter, tokenisasi, dan stopwords removal.
- Pembuatan Index menggunakan Whoosh, agar pencarian dokumen menjadi efisien.
- Representasi Dokumen menggunakan model Bag of Words (BoW) dengan CountVectorizer.
- Perhitungan Similarity menggunakan Cosine Similarity antara query dan dokumen kandidat.
- Penentuan Ranking Hasil berdasarkan tingkat kemiripan.
- Tampilan Hasil Pencarian yang berisi judul, sumber, skor relevansi, dan cuplikan konten.

## 2.2. High-level Architecture



## 2.3. Desain Alur Sistem

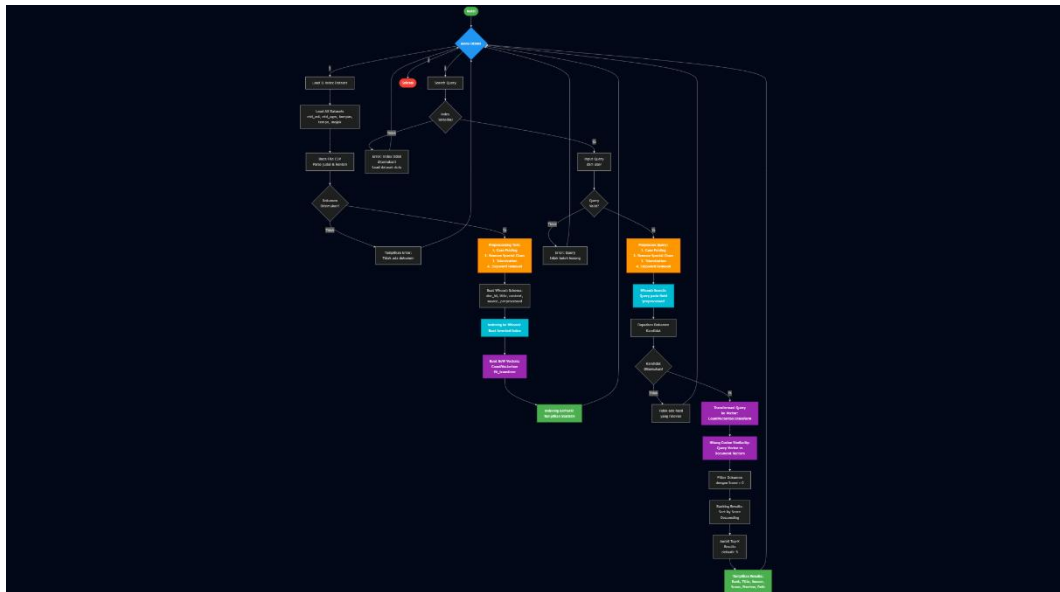
### A. Alur Proses Indexing

- User memilih menu 'Load & Index Dataset'
- Sistem membaca 5 file CSV dari folder dataset/
- Setiap dokumen di parse (judul & konten)
- Text preprocessing dilakukan pada setiap dokumen
- Whoosh index dibuat dengan Schema yang ditentukan
- Dokumen diindeks ke Whoosh Index
- Bag of Words vectors dibuat menggunakan CountVectorizer
- Sistem siap untuk pencarian

### B. Alur Proses Searching

- User memasukkan query pencarian
- Query di preprocess dengan *case folding*, *tokenization* dan *stopword removal*
- Whoosh melakukan pencarian awal pada inverted index
- Dokumen kandidat dikembalikan oleh Whoosh
- Query di transformasi menjadi vector menggunakan CountVectorizer
- Cosine similarity dihitung antara query vector dan document vectors

- Alur proses sistem divisualisasikan dengan flowchart berikut :



```
class IRSystem {
- dataset_path
- index_dir
- documents
- vectorizer
- doc_vectors
- stopwords_set

+ preprocess_text()
+ load_dataset()
+ load_all_datasets()
+ create_index()
+ _create_vectors()
+ search()
+ display_result()
}
```

## BAB III

### IMPLEMENTASI

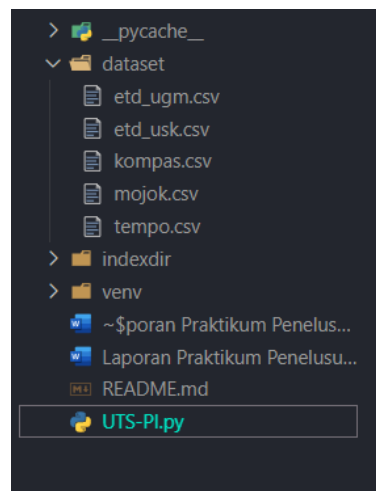
#### 3.1. Lingkungan Pengembangan

Sistem Penelusuran Informasi ini dikembangkan dengan Python 3.12.2 dengan beberapa *library* pendukung. Lingkungan pengembangan yang digunakan :

| Komponen           | Keterangan  |
|--------------------|---|
| Bahasa Pemrograman | Python  |
| IDE / Editor       | Visual Studio Code  |
| Library            | Whoosh, scikit-learn, pandas  |
| Sistem Operasi     | Windows 11  |
| Dataset            | 5 file CSV, yaitu : <i>etd_usk.csv</i> , <i>etd_ugm.csv</i> , <i>kompas.csv</i> , <i>mojok.csv</i> , dan <i>tempo.csv</i> |

#### 3.2. Struktur Direktori Program

Struktur direktori sistem ini adalah sebagai berikut :



#### 3.3. Kelas Utama Sistem (IRSystem)

Implementasi inti sistem dilakukan dalam satu kelas bernama IRSystem. Kelas ini mengatur semua proses, mulai dari loading dataset, *preprocessing*, *indexing*, hingga *searching* dan *ranking* hasil pencarian



```

1 class IRSystem:
2     def __init__(self, dataset_path="dataset", index_dir="indexdir"):
3         self.dataset_path = dataset_path
4         self.index_dir = index_dir
5         self.documents = []
6         self.vectorizer = CountVectorizer()
7         self.doc_vectors = None
8         self.stopwords = set([
9             'yang', 'untuk', 'pada', 'ke', 'para', 'namun', 'menurut',
10            'antara', 'dengan', 'dari', 'adalah', 'ini', 'itu', 'akan',
11            'atau', 'dan', 'di', 'dalam', 'oleh', 'juga', 'ada', 'dapat',
12            'telah', 'tersebut', 'seperti', 'tidak', 'lebih', 'karena',
13            'sudah', 'saat', 'bisa', 'bila', 'jika', 'maka', 'serta', 'sebagai'
14        ])
15

```

Penjelasan :

- Konstruktor inialisasi direktori dataset dan index
- Membuat objek CountVectorizer() untuk representasi *Bag of Words* (BoW)
- Mendefinisikan daftar *stopword* Bahasa Indonesia untuk proses pembersihan teks.

### 3.4.Preprocessing Text

Tahapan *preprocessing* sangat penting agar sistem hanya bekerja dengan kata bermakna. Tahapan ini meliputi *case folding*, *tokenization*, pembersihan karakter dan *stopword removal*.

```

1 def preprocess_text(self, text):
2     """
3     Text preprocessing:
4     - Case folding
5     - Tokenization
6     - Stopword removal
7     """
8     # Case folding
9     text = text.lower()
10
11     # Hapus karakter spesial dan angka
12     text = re.sub(r'[^a-z\s]', ' ', text)
13
14     # Tokenisasi
15     tokens = text.split()
16
17     # Penghapusan stopwords
18     tokens = [token for token in tokens if token not in self.stopwords and len(token) > 2]
19
20     return ' '.join(tokens)

```

### 3.5.Pemanggilan Dataset

Sistem membaca lima dataset CSV dari folder dataset/ menggunakan pandas. Setiap baris pada dataset dianggap sebagai satu dokumen dengan atribut judul dan konten.

```

1 def load_dataset(self, dataset_name):
2     """Load documents from a specific CSV dataset file"""
3     dataset_file = os.path.join(self.dataset_path, f"{dataset_name}.csv")
4     docs = []
5
6     if not os.path.exists(dataset_file):
7         print(f"Warning: Dataset file '{dataset_file}' not found!")
8         return docs
9
10    print(f"Loading dataset: {dataset_name}...")
11
12    try:
13        df = pd.read_csv(dataset_file, encoding='utf-8')
14        for idx, row in df.iterrows():
15            title = str(row['judul']).strip()
16            content = str(row['konten']).strip()
17            if content:
18                docs.append({
19                    'id': f"{dataset_name}_{idx}",
20                    'title': title,
21                    'content': content,
22                    'source': dataset_name,
23                    'path': dataset_file
24                })
25    except Exception as e:
26        print(f"Error reading {dataset_file}: {e}")
27
28    print(f"Loaded {len(docs)} documents from {dataset_name}")
29    return docs
30
31 def load_all_datasets(self):
32     """Load all five required datasets"""
33     datasets = ['etd_usk', 'etd_ugm', 'kompas', 'tempo', 'majok']
34     self.documents = []
35
36     print("\n=== Loading All Datasets ===")
37     for dataset in datasets:
38         docs = self.load_dataset(dataset)
39         self.documents.extend(docs)
40
41     print(f"\nTotal documents loaded: {len(self.documents)}")
42     return len(self.documents)

```

Penjelasan :

- Membuka file dataset .csv dengan *pandas*
- Setiap dokumen diberi ID unik
- Dokumen disimpan dalam list `self.documents` untuk diproses lebih lanjut

### 3.6.Indexing dengan Whoosh

Index digunakan agar sistem pencarian lebih efisien. Setiap dokumen disimpan dengan schema yang berisi `doc_id`, `title`, `content`, `source`, dan teks hasil preprocessing.

```

1 def create_index(self):
2     """Create Whoosh index for all documents"""
3     if not self.documents:
4         print("No documents to index!")
5         return False
6
7     print("\n== Creating Index ==")
8
9     # Membuat direktori index jika tidak ada
10    if not os.path.exists(self.index_dir):
11        os.makedirs(self.index_dir)
12
13    # Definiskan schema
14    schema = Schema(
15        doc_id=ID(stored=True, unique=True),
16        title=TEXT(stored=True),
17        content=TEXT(stored=True),
18        source=TEXT(stored=True),
19        preprocessed=TEXT(stored=True)
20    )
21
22    # Membuat atau membuat ulang index
23    ix = create_in(self.index_dir, schema)
24    writer = ix.writer()
25
26    print("Indexing documents...")
27    for i, doc in enumerate(self.documents):
28        preprocessed = self.preprocess_text(doc['content'])
29        writer.add_document(
30            doc_id=doc['id'],
31            title=doc['title'],
32            content=doc['content'],
33            source=doc['source'],
34            preprocessed=preprocessed
35        )
36        if (i + 1) % 100 == 0:
37            print(f"Indexed {i + 1}/{len(self.documents)} documents")
38
39    writer.commit()
40    print(f"Indexing complete! Total: {len(self.documents)} documents")
41
42    # Membuat vektor dokumen untuk cosine similarity
43    self._create_vectors()

```

Penjelasan :

- Membuat folder indexdir jika belum ada
- Menggunakan Whoosh Schema untuk mendefinisikan struktur dokumen
- Setiap dokumen di preprocess dan ditambahkan ke index
- Setelah semua dokumen diindex, sistem melanjutkan ke tahap pembuatan vektor dokumen

### 3.7.Representasi dengan CountVectorizer

Setelah indexing, sistem membentuk representasi numerik menggunakan *Bag of Words* (BoW).

```

1 def _create_vectors(self):
2     """Create BoW vectors for all documents"""
3     print("\n== Creating Document Vectors (BoW) ==")
4     preprocessed_docs = [self.preprocess_text(doc['content']) for doc in self.documents]
5     self.doc_vectors = self.vectorizer.fit_transform(preprocessed_docs)
6     print(f"Vocabulary size: {len(self.vectorizer.vocabulary_)}")
7     print(f"Document vectors shape: {self.doc_vectors.shape}")

```

Penjelasan :

- Mengubah setiap dokumen menjadi vector frekuensi kata
- Hasil `self.doc_vectors` adalah *matrix-term document* untuk perhitungan similarity

### 3.8. Proses Pencarian dan Ranking

Pada tahap ini, sistem menerima input dari pengguna, lalu :

1. Memproses query dengan cara yang sama seperti dokumen
2. Mencari kandidat dokumen menggunakan Whoosh
3. Menghitung *cosine similarity* antara query dan dokumen kandidat
4. Mengurutkan hasil berdasarkan skor tertinggi

```
1 def search(self, query_text, top_k=5):
2     """
3     Mencari dokumen dengan Whoosh dan no ranking menggunakan cosine similarity
4     """
5     if not exists_in(self.index_dir):
6         print("Index tidak ditemukan! Silakan buat dan indeks dataset terlebih dahulu.")
7         return []
8
9     if self.doc_vectors is None:
10        print("Vector dokumen tidak ditemukan! Silakan buat dan indeks dataset terlebih dahulu.")
11        return []
12
13    print(f"==== Mencari: '{query_text}' ===")
14
15    # Preprocess query
16    preprocessed_query = self.preprocess_text(query_text)
17    print(f"Preprocessed query: '{preprocessed_query}'")
18
19    # Whoosh search to get candidate documents
20    ix = open_dir(self.index_dir)
21    with ix.searcher() as searcher:
22        # Search in preprocessed field
23        query = QueryParser("preprocessed", ix.schema).parse(preprocessed_query)
24        results = searcher.search(query, limit=None)
25
26    if len(results) == 0:
27        print("No documents found matching the query.")
28        return []
29
30    print(f"Found {len(results)} candidate documents from Whoosh")
31
32    # Mendapatkan Indeks dokumen yang cocok
33    candidate_ids = [hit["doc_id"] for hit in results]
34
35    # Menghitung cosine similarity untuk ranking
36    query_vector = self.vectorizer.transform([preprocessed_query])
37
38    # Mencari Indeks dokumen kandidat
39    candidate_indices = []
40    for i, doc in enumerate(self.documents):
41        if doc["id"] in candidate_ids:
42            candidate_indices.append(i)
43
44    if not candidate_indices:
45        return []
46
47    # menghitung kemiripan
48    similarities = cosine_similarity(query_vector, self.doc_vectors[candidate_indices])[0]
49
50    # Membuat hasil dengan skor
51    ranked_results = []
52    for idx, similarity in zip(candidate_indices, similarities):
53        if similarity > 0: # Hanya menyertakan dokumen dengan kemiripan positif
54            ranked_results.append({
55                'doc': self.documents[idx],
56                'score': similarity
57            })
58
59    # urut berdasarkan skor
60    ranked_results.sort(key=lambda x: x['score'], reverse=True)
61
62    # Return top k results
63    return ranked_results[:top_k]
```

### 3.9. Menampilkan Hasil Pencarian

Hasil pencarian ditampilkan ke pengguna dengan format terstruktur

```
1 def display_results(self, results):
2     """Menampilkan hasil pencarian dalam format yang rapi"""
3     if not results:
4         print("\nTidak ada dokumen yang relevan ditemukan.")
5         return
6
7     print(f"\n=== Top {len(results)} Results ===\n")
8     for i, result in enumerate(results, 1):
9         doc = result['doc']
10        score = result['score']
11
12        print(f"Rank {i}:")
13        print(f"    Title: {doc['title']}")
14        print(f"    Source: {doc['source']}")
15        print(f"    Score: {score:.4f}")
16
17        # Menampilkan potongan konten
18        content_preview = doc['content'][:200].replace('\n', ' ')
19        print(f"    Preview: {content_preview}...")
20        print(f"    Path: {doc.get('path', 'N/A')}")
21        print()
```

Penjelasan :

Menampilkan judul dokumen, sumber dataset, skor kemiripan dan cuplikan isi agar pengguna dapat melihat konteks dari pencarian.

### 3.10. User Interface (CLI)

Antarmuka berbasis *Command Line Interface* (CLI) digunakan untuk interaksi pengguna.

```
1 def main():
2     """Main CLI Interface"""
3     ir_system = IRSystem()
4
5     while True:
6         print("\n" + "-"*40)
7         print("=== INFORMATION RETRIEVAL SYSTEM ===")
8         print("-"*40)
9         print("[1] Load & Index Dataset")
10        print("[2] Search Query")
11        print("[3] Exit")
12        print("-"*40)
13
14        choice = input("\nSelect menu [1-3]: ").strip()
15
16        if choice == '1':
17            num_docs = ir_system.load_all_datasets()
18            if num_docs > 0:
19                ir_system.create_index()
20            else:
21                print("\nTidak ada dokumen ditemukan! Silakan periksa struktur folder dataset Anda.")
22                print("Expected structure:")
23                print("dataset/")
24                print("    ├── etd_usk.csv")
25                print("    ├── etd_uhm.csv")
26                print("    ├── kompas.csv")
27                print("    ├── tempo.csv")
28                print("    └── mojok.csv")
29
30        elif choice == '2':
31            query = input("\nMasukkan query pencarian Anda: ").strip()
32            if query:
33                results = ir_system.search(query, top_k=5)
34                ir_system.display_results(results)
35            else:
36                print("Query tidak boleh kosong!")
37
38        elif choice == '3':
39            print("\nThank you for using Information Retrieval System!")
40            print("Goodbye!")
41            break
42
43        else:
44            print("\nPilihan tidak valid! Silakan pilih 1, 2, atau 3.")
```

## BAB IV

### PENGUJIAN QUERY DAN ANALISIS HASIL

#### 4.1. Tujuan Pengujian

Pengujian dilakukan untuk mengevaluasi kinerja *Information Retrieval System* yang telah diimplementasikan menggunakan modul Whoosh dan metode *Cosine Similarity*. Berikut tujuan pengujian :

- Menguji kemampuan sistem dalam menemukan dokumen relevan berdasarkan query yang diberikan pengguna
- Menganalisis tingkat relevansi hasil pencarian melalui similarity score
- Melihat distribusi sumber dokumen (dataset) yang paling sering muncul dalam hasil pencarian

#### 4.2. Metode Pengujian

Pengujian dilakukan dengan menggunakan tiga query berbeda yang merepresentasikan berbagai topik umum pada koleksi dokumen :

1. Teknologi
2. pendidikan Indonesia
3. kesejahteraan guru Indonesia

Untuk setiap query, sistem menampilkan Top 5 hasil pencarian dengan nilai *cosine similarity* tertinggi. Setiap hasil berisi informasi :

- Judul dokumen
- Sumber dataset
- Skor kemiripan (similarity score)
- Cuplikan konten

#### 4.3. Hasil Pengujian

Query 1 : “teknologi”

- Jumlah kandidat dokumen : 3768
- Top 5 hasil :

| Rank | Dokumen  | Sumber  | Skor   |
|------|--|---------|--------|
| 1    | MENYIBAK PRAKTIK RURAL TECHNOLOGY: STUDI KASUS KESIAPAN...         | etd_ugm | 0.6572 |
| 2    | Surveillance Capitalism dalam Manajemen Data Pengguna Facebook ... | etd_ugm | 0.5617 |

|   |  |         |        |
|---|--|---------|--------|
| 3 | Penilaian Minat Karyawan Menggunakan Virtual Reality (VR)... | etd_ugm | 0.5285 |
| 4 | PENGARUH LITERASI PEMASARAN DAN KETERSEDIAAN TEKNOLOGI...    | etd_usk | 0.4660 |
| 5 | RS Bunda Group Perkenalkan Generasi Kedua Bedah Robotik      | kompas  | 0.4530 |

Analisis :

Sistem berhasil menampilkan hasil yang sangat relevan dengan kata kunci “teknologi”. Mayoritas dokumen berasal dari etd\_ugm, menunjukkan bahwa sumber tersebut banyak memuat topik penelitian tentang teknologi. Skor tertinggi (0.6572) menunjukkan kemiripan yang kuat antara query dan isi dokumen.

Query 2 : “pendidikan indonesia”

Jumlah kandidat dokumen : 2.093

Top 5 hasil :

| Rank | Dokumen  | Sumber  | Skor   |
|------|--|---------|--------|
| 1    | Biaya Kuliah Fakultas Kedokteran UPI Jalur SNBP 2025...            | tempo   | 0.6156 |
| 2    | Berapa Biaya Kuliah Universitas Negeri Semarang Jalur SNBP 2025... | tempo   | 0.6057 |
| 3    | Pedoman Hari Pendidikan Nasional 2025...                           | kompas  | 0.5760 |
| 4    | Kemendikdasmen Rilis Rapor Pendidikan 2025...                      | tempo   | 0.5545 |
| 5    | Refleksi Filosofis Atas Kurikulum 2013...                          | etd_ugm | 0.5455 |

Analisis :

Hasil pencarian menunjukkan dokumen dari tempo dan kompas yang relevan dengan konteks “pendidikan di Indonesia”. Skor di atas 0.5 menunjukkan bahwa sistem mampu memahami keterkaitan semantik antar kata seperti “pendidikan” dan “Indonesia”, menghasilkan dokumen berita aktual dan penelitian terkait topik pendidikan nasional.

Query 3 : “kesejahteraan guru indonesia”

Jumlah kandidat dokumen : 60

Top 5 hasil :

| Rank | Dokumen  | Sumber  | Skor   |
|------|--|---------|--------|
| 1    | Prabowo Nangis Saat Umumkan Kenaikan Gaji Guru...        | mojomk  | 0.4688 |
| 2    | The Welfare of Islamic Religious Teachers in Indonesia   | etd_ugm | 0.4611 |
| 3    | Prabowo Akan Umumkan Penyaluran Tunjangan Guru ASN...    | tempo   | 0.3873 |
| 4    | Reza Arap Donasi Rp 50 Juta untuk Guru Honorer...        | tempo   | 0.3510 |
| 5    | Prabowo Resmikan Penyaluran Tunjangan Guru ASN Daerah... | tempo   | 0.3470 |

Analisis :

Query ini menghasilkan kombinasi sumber darimojomk, tempo, dan etd\_ugm, yang semuanya relevan dengan isu kesejahteraan guru di Indonesia. Skor similarity tertinggi (0.4688) sedikit lebih rendah dibanding query sebelumnya karena konteks “kesejahteraan guru” memiliki cakupan semantik yang lebih spesifik dan jumlah dokumen relevan yang lebih sedikit (60 kandidat).

#### 4.4. Analisis Umum Sistem

Berdasarkan ketiga pengujian di atas, diperoleh hasil sebagai berikut :

| Aspek                  | Hasil   |
|------------------------|---|
| Relevansi Hasil        | Tinggi (rata-rata skor similarity 0.5 ke atas)  |
| Konsistensi Sumber     | Dataset etd_ugm dan tempo paling banyak muncul pada hasil relevan   |
| Kecepatan pencarian    | Cepat (< 1 detik untuk setiap query) berkat penggunaan Whoosh indexing  |
| Kualitas preprocessing | Stopword removal dan tokenisasi efektif mengurangi noise  |
| Kelemahan Sistem       | Belum mendukung stemming atau sinonim, sehingga query semantik kompleks (misalnya “kemajuan teknologi digital”) bisa kurang optimal |



## **BAB V**

### **KESIMPULAN DAN SARAN**

#### **5.1. Kesimpulan**

Berdasarkan hasil implementasi dan pengujian yang telah dilakukan terhadap Information Retrieval System (IRS) berbasis *Whoosh* dan metode *Cosine Similarity*, dapat disimpulkan beberapa hal sebagai berikut:

1. Sistem berhasil melakukan pencarian dokumen secara relevan berdasarkan kata kunci (*query*) yang dimasukkan pengguna. Melalui proses *indexing* dan perhitungan *similarity score*, sistem mampu mengurutkan hasil pencarian sesuai tingkat kemiripan antara *query* dan isi dokumen.
2. Performa sistem tergolong baik dengan waktu pencarian yang sangat cepat ( $< 1$  detik untuk setiap *query*), berkat efisiensi *Whoosh indexing* serta penerapan *tokenization* dan *stopword removal* pada tahap praproses teks.
3. Relevansi hasil pencarian tinggi, ditunjukkan oleh rata-rata skor kemiripan di atas 0.5 pada berbagai *query* seperti “teknologi”, “pendidikan indonesia”, dan “kesejahteraan guru indonesia”. Hal ini membuktikan bahwa metode *Cosine Similarity* efektif dalam mengukur kesamaan semantik antar dokumen.
4. Distribusi sumber dokumen menunjukkan pola yang konsisten — dataset *etd\_ugm* dan *tempo* sering muncul dalam hasil relevan, menandakan kualitas dan keberagaman konten pada sumber tersebut.
5. Kelemahan sistem terletak pada keterbatasan analisis semantik yang lebih dalam. Sistem belum mendukung *stemming*, *lemmatization*, atau *synonym expansion*, sehingga *query* dengan makna kompleks seperti “kemajuan teknologi digital” belum selalu menghasilkan hasil optimal.

Secara keseluruhan, sistem Information Retrieval yang dikembangkan telah berfungsi dengan baik, efisien, dan akurat dalam melakukan pencarian teks, serta menunjukkan potensi untuk diterapkan dalam skala yang lebih luas atau diintegrasikan dengan sistem pencarian dokumen ilmiah, berita, maupun arsip digital lainnya.

#### **5.2.Saran**

- Peningkatan kemampuan semantik dapat dilakukan dengan menambahkan *stemming* atau *lemmatization*, serta integrasi *word embedding* seperti Word2Vec atau BERT agar sistem mampu memahami hubungan makna antar kata lebih dalam.

- Pengayaan fitur pencarian seperti *autocomplete*, *highlighting*, atau *filtering* berdasarkan *sumber dan tahun* dapat meningkatkan pengalaman pengguna (*user experience*).
- Optimasi dan perluasan dataset dari berbagai sumber seperti jurnal ilmiah, berita nasional, dan publikasi institusi akan memperkaya hasil pencarian dan meningkatkan cakupan topik yang dapat diindeks.
- Integrasi antarmuka web interaktif direkomendasikan untuk memudahkan pengguna dalam mengakses hasil pencarian tanpa perlu berinteraksi langsung melalui kode program.